



Tools for analyzing binary formats



Petr Pučil

admin and developer of **Kaitai Struct** (a tool for working with binary formats)



Contents

1. Binary formats

Text vs. binary files
Structure of binary formats

2. Binary analysis

Parsing and serialization
Getting a parser

3. Kaitai Struct

History
Features

4. Other tools

Kaitai collaboration with other tools

1

Binary formats



1 / Text vs. binary file

Text file

- the basic unit is a **character**
- divided into **lines** - sequence of printable and whitespace characters + end of line marker (LF, CR, CR LF)
- **Readable** and easy to edit in plain text editors (Notepad, VS Code, Notepad++, Vim)
- .txt, .rtf, .csv, .json, .xml
- Characters are encoded into bytes according to the **character set** (e.g. ASCII, UTF-8)

Binary file

- the basic unit is a **byte** (1 byte = 8 bits)
- a sequence of bytes from start to end
- **Unreadable** in a plain text editor
- Image (.jpg), video (.mp4), font (.ttf), database (.sqlite), archive (.zip), executable code (ELF, Windows .exe)
- Integer values are in the **same format** as in the computer's memory



1 / Structure of binary formats

1 x 1 pixel GIF

Parsing a binary file

```

0001 0203 0405 0607 0809 0A0B 0C0D 0E0F | 0123456789ABCDEF
00 4749 4638 3961 0100 0100 8000 00FF 8000 | GIF89a....€..`€.
10 FFFF FF2C 0000 0000 0100 0100 0002 0244 | ',.....D
20 0100 3B | ..;

```

Displaying the binary representation of a GIF file in the PSpad editor

	Offset (dec)	Byte sequence (hex)	Byte sequence meaning
GIF signature + version	00	47 49 46	ASCII string 'GIF' – "magic" file number (GIF signature)
	03	38 39 61	ASCII string '89a' – GIF version
Logical Screen Descriptor	06	01 00	Logical screen width: one pixel
	08	01 00	Logical screen height: one pixel
	10	80	Bit field: enable global color palette
	11	00	Background color index
Global Color Table	12	00	Pixel width to height ratio: 1:1
	13	FF 80 00	First color in the palette: orange (#ff8000)
	16	FF FF FF	Second color in the palette: white (#ffffff)

Image data

Offset (dec)	Byte sequence (hex)	Byte sequence meaning
19	2C	Image descriptor start marker
20	00 00	x-position of the left edge of the image: column zero
22	00 00	y-position of the top edge of the image: line zero
24	01 00	Image width: one pixel
26	01 00	Image height: one pixel
28	00	Bit field: no color palette and no interlacing
29	02	Initial LZW code size in bits minus 1
30	02	LZW-encoded block size in bytes
31	44 01	Encoded image data (one pixel with zero index)
33	00	Terminator of the LZW-encoded block
34	3B	GIF file terminator: the ';' character



1 / Structure of binary formats

1 x 1 pixel GIF

Parsing a binary file

```

0001 0203 0405 0607 0809 0A0B 0C0D 0E0F | 0123456789ABCDEF
00 4749 4638 3961 0100 0100 8000 00FF 8000 | GIF89a....€.. €.
10 FFFF FF2C 0000 0000 0100 0100 0002 0244 | .....D
20 0100 3B | ..;

```

Displaying the binary representation of a GIF file in the PSpad editor

	Offset (dec)	Byte sequence (hex)	Byte sequence meaning
GIF signature + version	00	47 49 46	ASCII string 'GIF' – "magic" file number (GIF signature)
	03	38 39 61	ASCII string '89a' – GIF version
Logical Screen Descriptor	06	01 00	Logical screen width: one pixel
	08	01 00	Logical screen height: one pixel
	10	80	Bit field: enable global color palette
	11	00	Background color index
Global Color Table	12	00	Pixel width to height ratio: 1:1
	13	FF 80 00	First color in the palette: orange (#ff8000)
	16	FF FF FF	Second color in the palette: white (#ffffff)

Image data

Offset (dec)	Byte sequence (hex)	Byte sequence meaning
19	2C	Image descriptor start marker
20	00 00	x-position of the left edge of the image: column zero
22	00 00	y-position of the top edge of the image: line zero
24	01 00	Image width: one pixel
26	01 00	Image height: one pixel
28	00	Bit field: no color palette and no interlacing
29	02	Initial LZW code size in bits minus 1
30	02	LZW-encoded block size in bytes
31	44 01	Encoded image data (one pixel with zero index)
33	00	Terminator of the LZW-encoded block
34	3B	GIF file terminator: the ';' character



1 / Structure of binary formats

1 x 1 pixel GIF

Parsing a binary file

```

0001 0203 0405 0607 0809 0A0B 0C0D 0E0F | 0123456789ABCDEF
00 4749 4638 3961 0100 0100 8000 00FF 8000 | GIF89a....€..`€.
10 FFFF FF2C 0000 0000 0100 0100 0002 0244 | ',.....D
20 0100 3B | ..;

```

Displaying the binary representation of a GIF file in the PSpad editor

	Offset (dec)	Byte sequence (hex)	Byte sequence meaning
GIF signature + version	00	47 49 46	ASCII string 'GIF' – "magic" file number (GIF signature)
	03	38 39 61	ASCII string '89a' – GIF version
Logical Screen Descriptor	06	01 00	Logical screen width: one pixel
	08	01 00	Logical screen height: one pixel
	10	80	Bit field: enable global color palette
	11	00	Background color index
	12	00	Pixel width to height ratio: 1:1
Global Color Table	13	FF 80 00	First color in the palette: orange (#ff8000)
	16	FF FF FF	Second color in the palette: white (#ffffff)

Image data

Offset (dec)	Byte sequence (hex)	Byte sequence meaning
19	2C	Image descriptor start marker
20	00 00	x-position of the left edge of the image: column zero
22	00 00	y-position of the top edge of the image: line zero
24	01 00	Image width: one pixel
26	01 00	Image height: one pixel
28	00	Bit field: no color palette and no interlacing
29	02	Initial LZW code size in bits minus 1
30	02	LZW-encoded block size in bytes
31	44 01	Encoded image data (one pixel with zero index)
33	00	Terminator of the LZW-encoded block
34	3B	GIF file terminator: the ';' character



1 / Structure of binary formats

1 x 1 pixel GIF

Parsing a binary file

```

0001 0203 0405 0607 0809 0A0B 0C0D 0E0F | 0123456789ABCDEF
00 4749 4638 3961 0100 0100 8000 00FF 8000 | GIF89a....€..`€.
10 FFFF FF2C 0000 0000 0100 0100 0002 0244 | ',.....D
20 0100 3B | ..;

```

Displaying the binary representation of a GIF file in the PSpad editor

	Offset (dec)	Byte sequence (hex)	Byte sequence meaning
GIF signature + version	00	47 49 46	ASCII string 'GIF' – "magic" file number (GIF signature)
	03	38 39 61	ASCII string '89a' – GIF version
Logical Screen Descriptor	06	01 00	Logical screen width: one pixel
	08	01 00	Logical screen height: one pixel
	10	80	Bit field: enable global color palette
	11	00	Background color index
Global Color Table	12	00	Pixel width to height ratio: 1:1
	13	FF 80 00	First color in the palette: orange (#ff8000)
	16	FF FF FF	Second color in the palette: white (#ffffff)

Image data

Offset (dec)	Byte sequence (hex)	Byte sequence meaning
19	2C	Image descriptor start marker
20	00 00	x-position of the left edge of the image: column zero
22	00 00	y-position of the top edge of the image: line zero
24	01 00	Image width: one pixel
26	01 00	Image height: one pixel
28	00	Bit field: no color palette and no interlacing
29	02	Initial LZW code size in bits minus 1
30	02	LZW-encoded block size in bytes
31	44 01	Encoded image data (one pixel with zero index)
33	00	Terminator of the LZW-encoded block
34	3B	GIF file terminator: the ';' character



1 / Structure of binary formats

1 x 1 pixel GIF

Parsing a binary file

```

0001 0203 0405 0607 0809 0A0B 0C0D 0E0F | 0123456789ABCDEF
00 4749 4638 3961 0100 0100 8000 00FF 8000 | GIF89a....€..`€.
10 FFFF FF2C 0000 0000 0100 0100 0002 0244 | ',.....D
20 0100 3B | ..;

```

Displaying the binary representation of a GIF file in the PSpad editor

	Offset (dec)	Byte sequence (hex)	Byte sequence meaning
GIF signature + version	00	47 49 46	ASCII string 'GIF' – "magic" file number (GIF signature)
	03	38 39 61	ASCII string '89a' – GIF version
Logical Screen Descriptor	06	01 00	Logical screen width: one pixel
	08	01 00	Logical screen height: one pixel
	10	80	Bit field: enable global color palette
	11	00	Background color index
	12	00	Pixel width to height ratio: 1:1
Global Color Table	13	FF 80 00	First color in the palette: orange (#ff8000)
	16	FF FF FF	Second color in the palette: white (#ffffff)

Image data

Offset (dec)	Byte sequence (hex)	Byte sequence meaning
19	2C	Image descriptor start marker
20	00 00	x-position of the left edge of the image: column zero
22	00 00	y-position of the top edge of the image: line zero
24	01 00	Image width: one pixel
26	01 00	Image height: one pixel
28	00	Bit field: no color palette and no interlacing
29	02	Initial LZW code size in bits minus 1
30	02	LZW-encoded block size in bytes
31	44 01	Encoded image data (one pixel with zero index)
33	00	Terminator of the LZW-encoded block
34	3B	GIF file terminator: the ';' character



1 / Structure of binary formats

1 x 1 pixel GIF

Parsing a binary file

```

0001 0203 0405 0607 0809 0A0B 0C0D 0E0F | 0123456789ABCDEF
00 4749 4638 3961 0100 0100 8000 00FF 8000 | GIF89a....€..`€.
10 FFFF FF2C 0000 0000 0100 0100 0002 0244 | ',.....D
20 0100 3B | ;

```

Displaying the binary representation of a GIF file in the PSpad editor

	Offset (dec)	Byte sequence (hex)	Byte sequence meaning
GIF signature + version	00	47 49 46	ASCII string 'GIF' – "magic" file number (GIF signature)
	03	38 39 61	ASCII string '89a' – GIF version
Logical Screen Descriptor	06	01 00	Logical screen width: one pixel
	08	01 00	Logical screen height: one pixel
	10	80	Bit field: enable global color palette
	11	00	Background color index
	12	00	Pixel width to height ratio: 1:1
Global Color Table	13	FF 80 00	First color in the palette: orange (#ff8000)
	16	FF FF FF	Second color in the palette: white (#ffffff)

Image data

Offset (dec)	Byte sequence (hex)	Byte sequence meaning
19	2C	Image descriptor start marker
20	00 00	x-position of the left edge of the image: column zero
22	00 00	y-position of the top edge of the image: line zero
24	01 00	Image width: one pixel
26	01 00	Image height: one pixel
28	00	Bit field: no color palette and no interlacing
29	02	Initial LZW code size in bits minus 1
30	02	LZW-encoded block size in bytes
31	44 01	Encoded image data (one pixel with zero index)
33	00	Terminator of the LZW-encoded block
34	3B	GIF file terminator: the ';' character

2

Binary analysis



2 / Parsing and serialization

GIF format

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F		
00000000	47	49	46	38	37	61	20	02	06	00	e7	00	00	18	b5	f7	GIF87a	...	ç	..	μ	+												
00000010	18	c6	ef	18	ce	ef	18	d6	ef	18	de	ef	18	e7	ef	18	.	Æ	ï	.	Î	.	Ö	ï	.	Ð	ï	.	ç	ï	.			
00000020	ef	ef	18	f7	e7	18	ff	21	18	ff	29	18	ff	de	18	ff	ï	ï	.	÷	ç	.	ÿ	!	.	ÿ)	.	ÿ	Þ	.	ÿ		
00000030	e7	21	31	f7	21	39	ff	21	42	ff	21	4a	ff	21	52	ff	ç	!	1	÷	!	9	ÿ	!	B	ÿ	!	J	ÿ	!	R	ÿ		
00000040	21	5a	f7	21	5a	ff	21	63	f7	21	6b	f7	21	73	f7	21	!	Z	÷	!	Z	ÿ	!	c	÷	!	k	÷	!	s	÷	!		
00000050	7b	f7	21	84	f7	21	8c	f7	21	94	f7	21	9c	f7	21	a5	{	÷	!	.	÷	!	.	÷	!	.	÷	!	.	÷	!	¥		
00000060	f7	21	ad	f7	21	b5	f7	21	ff	18	21	ff	29	21	ff	31	÷	!	.	÷	!	μ	÷	!	ÿ	.	ÿ	!	ÿ	!	ÿ	!		
00000070	21	ff	39	21	ff	42	21	ff	4a	21	ff	52	21	ff	5a	21	!	ÿ	9	!	ÿ	B	!	ÿ	J	!	ÿ	R	!	ÿ	Z	!		
00000080	ff	63	21	ff	6b	21	ff	73	21	ff	7b	21	ff	84	21	ff	ÿ	c	!	ÿ	k	!	ÿ	s	!	ÿ	{	!	ÿ	.	ÿ			



serialization
(writing)

```

header [Header]
├─ LogicalScreenDescriptor [LogicalScreenDescriptor]
│   ├── screenWidth = 0x220 = 544
│   ├── screenHeight = 0x6 = 6
│   ├── flags = 0xE7 = 231
│   ├── bgColorIndex = 0x0 = 0
│   ├── pixelAspectRatio = 0x0 = 0
│   ├── hasColorTable = true
│   └─ colorTableSize = 0x100 = 256
├─ globalColorTable [GlobalColorTable]
│   └─ entries
│       ├── 0 [ColorTableEntry]
│       ├── 1 [ColorTableEntry]
│       ├── 2 [ColorTableEntry]
│       └─ 3 [ColorTableEntry]

```



parsing
(reading)



2 / Getting a parser

Dedicated format library

Own parser

Parser combinator



Construct (Python)
Binary-parser (JavaScript)
BinData (Ruby)

Parser generator

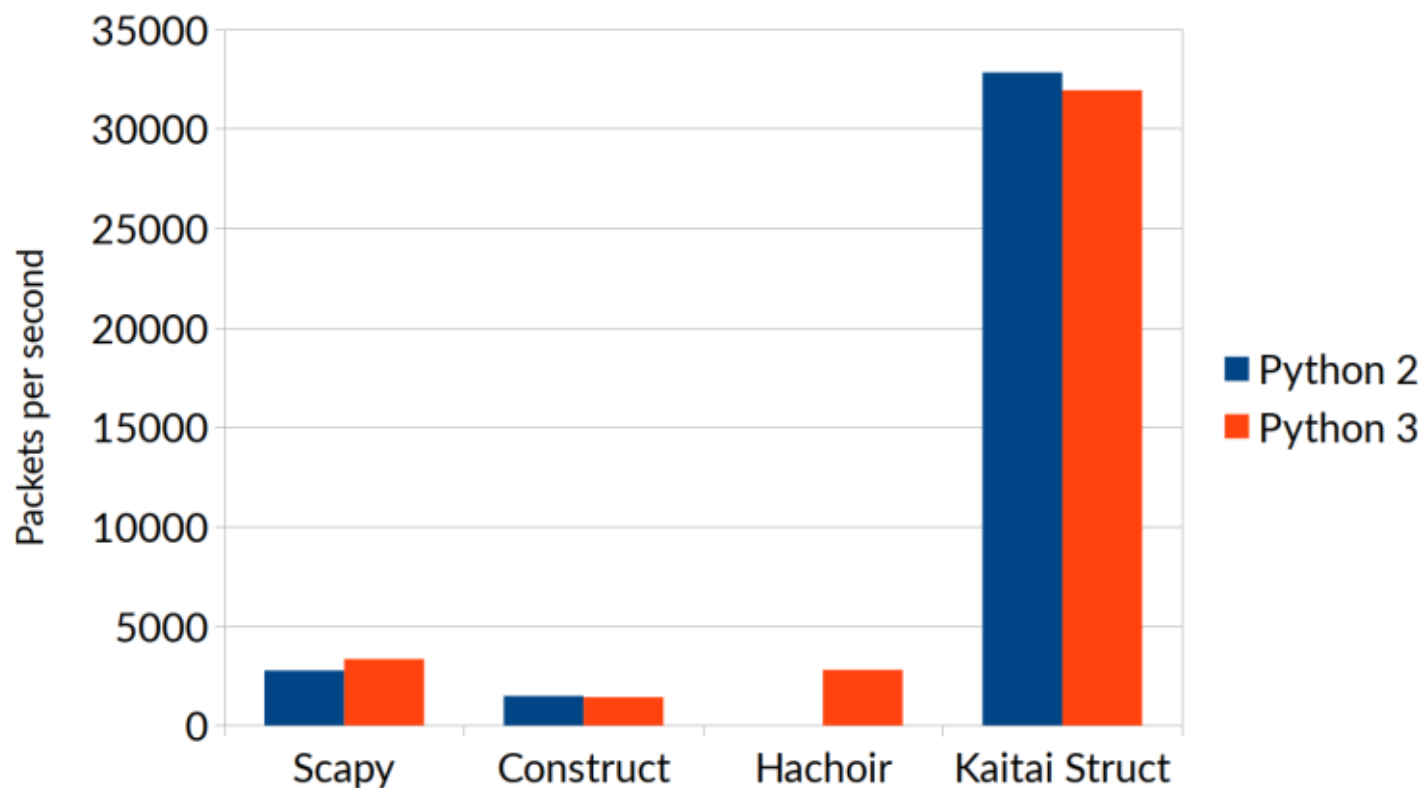


Kaitai Struct
Spicy (C++)
EverParse (C)
Apache Daffodil (Java, C)



2 / Parser combinators vs. generators

<https://pythonistac.wordpress.com/2017/03/09/python-network-packet-dissection-frameworks-shootout-scapy-vs-construct-vs-hachoir-vs-kaitai-struct/>



Comparison of the speed of combinators (Scapy, Construct and Hachoir) and the Kaitai Struct generator in processing network packets



The screenshot shows the Kaitai Struct website homepage. The navigation bar includes links for 'Kaitai Struct', 'What is it?', 'Quick Start', 'Download', 'News', 'Format Gallery', 'Try it - Web IDE', and 'Documentation'. The main content area features a large header with the Kaitai Struct logo and the title 'Kaitai Struct'. Below the title is a tagline: 'A new way to develop parsers for binary structures.' The page is divided into several sections, each with an icon and a title:

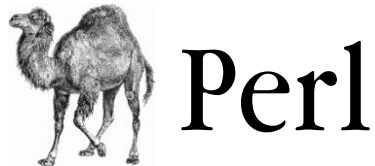
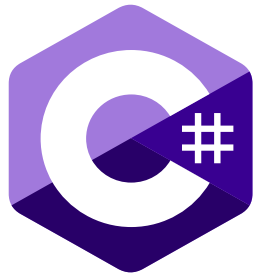
- Declarative:** describe the very structure of the data, not how you read or write it. (Icon: tree structure)
- Language-neutral:** write once, use in all supported languages: (Icon: Kaitai Struct logo with 'A' and a refresh arrow)
 - C++/STL
 - C#
 - Go
 - Java
 - JavaScript
 - Lua
 - Nim
 - Perl
 - PHP
 - Python
 - Ruby

entry-level support
- Packed with tools and samples:** includes a compiler, an IDE, a visualizer and massive library of popular formats. (Icon: battery)
- Free & open source:** feel free to use, modify and join the project. (Icon: open lock)

At the bottom left, there is a box indicating '0.10 released 2022-07-08' with a 'Download' button.

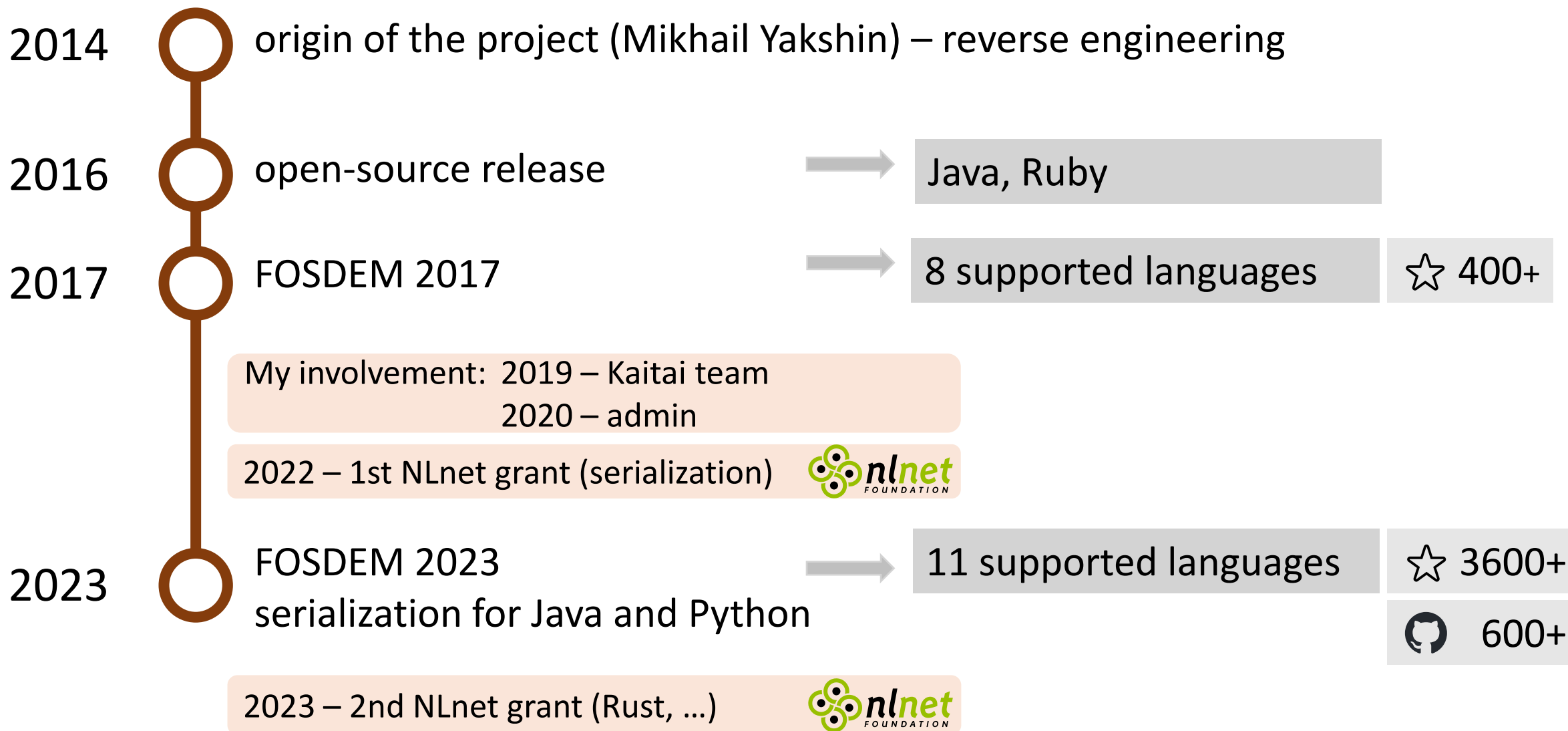
力 3 / What is Kaitai Struct?

- parser generator for 11 programming languages
- other parsing targets – Construct, Awkward Arrays
- serialization in 2 languages – Java and Python
- declarative language (.ksy) for specifying arbitrary binary formats





3 / History of Kaitai Struct





3 / Kaitai workflow

1. Compilation

hello_world.ksy

```
meta:  
  id: hello_world  
seq:  
  - id: one  
    type: u1
```

kaitai-struct-compiler

hello_world.py

```
class HelloWorld(KaitaiStruct):  
  # ...  
  def _read(self):  
    self.one = self._io.read_u1()
```

2. Parsing

input binary file

```
sample.bin  
0 1 2 3 4 5 6 7 01234567  
00000000 ff 01 y.
```

hello_world.py

```
class HelloWorld(KaitaiStruct):  
  # ...  
  def _read(self):  
    self.one = self._io.read_u1()
```

parsed data

```
└ one = 0xFF = 255
```

kaitaistruct.py (runtime library)

```
class KaitaiStream:  
  # ...  
  def read_u1(self):  
    return struct.unpack('B', self.read_bytes(1))[0]
```



3 / Kaitai features

- 1 format spec = 11 parsers
- standard way to describe binary formats
- library of format specifications
- GraphViz diagram
- .ksy language simplicity
- serialization
- testing
- visualization and dumping tools
 - console visualizer (ksv)
 - ksdump
 - Web IDE
- distribution



3 / Kaitai features

→ 1 .ksy specification = 11 parsers

.ksy spec

```
meta:  
  id: hello_world  
seq:  
  - id: one  
    type: u1
```

```
public class HelloWorld extends KaitaiStruct {  
  // ...  
  private void _read() {  
    this.one = this._io.readU1();  
  }  
}
```

Java

```
class HelloWorld(KaitaiStruct):  
  # ...  
  def _read(self):  
    self.one = self._io.read_u1()
```

Python

```
class HelloWorld < Kaitai::Struct::Struct  
  # ...  
  def _read  
    @one = @_io.read_u1  
  end
```

Ruby

and others (C++, C#, Go, JavaScript, Lua, Nim, Perl, PHP)

→ standard way to describe binary formats

GIF

18. Logical Screen Descriptor.

a. Description. The Logical Screen Descriptor contains the parameters necessary to define the area of the display device within which the images will be rendered. The coordinates in this block are given with respect to the top-left corner of the virtual screen; they do not necessarily refer to absolute coordinates on the display device. This implies that they could refer to window coordinates in a window-based environment or printer coordinates when a printer is used.

This block is REQUIRED; exactly one Logical Screen Descriptor must be present per Data Stream.

b. Required Version. Not applicable. This block is not subject to a version number. This block must appear immediately after the Header.

c. Syntax.

	7 6 5 4 3 2 1 0	Field Name	Type
0	-----	Logical Screen Width	Unsigned
1	-----		
2	-----	Logical Screen Height	Unsigned
3	-----		
4		<Packed Fields>	See below
5	-----	Background Color Index	Byte
6	-----	Pixel Aspect Ratio	Byte

Microsoft Word .doc

2.9.161 OcxInfo

The **OcxInfo** structure specifies an **OLE control** (such as a checkbox, radio button, and so on) in the document. The data that is contained in **OcxInfo** structures SHOULD [<229>](#) be ignored.

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31								
dwCookie																																							
ifld																																							
hAccel																																							
cAccel																A	B	C	D	E	F	G	H																
idoc																reserved																							

dwCookie (4 bytes): An integer value that specifies the index location of this **OcxInfo** in the [RgxOcxInfo](#) array. This value MUST be unique for all **OcxInfo** structures in the document.

ifld (4 bytes): An unsigned integer value that specifies an index location in the [PlcFld](#) structure. The value MUST be a valid [FLD](#) index in the correct **PlcFld** structure.







The PlcFld that is used is dependent on the value of **idoc**, as specified following.







Value	Location
1	The Main Document (FibRgFcLcb97.fcPlcfFldMom).
2	The Header Document (FibRgFcLcb97.fcPlcfFldHdr).
3	The Footnote Document (FibRgFcLcb97.fcPlcfFldFtn).
4	The Textbox Document (FibRgFcLcb97.fcPlcfFldTxbx).
6	The Endnote Document (FibRgFcLcb97.fcPlcfFldEdn).
7	The Comment Document (FibRgFcLcb97.fcPlcfFldAtn).
8	The Header Textbox Document (FibRgFcLcb97.fcPlcfHdrtxbxTxt).

➔ library of .ksy format specifications

185 format specifications by 81 contributors

formats.kaitai.io

-  3D Models
gltf_binary , quake_md1
-  Archive Files
android_bootldr_asus , android_bootldr_huawei , android_bootldr_qcom , android_dto , android_img , android_sparse , chrome_pak , cpio_old_le , gzip , lzh , mozilla_mar , phar_without_stub , rar , rpm , xar , zip , zisofs
-  Commonly Used Data Types
bcd , bytes_with_io , dos_datetime , riff , utf8_string , vlq_base128_be , vlq_base128_le
-  DOS-specific
dos_datetime , dos_mz , mbr_partition_table , vfat
-  Filesystems
android_super , apm_partition_table , apple_single_double , btrfs_stream , cramfs , ext2 , gpt_partition_table , iso9660 , luks , lvm2 , mbr_partition_table , tr_dos_image , vdi , vfat , vmware_vmdk , zisofs , zx_spectrum_tap
-  Fonts
grub2_font , pcf_font , ttf

-  Android-specific
android_bootldr_asus , android_bootldr_huawei , android_bootldr_qcom , android_img , android_nanoapp_header , android_opengl_shaders_cache , android_sparse , android_super , dex
-  CAD
monomakh_sapr_chg
-  Databases
dbf , gettext_mo , sqlite3 , tsm
-  Executables and Byte-code
android_nanoapp_header , dex , dos_mz , elf , java_class , mach_o , mach_o_fat , microsoft_pe , python_pyc_27 , swf , uefi_te
-  Firmware
andes_firmware , broadcom_trx , ines , uefi_te , uimage
-  Game Data Files
allegro_dat , doom_wad , dune_2_pak , fallout2_dat , fallout_dat , ftl_dat , gran_turismo_vol , heaps_pak , heroes_of_might_and_magic_agg , heroes_of_might_and_magic_bmp , minecraft_nbt , quake_md1 , quake_pak , renderware_binary_stream , saints_row_2_vpp_pc , warcraft_2_pud

→ library of .ksy format specifications

185 format specifications by 81 contributors

formats.kaitai.io

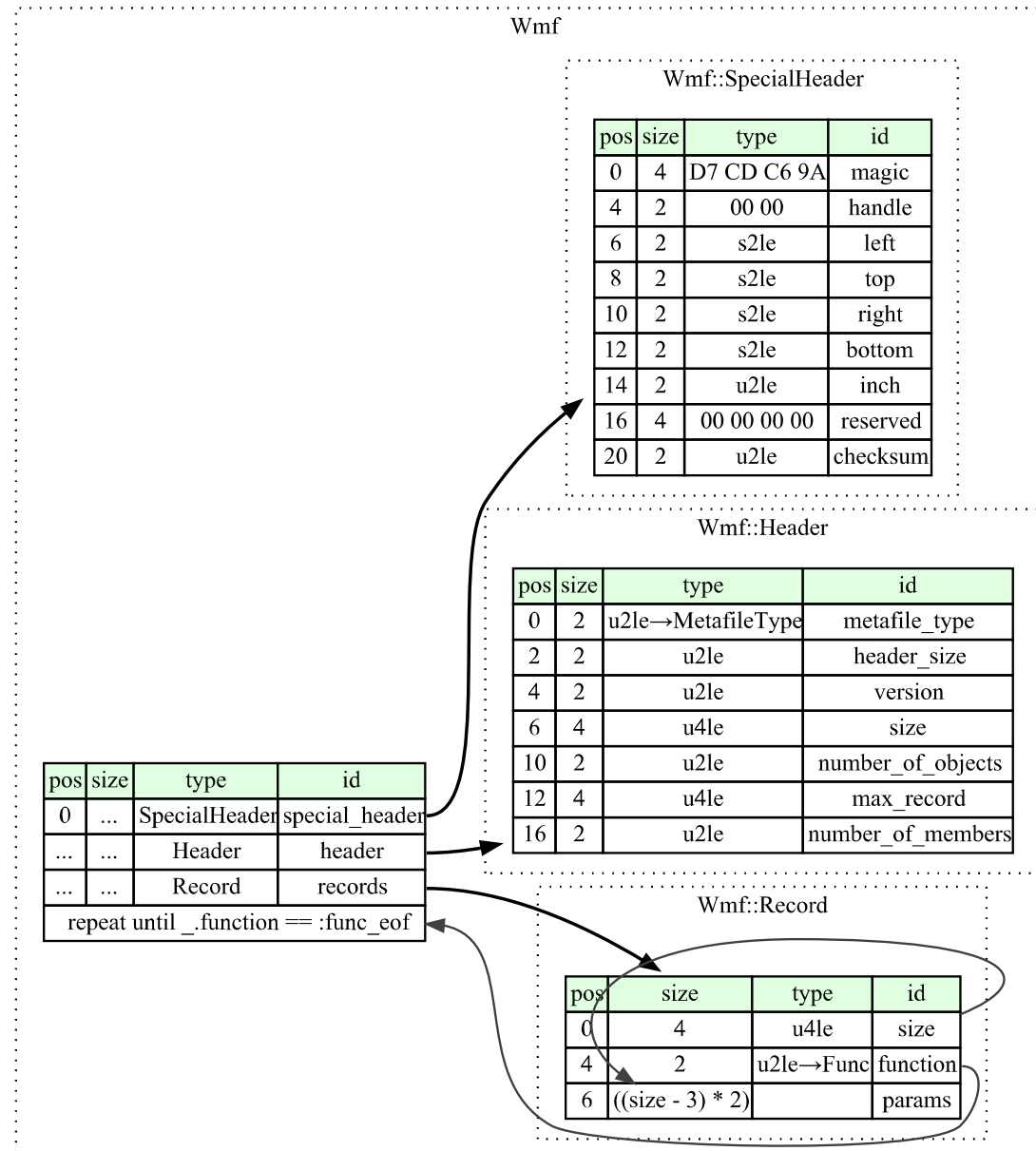
-  Geospatial (Maps)
`shapefile_index` , `shapefile_main`
-  Image Files
`bmp` , `dicom` , `exif` , `gif` , `gimp_brush` , `icc_4` , `ico` , `jpeg` , `nitf` , `pcx` , `pcx_dcx` , `png` , `psx_tim` , `tga` , `wmf` , `xwd`
-  Logs
`aix_utmp` , `glibc_utmp` , `hashcat_restore` , `mcap` , `sudoers_ts` , `systemd_journal` , `windows_evt_log`
-  macOS-specific
`apm_partition_table` , `apple_single_double` , `compressed_resource` , `dcmp_0` , `dcmp_1` , `dcmp_2` , `dcmp_variable_length_integer` , `ds_store` , `mac_os_resource_snd` , `resource_fork`
-  Networking Protocols
`bitcoin_transaction` , `dime_message` , `dns_packet` , `ethernet_frame` , `hccap` , `hccapx` , `icmp_packet` , `ipv4_packet` , `ipv6_packet` , `microsoft_network_monitor_v2` , `packet_ppi` , `pcap` , `protocol_body` , `rtcp_payload` , `rtp_packet` , `rtpdump` , `some_ip` , `some_ip_container` , `some_ip_sd` , `some_ip_sd_entries` , `some_ip_sd_options` , `tcp_segment` , `tls_client_hello` , `udp_datagram` , `websocket`
-  Security
`efivar_signature_list` , `openpgp_message` , `ssh_public_key`
-  Windows-specific
`avi` , `bmp` , `ico` , `microsoft_pe` , `regf` , `wav` , `windows_evt_log` , `windows_lnk_file` , `windows_minidump` , `windows_resource_file` , `windows_shell_items` , `windows_systemtime` , `wmf`

-  Hardware Protocols
`dtb` , `edid` , `mifare_classic`
-  GNU/Linux-specific
`btrfs_stream` , `cramfs` , `dtb` , `elf` , `ext2` , `gettext_mo` , `glibc_utmp` , `luks` , `lvm2` , `sudoers_ts` , `systemd_journal`
-  CPU / Machine Code Disassembly
`code_6502`
-  Multimedia Files
`android_opengl_shaders_cache` , `au` , `avi` , `blender_blend` , `creative_voice_file` , `fasttracker_xm_module` , `genmidi_op2` , `id3v1_1` , `id3v2_3` , `id3v2_4` , `magicavoxel_vox` , `ogg` , `quicktime_mov` , `s3m` , `standard_midi_file` , `stl` , `swf` , `vp8_ivf` , `wav`
-  Scientific Applications
`avantes_roh60` , `nt_mdt` , `nt_mdt_pal` , `specpr`
-  Serialization Protocols
`asn1_der` , `bson` , `chrome_pak` , `dtb` , `google_protobuf` , `microsoft_cfb` , `minecraft_nbt` , `msgpack` , `php_serialized_value` , `python_pickle` , `ruby_marshall`



3 / Kaitai features

→ GraphViz diagram





3 / Kaitai features

➔ .ksy language simplicity

but powerful: unaligned bit types, type switch, byte processing, imports, ...

```

meta:
  id: ftl_dat
  endian: le little endian
seq:
  - id: num_files attribute name
    type: u4 unsigned 4-byte integer
  - id: files
    type: file
    repeat: expr
    repeat-expr: num_files
                    number of repetitions

```

```

types:
  file:
    seq:
      - id: ofs_data
        type: u4
    instances:
      data: byte offset
      pos: ofs_data
      type: file_data
      if: ofs_data != 0
          expression language

```

```

file_data:
  seq:
    - id: len_file
      type: u4
    - id: len_filename
      type: u4
    - id: filename
      size: len_filename
      type: str char. string
      encoding: UTF-8
    - id: body
      size: len_file
no type => byte array

```



→ serialization



1. **editing** an existing file
2. creating a **new** file

Areas of application

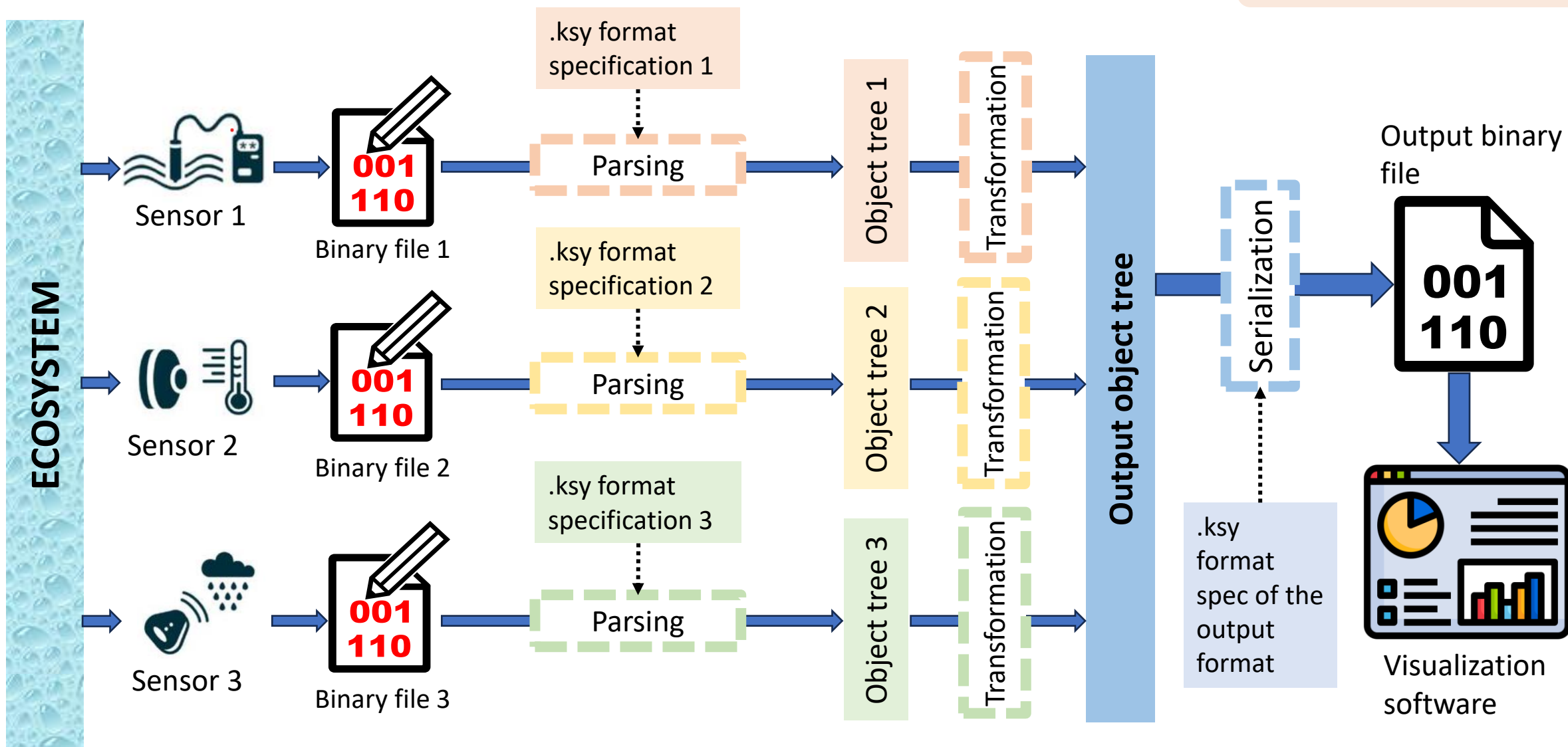
- format conversions (parse → transform → serialize to another format)
- fuzzing
- video games modding



3 / Kaitai features

→ serialization

format conversions





3 / Kaitai features

→ testing

292 tests

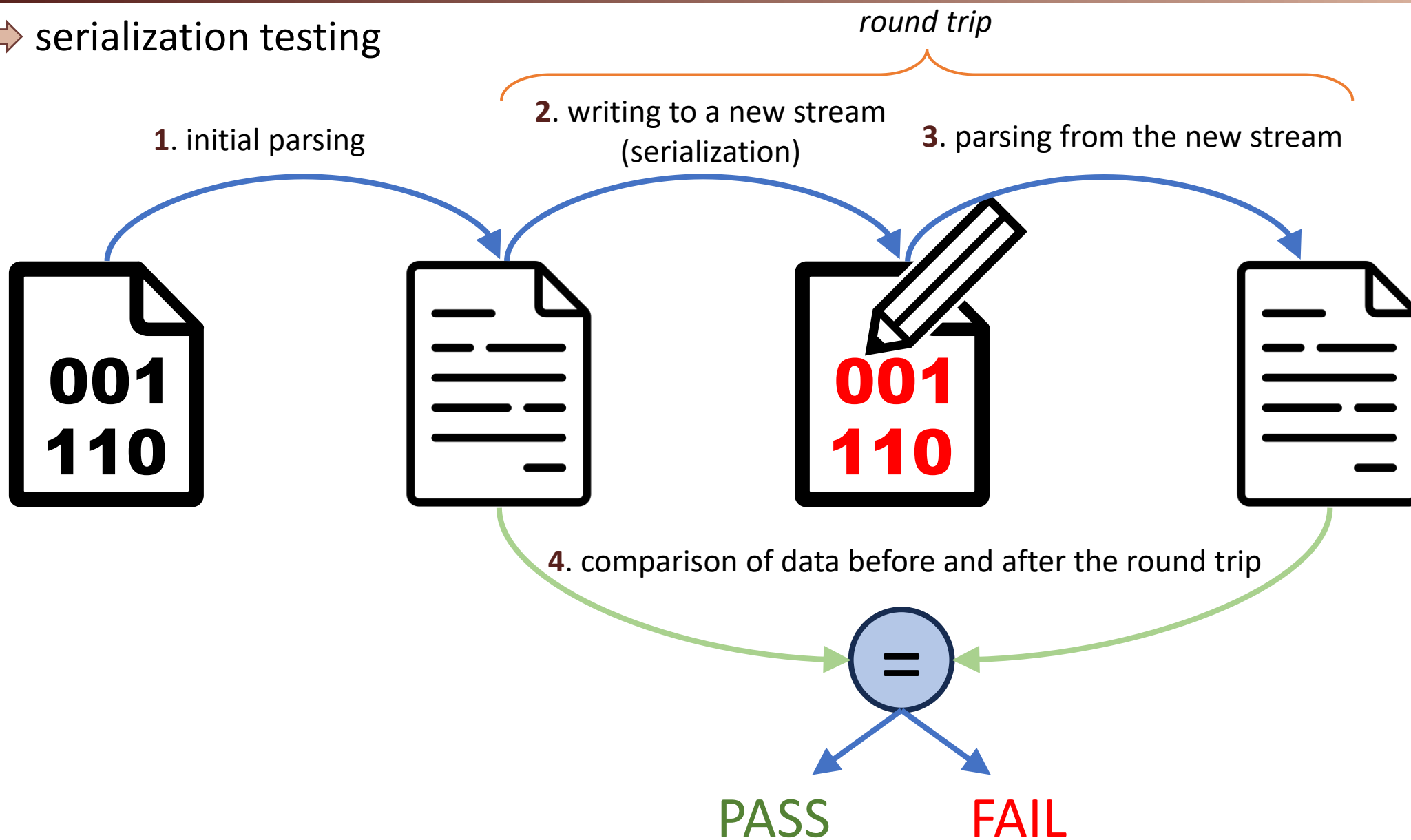
The screenshot shows the Kaitai Struct CI website interface. At the top, there's a browser tab for 'Kaitai Struct CI' and a URL 'ci.kaitai.io'. Below the header, there are search bars for 'Search test...' and 'Search target...'. The main content is a table with columns for 'Test \ Target', 'cpp_stl_98', 'cpp_stl_11', 'csharp', 'graphviz', 'go', 'java', 'javascript', and 'lisp'. The table lists various tests and their results across these targets. For example, 'BodUserTypeBe' passed on all targets except 'graphviz' where it failed with 'format_build_failed'. 'BitsShiftByB64Le' failed on the 'javascript' target.

Test \ Target	cpp_stl_98	cpp_stl_11	csharp	graphviz	go	java	javascript	lisp
Last update	9d ago Job Artifacts	9d ago Job Artifacts	9d ago Job Artifacts	9d ago Job Artifacts	9d ago Job Artifacts	9d ago Job Artifacts	9d ago Job Artifacts	9d ago Job Artifacts
Rating	91.9%	92.8%	91.1%	40%	77.4%	94.5%	91.9%	91.9%
BodUserTypeBe	passed	passed	passed	format_build_failed	passed	passed	passed	passed
BodUserTypeLe	passed	passed	passed	format_build_failed	passed	passed	passed	passed
BitsByteAligned	passed	passed	passed	passed	passed	passed	passed	passed
BitsEnum	passed	passed	passed	passed	passed	passed	passed	passed
BitsSeqEndianCombo	passed	passed	passed	format_build_failed	passed	passed	passed	passed
BitsShiftByB32Le	passed	passed	passed	passed	passed	passed	passed	passed
BitsShiftByB64Le	passed	passed	passed	passed	passed	passed	failed	passed
BitsSignedResB32Be	passed	passed	passed	passed	passed	passed	passed	passed



3 / Kaitai features

→ serialization testing





3 / Kaitai features

→ console visualizer (ksv)

github.com/kaitai-io/kaitai_struct_visualizer

```

[-] [root]
[.] magic = 89 50 4e 47 0d 0a 1a 0a
[.] ihdr_len = 13
[.] ihdr_type = 49 48 44 52
[-] ihdr
[.] width = 300
[.] height = 300
[.] bit_depth = 8
[.] color_type = color_type_truecolor
[.] compression_method = 0
[.] filter_method = 0
[.] interlace_method = 0
[.] ihdr_crc = f6 1f 19 22
[-] chunks (2 = 0x2 entries)
[-] 0
[.] len = 919
[.] type = "IDAT"
[.] body = 78 9c ed d9 31 8a c3 40 14 44 c1 1e e3 fb 5f...
[.] crc = 21 cb 54 d1
[-] 1
[.] len = 0
[.] type = "IEND"
[.] body =
[.] crc = ae 42 60 82
00000000: 89 50 4e 47 0d 0a 1a 0a 00 00 00 0d 49 48 44 52 | .PNG.....IHDR
00000010: 00 00 01 2c 00 00 01 2c 08 02 00 00 00 f6 1f 19 | ...,.,.....
00000020: 22 00 00 03 97 49 44 41 54 78 9c ed d9 31 8a c3 | "....IDATx...1..
00000030: 40 14 44 c1 1e e3 fb 5f 59 8a 9d 09 1c bc 40 55 | @.D...._Y.....@U
00000040: 6c b4 20 70 f2 68 98 7f b6 6b bb ce ef df b6 f3 | \. p.h...k.....
00000050: e8 9f f3 ad 6f 7d fb e7 b7 9f 01 a9 ef 4e fd 13 | .....o}.....N..
00000060: e0 dd 44 08 31 11 42 4c 84 10 13 21 c4 bc 8e 42 | ..D.1.BL...!...B
00000070: cc 12 42 4c 84 10 13 21 c4 44 08 31 11 42 4c 84 | ..BL...!.D.1.BL.
00000080: 10 73 a2 80 98 25 84 98 08 21 26 42 88 89 10 62 | .s...%...!&B...b
00000090: 22 84 98 08 21 e6 44 01 31 4b 08 31 11 42 4c 84 | "...!.D.1K.1.BL.
000000a0: 10 13 21 c4 44 08 31 af a3 10 b3 84 10 13 21 c4 | ..!.D.1.....!.
000000b0: 44 08 31 11 42 4c 84 10 13 21 c4 9c 28 20 66 09 | D.1.BL...!..( f.
000000c0: 21 26 42 88 89 10 62 22 84 98 08 21 26 42 88 39 | !&B...b"...!&B.9
000000d0: 51 40 cc 12 42 4c 84 10 13 21 c4 44 08 31 11 42 | Q@..BL...!.D.1.B
000000e0: cc eb 28 c4 2c 21 c4 44 08 31 11 42 4c 84 10 13 | ..(.,!.D.1.BL...
000000f0: 21 c4 44 08 31 27 0a 88 59 42 88 89 10 62 22 84 | !.D.1'..YB...b".
00000100: 98 08 21 26 42 88 89 10 62 4e 14 10 b3 84 10 13 | ..!&B...bN.....
00000110: 21 c4 44 08 31 11 42 4c 84 10 f3 3a 0a 31 4b 08 | !.D.1.BL...:..1K.
00000120: 31 11 42 4c 84 10 13 21 c4 44 08 31 11 42 cc 89 | 1.BL...!.D.1.B..
00000130: 02 62 96 10 62 22 84 98 08 21 26 42 88 89 10 62 | .b..b"...!&B...b
00000140: 22 84 98 13 05 c4 2c 21 c4 44 08 31 11 42 4c 84 | "...,,!.D.1.BL.
00000150: 10 13 21 c4 bc 8e 42 cc 12 42 4c 84 10 13 21 c4 | ..!...B..BL...!.
00000160: 44 08 31 11 42 4c 84 10 73 a2 80 98 25 84 98 08 | D.1.BL..s...%...
00000170: 21 26 42 88 89 10 62 22 84 98 08 21 e6 44 01 31 | !&B...b"...!.D.1
00000180: 4b 08 31 11 42 4c 84 10 13 21 c4 44 08 31 af a3 | K.1.BL...!.D.1..
00000190: 10 b3 84 10 13 21 c4 44 08 31 11 42 4c 84 10 13 | .....!.D.1.BL...
000001a0: 21 c4 9c 28 20 66 09 21 26 42 88 89 10 62 22 84 | !..( f.!&B...b".
000001b0: 98 08 21 26 42 88 39 51 40 cc 12 42 4c 84 10 13 | ..!&B.9Q@..BL...

```

→ ksdump (JSON output)

github.com/kaitai-io/kaitai_struct_visualizer

```
{
  "magic": "89 50 4E 47 0D 0A 1A 0A",
  "ihdr_len": 13,
  "ihdr_type": "49 48 44 52",
  "ihdr": {
    "width": 300,
    "height": 300,
    "bit_depth": 8,
    "color_type": "color_type_truecolor",
    "compression_method": 0,
    "filter_method": 0,
    "interlace_method": 0
  },
  "ihdr_crc": "F6 1F 19 22",
  "chunks": [
    {
      "len": 919,
      "type": "IDAT",
      "body": "78 9C ED D9 31 8A C3 40 14 44 C1 1E E3 FB 5F...",
      "crc": "21 CB 54 D1"
    },
    {
      "len": 0,
      "type": "IEND",
      "body": "",
      "crc": "AE 42 60 82"
    }
  ]
}
```

3 / Kaitai features

→ Web IDE

ide.kaitai.io

The screenshot displays the Kaitai Web IDE interface with several key components highlighted:

- .ksy format spec:** A code editor showing the definition for the PNG format in `formats/image/png.ksy`. The spec includes fields like `id`, `title`, `file-extension`, `license`, `ks-version`, `endian`, `doc`, and `seq`.
- object tree:** A tree view showing the parsed structure of the `pnggrad8rgb.png` file. It includes fields like `magic`, `ihdrLen`, `ihdrType`, `ihdr` (with sub-fields like `width`, `height`, `bitDepth`, `colorType`, `compressionMethod`, `filterMethod`, `interlaceMethod`, `ihdrCrc`), and `chunks`.
- hex viewer:** A hex viewer showing the raw bytes of the `pnggrad8rgb.png` file. The selected range is `0x14 - 0x17`, which corresponds to the `height` field in the object tree.
- converter:** A panel showing the conversion of the selected hex data to various data types. The selected data is `0x14 - 0x17`, which is converted to `300` in decimal, `300` in hexadecimal, and `1.8332002582610585e-12` in float.

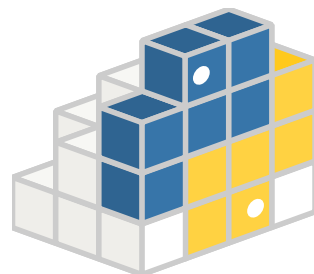


3 / Kaitai features

→ distribution – runtime libraries



npm – JavaScript



PyPI – Python



NuGet – C#



Maven Central
Repository – Java



RubyGems – Ruby



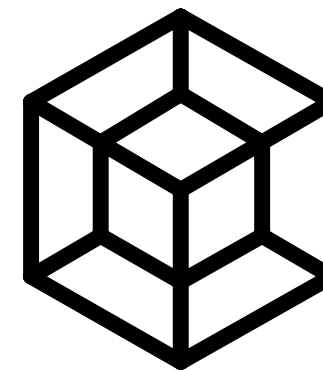
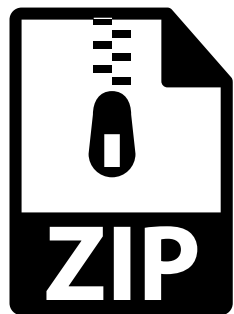
Packagist – PHP



3 / Kaitai features

→ distribution – compiler

<https://repology.org/project/kaitai-struct/versions>



WinGet



Homebrew



Void Linux



AUR



nixpkgs

4

Other binary analysis tools

Dirk Loss

list of parsing tools, hex editors...

<https://github.com/dloss/binary-parsing>

dloss/binary-parsing



A list of generic tools for parsing binary data structures, such as file formats, network protocols or bitstreams

11

Contributors

2

Issues

781

Stars

62

Forks





Awesome binary parsing [↗](#)

A list of generic tools for parsing binary data structures, such as file formats, network protocols or bitstreams.

Parser generators, parsing libraries and frameworks [↗](#)

- [Kaitai Struct](#) (DSL): declarative language used for describe various binary data structures, laid out in files or in memory
- [Nom](#) (Rust): Rust parser combinator framework
- [Hammer](#) (C): bit-oriented parsing library
- [Construct](#) (Python): library for parsing and building of data structures (binary or textual). Define your data structures in a declarative manner
- [Spicy](#) (DSL, C/C++, Zeek): a next-generation parser generator for network protocols and file formats
- [Hachoir](#) (Python): view and edit a binary stream field by field. Long [list of parsers](#) for all



4 / Binary analysis tools

decoding and encoding.

- [EverParse](#): a framework for generating verified secure parsers and formatters from domain-specific format specification languages

Stand-alone software [↗](#)

Hex editors with grammars [↗](#)

- [Synalyze It!](#)
- [Hexinator](#)
- [010 Editor](#)
- [Kiewtai](#): plugin for the Hiew hex editor that makes the Kaitai parsers available
- [Hobbits](#): multi-platform GUI for bit-based analysis, processing, and visualization. Has a Kaitai plugin.
- [ImHex](#): A Hex Editor for Reverse Engineers, Programmers and people who value their retinas when working at 3 AM.

Wireshark [↗](#)



4 / Binary analysis tools

- [Hobbits](#): multi-platform GUI for bit-based analysis, processing, and visualization. Has a Kaitai plugin.
- [ImHex](#): A Hex Editor for Reverse Engineers, Programmers and people who value their retinas when working at 3 AM.

Wireshark

[Wireshark](#) is a network protocol analyzer that includes [dissectors](#) for over [two thousand](#) protocols.

- [TShark](#): command line version, can easily be called from shell scripts.
- [Wireshark Generic Dissector](#): add-on, allows dissection of a protocol based on a text description of the protocol elements
- [Wireshark Lua](#): dissectors can be written in Lua ([Examples](#))
- [pyreshark](#): plugin providing a simple interface for writing Wireshark dissectors in Python
- [Sharktools](#) (Python, Matlab): Tools for programmatic parsing of packet captures using Wireshark functionality

Other Stand-alone Software



4 / Binary analysis tools

- [Wireshark Lua](#): dissectors can be written in Lua ([Examples](#))
- [pyreshark](#): plugin providing a simple interface for writing Wireshark dissectors in Python
- [Sharktools](#) (Python, Matlab): Tools for programmatic parsing of packet captures using Wireshark functionality

Other Stand-alone Software [↗](#)

- [GNU poke](#): The extensible editor for structured binary data
- [Netzob](#): open source tool for reverse engineering, traffic generation and fuzzing of communication protocols
- [Cat Karat Packet Builder](#): packet generation tool that allows to build custom packets for firewall or target testing
- [radare2](#) (C, with bindings/pipe for almost all languages): Unix-like reverse engineering framework and commandline tools. See [Parsing a fileformat with radare2](#) and [Types](#).
- [Veles](#): open source tool for binary analysis

Research papers [↗](#)

- [LangSec Platform](#): Towards a Platform to Compare Binary Parser Generators. Olivier
Levillain, Sébastien Naud, Aïme Talu, Pascal Boyer ([Video](#))



Parser generators, parsing libraries and frameworks [↗](#)

- [Kaitai Struct](#) (DSL): declarative language used for describe various binary data structures, laid out in files or in memory
- [Nom](#) (Rust): Rust parser combinator framework
- [Hammer](#) (C): bit-oriented parsing library
- [Construct](#) (Python): library for parsing and building of data structures (binary or textual). Define your data structures in a declarative manner
- [Spicy](#) (DSL, C/C++, Zeek): a next-generation parser generator for network protocols and file formats
- [Hachoir](#) (Python): view and edit a binary stream field by field. Long [list of parsers](#) for all kinds of formats
- [RecordFlux](#): toolset for the formal specification of messages and the generation of verifiable binary parsers and message generators (Ada-inspired).
- [DataScript Tools](#) (DSL): DataScript is a formal language for modelling binary datatypes,



4 / Binary analysis tools

decoding and encoding.

- [EverParse](#): a framework for generating verified secure parsers and formatters from domain-specific format specification languages

Stand-alone software [↗](#)

Hex editors with grammars [↗](#)

- [Synalyze It!](#)
- [Hexinator](#)
- [010 Editor](#)
- [Kiewtai](#): plugin for the Hiew hex editor that makes the Kaitai parsers available
- [Hobbits](#): multi-platform GUI for bit-based analysis, processing, and visualization. Has a Kaitai plugin.
- [ImHex](#): A Hex Editor for Reverse Engineers, Programmers and people who value their retinas when working at 3 AM.

Wireshark [↗](#)



fq

Mattias Wadman

inspired by `jq`, a tool for working with JSON data
dealing with multimedia formats

github.com/wader/fq

wader/fq

jq for binary formats - tool, language and decoders
for working with binary and text formats



26

Contributors

19

Used by

9k

Stars

214

Forks





BANG

Armijn Hemel
firmware analysis

- provenance detection
- security scans

github.com/armijnhemel/binaryanalysis-ng

armijnhemel/
binaryanalysis-ng

Binary Analysis Next Generation (BANG)



6

Contributors

18

Issues

428

Stars

66

Forks



Corkami posters

Ange Albertini

MP4

00x	00	00	00	1C	f	t	y	p	i	s	o	m	00	00	02	00
01x	i	s	o	m	i	s	o	2	m	p	4	1	00	00	01	16
02x	m	o	o	v	00	00	01	0E	t	r	a	k	00	00	01	06
03x	m	d	i	a	00	00	00	FE	m	i	n	f	00	00	00	F6
04x	s	t	b	l	00	00	00	92	s	t	s	d	00	00	00	00
05x	00	00	00	01	00	00	00	82	m	p	4	v	00	00	00	00
06x	00	00	00	01	00	00	00	00	00	00	00	00	00	00	00	00
07x	00	00	00	00	00	0D	00	07	00	48	00	00	00	48	00	00
08x	00	00	00	00	00	01	06	F	F	M	P	e	g	00	00	00
09x	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00
0Ax	00	00	00	00	00	00	00	18	FF	FF	00	00	00	2C	e	s
0Bx	d	s	00	00	00	00	03	80	80	80	18	00	01	00	04	80
0Cx	80	80	0D	6D	11	00	00	00	00	41	A0	00	00	41	A0	00
0Dx	06	80	80	80	01	02	00	00	00	18	s	t	t	s	00	00
0Ex	00	00	00	00	00	01	00	00	00	01	00	00	02	00	00	00
0Fx	00	1C	s	t	s	c	00	00	00	00	00	00	00	01	00	00
10x	00	01	00	00	00	01	00	00	00	01	00	00	00	14	s	t
11x	s	z	00	00	00	00	00	00	00	54	00	00	00	01	00	00
12x	00	14	s	t	c	o	00	00	00	00	00	00	00	01	00	00
13x	01	3A	00	00	00	5C	m	d	a	t	89	P	N	G	\r	\n
18x	D7	F7	00	00	00	00	I	E	N	D	AE	42	60	82	E0	F

AN MP4 FILE

0	Size 1C	Type ftyp	FILE TYPE	Major Brand isom	Minor Version 2.00	Compat. Brand isom iso2 mp41
1C	Size 116	Type moov	MOVIE			
24	Size 10E	Type trak	TRACK			
2C	Size 106	Type mdia	MEDIA			
34	Size FE	Type minf	MEDIA INFORMATION			
3C	Size F6	Type stbl	SAMPLE TABLE			
44	Size 92	Type stsd	SAMPLE DESCRIPTION	Version:0 Flags:0 Entries:1		
54	Size 82	Type mp4v	Mp4 VISUAL SAMPLE	Size:0 Track:1 Reserved Reserved Index:0 Uncomp:0 Revision:0 Reserved Width:13 Height:7 Res:72x72 Size:0 FrameCount:1 EncoderLen:6 Encoder 24b Marker		
AA	Size 2C	Type esds	ELEMENTARY STREAM DESCRIPTOR	Version:0 Flags:0 Tag:Es Len:27 TrackId:1 Priority:0 Tag:Cfg Len:13 Obj:PNG Stream:Video Buffer:0 Avg:16,8Kps Max:16,8Kps Tag:S1 Len:1 Marker		
D6	Size 18	Type stts	TIME TO SAMPLE	Version:0 Flags:0 Entries:1 SampleCount:1 SampleDuration:200		
EE	Size 1C	Type stsc	SAMPLE TO CHUNK ENTRIES	Version:0 Flags:0 Size:1 Entries:1 FirstChunk:1 SamplesPerChunk:1 SampleID:1		
10A	Size 14	Type stsz	SAMPLE SIZE	Version:0 Flags:0 Size:54 Entries:1		
11E	Size 14	Type stco	CHUNK OFFSET	Version:0 Flags:0 Entries:1 Entry:13A		
132	Size 5C	Type mdat	MEDIA DATA			
13A	PNG image		Signature	89 PNG\r\n\r\n		
			Header	0000000D IHDR 00 00 00 0D 00 00 00 07 01 00 00 00 00 FB0BFAB7		
			Data	0000001B IDAT 78 01 63 F8 FF 83 61 E5 0A 86 AE 15 0C 9D 1D 0C AB 5F 00 11 50 04 00 6D 06 0A A9 112AD7F7		
			Img End	00000000 IEND AE426082		

ANGE ALBERTINI 2022 (CC BY 4.0)

This "atom/box" structure is the ISO base media file format (also used in MOV, JP2, HEIF...)




Don't write parsers by hand, use a suitable binary tool.

There are many to choose from
and they will save you a lot of time.




Thanks!

 <https://kaitai.io/>

 <https://github.com/kaitai-io>

 https://gitter.im/kaitai_struct/Lobby

 [@kaitai_io](https://twitter.com/kaitai_io)