# Binary Data in Ruby

David Grayson
Las Vegas Ruby Meetup
2014-05-07

"\xf9\xbe\xb4\xd9\x1d\x01\x00\x00\x01\x00\x00\x00\x00\x00\x00\x00\x00\x00
0\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00
0\x00\x00\x00\x00\x00\x00\x3b\xa3\xed\xfd\x7a\x7b\x12\xb2\x7a\xc7\x2c\x3e
7\x76\x8f\x61\x7f\xc8\x1b\xc3\x88\x8a\x51\x32\x3a\x9f\xb8\xaa\x4b\x1e\x5e
a\x29\xab\x5f\x49\xff\xff\x00\x1d\x1d\xac\x2b\x7c\x01\x01\x00\x00\x00\x01
0\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00
0\x00\x00\x00\x0                                        \xff\xff\x4d\x04
f\xff\x00\x1d\x0                                        \x73\x20\x30\x33
f\x4a\x61\x6e\x2                                        \x65\x6c\x6c\x6f
2\x20\x6f\x6e\x2                                        \x65\x63\x6f\x6e
4\x20\x62\x61\x69\x6c\x6f\x75\x74\x20\x66\x6f\x72\x20\x62\x61\x6e\x6b\x73
f\xff\xff\xff\x01\x00\xf2\x05\x2a\x01\x00\x00\x00\x43\x41\x04\x67\x8a\xfd
0\xfe\x55\x48\x27\x19\x67\xf1\xa6\x71\x30\xb7\x10\x5c\xd6\xa8\x28\xe0\x39
9\xa6\x79\x62\xe0\xea\x1f\x61\xde\xb6\x49\xf6\xbc\x3f\x4c\xef\x38\xc4\xf3
5\x04\xe5\x1e\xc1\x                                      57\x8a\x4c\x70\x2b
b\xf1\x1d\x5f\xac\x                                      00\x00\x00\x01\x00
0\x00\x6f\xe2\x8c\x                                      ae\x63\xf7\x4f\x93
e\x83\x65\xe1\x5a\x                                      00\x98\x20\x51\xfd
e\x4b\xa7\x44\xbb\x                                      a3\xc3\x54\x0b\xf7
1\xcd\xb6\x06\xe8\x                                      ff\x00\x1d\x01\xe3
2\x99\x01\x01\x00\x00\x00\x01\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00
0\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00
0\x00\xff\xff\xff\xff\x07\x04\xff\xff\x00\x1d\x01\x04\xff\xff\xff\xff\x01
0\xf2\x05\x2a\x01\x00\x00\x00\x43\x41\x04\x96\xb5\x38\xe8\x53\x51\x9c\x72
a\x2c\x91\xe6\x1e\xc1\x16\x00\xae\x13\x90\x81\x3a\x62\x7c\x66\xfb\x8b\xe7
4\x7b\xe6\x3c\x52\xda\x75\x89\x37\x95\x15\xd4\xe0\xa6\x04\xf8\x14\x17\x81

# My applications of binary data

- Brawlsnapshots.com (2008)
  - Extract metadata from user-uploaded stage file
- redstone-bot2 (2012)
  - Minecraft bot written in Ruby
- Ruby ECDSA gem (2014)
  - Supports standard binary data formats

# Outline

- Quick introduction to binary data
  - Bytes
  - ASCII
- Binary data in Ruby
  - Strings
  - Getting bytes and integers from a binary string
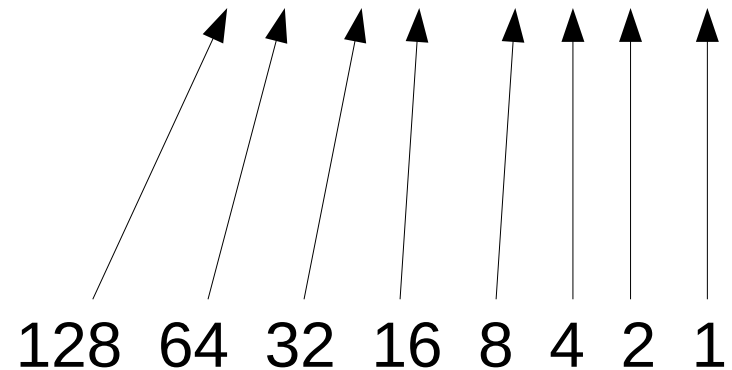  - Bit fields
  - IO objects

This talk focuses on reading information from binary data rather than writing binary data, but Ruby supports both!

# Definition: Byte

1. A number between 0 and 255

2. A storage location that can hold such a number

# A byte has 8 bits

```
Bits of byte 0:    0000 0000
Bits of byte 1:    0000 0001
Bits of byte 2:    0000 0010
Bits of byte 132:  1000 0100
```

128 64 32 16 8 4 2 1

# List of all bytes in decimal

| | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 16 | 32 | 48 | 64 | 80 | 96 | 112 | 128 | 144 | 160 | 176 | 192 | 208 | 224 | 240 |
| 1 | 17 | 33 | 49 | 65 | 81 | 97 | 113 | 129 | 145 | 161 | 177 | 193 | 209 | 225 | 241 |
| 2 | 18 | 34 | 50 | 66 | 82 | 98 | 114 | 130 | 146 | 162 | 178 | 194 | 210 | 226 | 242 |
| 3 | 19 | 35 | 51 | 67 | 83 | 99 | 115 | 131 | 147 | 163 | 179 | 195 | 211 | 227 | 243 |
| 4 | 20 | 36 | 52 | 68 | 84 | 100 | 116 | 132 | 148 | 164 | 180 | 196 | 212 | 228 | 244 |
| 5 | 21 | 37 | 53 | 69 | 85 | 101 | 117 | 133 | 149 | 165 | 181 | 197 | 213 | 229 | 245 |
| 6 | 22 | 38 | 54 | 70 | 86 | 102 | 118 | 134 | 150 | 166 | 182 | 198 | 214 | 230 | 246 |
| 7 | 23 | 39 | 55 | 71 | 87 | 103 | 119 | 135 | 151 | 167 | 183 | 199 | 215 | 231 | 247 |
| 8 | 24 | 40 | 56 | 72 | 88 | 104 | 120 | 136 | 152 | 168 | 184 | 200 | 216 | 232 | 248 |
| 9 | 25 | 41 | 57 | 73 | 89 | 105 | 121 | 137 | 153 | 169 | 185 | 201 | 217 | 233 | 249 |
| 10 | 26 | 42 | 58 | 74 | 90 | 106 | 122 | 138 | 154 | 170 | 186 | 202 | 218 | 234 | 250 |
| 11 | 27 | 43 | 59 | 75 | 91 | 107 | 123 | 139 | 155 | 171 | 187 | 203 | 219 | 235 | 251 |
| 12 | 28 | 44 | 60 | 76 | 92 | 108 | 124 | 140 | 156 | 172 | 188 | 204 | 220 | 236 | 252 |
| 13 | 29 | 45 | 61 | 77 | 93 | 109 | 125 | 141 | 157 | 173 | 189 | 205 | 221 | 237 | 253 |
| 14 | 30 | 46 | 62 | 78 | 94 | 110 | 126 | 142 | 158 | 174 | 190 | 206 | 222 | 238 | 254 |
| 15 | 31 | 47 | 63 | 79 | 95 | 111 | 127 | 143 | 159 | 175 | 191 | 207 | 223 | 239 | 255 |

# List of all bytes in hexadecimal

```
00  10  20  30  40  50  60  70  80  90  a0  b0  c0  d0  e0  f0
01  11  21  31  41  51  61  71  81  91  a1  b1  c1  d1  e1  f1
02  12  22  32  42  52  62  72  82  92  a2  b2  c2  d2  e2  f2
03  13  23  33  43  53  63  73  83  93  a3  b3  c3  d3  e3  f3
04  14  24  34  44  54  64  74  84  94  a4  b4  c4  d4  e4  f4
05  15  25  35  45  55  65  75  85  95  a5  b5  c5  d5  e5  f5
06  16  26  36  46  56  66  76  86  96  a6  b6  c6  d6  e6  f6
07  17  27  37  47  57  67  77  87  97  a7  b7  c7  d7  e7  f7
08  18  28  38  48  58  68  78  88  98  a8  b8  c8  d8  e8  f8
09  19  29  39  49  59  69  79  89  99  a9  b9  c9  d9  e9  f9
0a  1a  2a  3a  4a  5a  6a  7a  8a  9a  aa  ba  ca  da  ea  fa
0b  1b  2b  3b  4b  5b  6b  7b  8b  9b  ab  bb  cb  db  eb  fb
0c  1c  2c  3c  4c  5c  6c  7c  8c  9c  ac  bc  cc  dc  ec  fc
0d  1d  2d  3d  4d  5d  6d  7d  8d  9d  ad  bd  cd  dd  ed  fd
0e  1e  2e  3e  4e  5e  6e  7e  8e  9e  ae  be  ce  de  ee  fe
0f  1f  2f  3f  4f  5f  6f  7f  8f  9f  af  bf  cf  df  ef  ff
```

# ASCII code

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 00 | 10 | 20 | 30 0 | 40 @ | 50 P | 60 ` | 70 p |
| 01 | 11 | 21 ! | 31 1 | 41 A | 51 Q | 61 a | 71 q |
| 02 | 12 | 22 \ | 32 2 | 42 B | 52 R | 62 b | 72 r |
| 03 | 13 | 23 # | 33 3 | 43 C | 53 S | 63 c | 73 s |
| 04 | 14 | 24 $ | 34 4 | 44 D | 54 T | 64 d | 74 t |
| 05 | 15 | 25 % | 35 5 | 45 E | 55 U | 65 e | 75 u |
| 06 | 16 | 26 & | 36 6 | 46 F | 56 V | 66 f | 76 v |
| 07 \a | 17 | 27 ' | 37 7 | 47 G | 57 W | 67 g | 77 w |
| 08 \b | 18 | 28 ( | 38 8 | 48 H | 58 X | 68 h | 78 x |
| 09 \t | 19 | 29 ) | 39 9 | 49 I | 59 Y | 69 i | 79 y |
| 0a \n | 1a | 2a * | 3a : | 4a J | 5a Z | 6a j | 7a z |
| 0b \v | 1b \e | 2b + | 3b ; | 4b K | 5b [ | 6b k | 7b { |
| 0c \f | 1c | 2c , | 3c < | 4c L | 5c \ | 6c l | 7c | |
| 0d \r | 1d | 2d - | 3d = | 4d M | 5d ] | 6d m | 7d } |
| 0e | 1e | 2e . | 3e > | 4e N | 5e ^ | 6e n | 7e ~ |
| 0f | 1f | 2f / | 3f ? | 4f O | 5f _ | 6f o | 7f |

Not shown: ASCII characters without escape sequences in Ruby

# Text files generally use ASCII

Text editor
display:

LVRUG!

Actual bytes in
text file:

0x4c, 0x56, 0x52, 0x55, 0x47, 0x21

But there are many other encodings, like UTF-8, UTF-16.

# Definition: text

- A sequence of characters

# Definition: binary data

- A sequence of bytes that is not text

# General properties of binary data

- It has bytes that are not ASCII characters, like 0x80–0xFF  (128 to 255)

- It stores data more compactly than text.

- It looks like junk in a text editor.

# Text editors

# Hex editors

# Recap of binary data

- A byte is a number between 0 and 255.

- A byte can be written as 2 hex digits.

- A byte can be written as 8 bits.

- ASCII is a popular mapping between *bytes* and *characters*.

# Binary data in Ruby

- Typically stored in a String

```
str = "\x12\x34\xFE"
```

# Fix your string's encoding!

```ruby
"\x12\x34\xfe".force_encoding('BINARY')


# coding: ASCII-8BIT
"\x12\x34\xfe"        # ← first line of file


def some_public_method(str)
  str = str.dup.force_encoding('BINARY')
  # ...
end
```

# Integer literals in Ruby

- Three ways to write the same number:

```
181        # decimal
0xb5       # hex
0b10110101 # binary
```

# Inspecting a string's bytes

```ruby
str = "\x0d\x0e\x00\x40"

str.bytes.to_a    # => [13, 14, 0, 64]

str.inspect       # => "\r\x0E\x00@"

str.hex_inspect   # => "\x0d\x0e\x00\x40"
```

```ruby
class String
  def hex_inspect
    '"' + each_byte.map { |b| '\x%02x' % b }.join + '"'
  end
end
```

https://gist.github.com/DavidEGrayson/10093790

# Getting a single byte from a string

```
str = "\x0d\x0e\x00\x40"

byte1 = str[1].ord   # => 14
```

# Unpacking integers from a string

- Main tool: String#unpack

- Big endian vs. little endian

```ruby
# two 8-bit unsigned integers
"\x34\x12".unpack('cc')  # => [0x34, 0x12]


# one 16-bit unsigned int, little-endian
"\x34\x12".unpack('S')   # => [0x1234]


# one 16-bit unsigned int, big-endian
"\x34\x12".unpack('n')   # => [0x3412]
```

http://ruby-doc.org/core/String.html#method-i-unpack

# Converting integers to hex

```
'm = %02x' % [14]     # => "m = 0e"
```

# Bit fields

- Bits and groups of bits inside a byte can have individual meaning

**Table 9-2. Format of Setup Data**

| Offset | Field | Size | Value | Description |
|--------|-------|------|-------|-------------|
| 0 | *bmRequestType* | 1 | Bitmap | Characteristics of request:<br><br>D7: Data transfer direction<br>0 = Host-to-device<br>1 = Device-to-host<br><br>D6...5: Type<br>0 = Standard<br>1 = Class<br>2 = Vendor<br>3 = Reserved<br><br>D4...0: Recipient<br>0 = Device<br>1 = Interface<br>2 = Endpoint<br>3 = Other<br>4...31 = Reserved |

# Getting bit fields from an integer

- Binary operators:
  - <<    bitwise shift left
  - \>>    bitwise shift right
  - &    bitwise and
  - |    bitwise or

- Example: extract a 2-bit bitfield in bits 5-6

```
x = 0b11000000      # same as 0xC0
(x >> 5) & 0b11     # => 0b10
```

# Read binary data from a file

```ruby
f = File.open('foo.dat', 'rb')

f.read(2)  # => 2-byte string
f.read(10) # => 10-byte string

f.close
```

More methods you can use on files & sockets: http://www.ruby-doc.org/core/IO.html

# Converting a String to an IO object

```ruby
require 'stringio'
io = StringIO.new "\x0d\x0e\x00\x40"

io.read(2)  # => "\x0d\x0e"
io.read(1)  # => "\x00"
io.read(1)  # => "\x40"
io.read(1)  # => nil
```
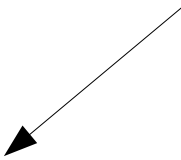
# Reading variable-length binary data

```ruby
"\x02\x00\x12\x34\xAA\xBB"

length = io.read(2).unpack('S')[0]
array = length.times.map do
  io.read(2).unpack('S')[0]
end
```

Extending or refining the IO class allows clearer code!

```ruby
length = io.read_uint16
array = length.times.map do
  io.read_uint16
end
```

# Example
# IO extension

```ruby
module DataReader
  def read_uint8
    read(1).ord
  end

  def read_int8
    read(1).unpack('c')[0]
  end

  def read_int16
    read(2).unpack('s')[0]
  end

  def read_uint16
    read(2).unpack('S')[0]
  end

  def read_uint16_array
    length = read_uint16
    length.times.map { read_uint16 }
  end
end

io_object.extend DataReader
```

# Example: reading Minecraft Entity Properties packet (simplified)

```ruby
def receive_data(stream)
  @eid = stream.read_int
  property_count = stream.read_int
  @properties = property_count.times.map do
    key = stream.read_string
    value = stream.read_double
    [key, value]
  end
end
```

# Creating binary data

- Array#pack

- String#concat (or just adding strings)

- IO#write

- Be careful to open files in binary mode:

  - `File.open(name, 'wb')`

# End