

# A gentle introduction to the: Succinct Data Structure Library

Algorithm Engineering A.Y. 2021/22

Università di Pisa

Giorgio Vinciguerra

Teaching assistant

# Lab lectures

- **Not** part of the syllabus but useful for your future career
- Tools for working with big data: Compressed/succinct data structures
  - Complex algorithms become possible in memory, no clusters and no disks
  - The libraries we will see offer the same API of standard data structures

# From theory to practice

- Integer coding via Elias- $\gamma$  and  $-\delta$  codes (§11.1 of the notes)
- Plain bitvectors with rank/select support (§15.1.1)
- Compressed bitvectors via Elias-Fano coding (§11.6 and §15.1.2)
- Pointerless programming (§15.1)

# The Succinct Data Structure Library (SDSL)

<https://github.com/simongog/sdsl-lite>

(see also the fork <https://github.com/xxsds/sdsl-lite>)

- An **easy-to-use, highly-efficient, configurable**, and **extensible** library of succinct data structures for researchers and practitioners
- Implements highlights of 40 research articles
- Faithful to the original proposals while using modern instruction sets
- Written in C++11, API familiar to C++ STL users
- Wrappers in other languages (e.g. Python) available online

# SDSL resources

Cheatsheet: <http://simongog.github.io/assets/data/sdsl-cheatsheet.pdf>

API Docs: <http://algo2.iti.kit.edu/gog/docs/html/index.html>

Examples: <https://github.com/simongog/sdsl-lite/examples>



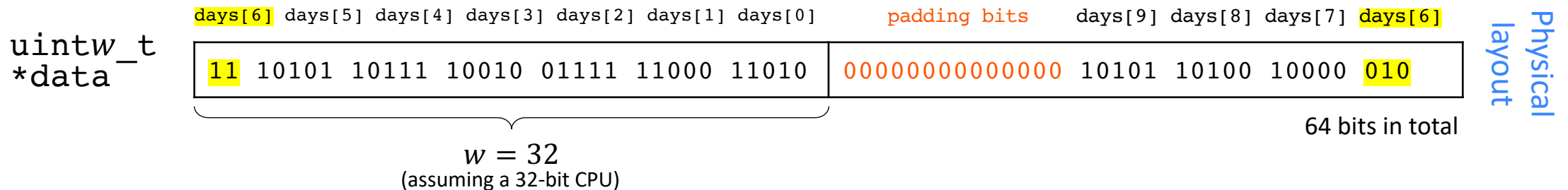
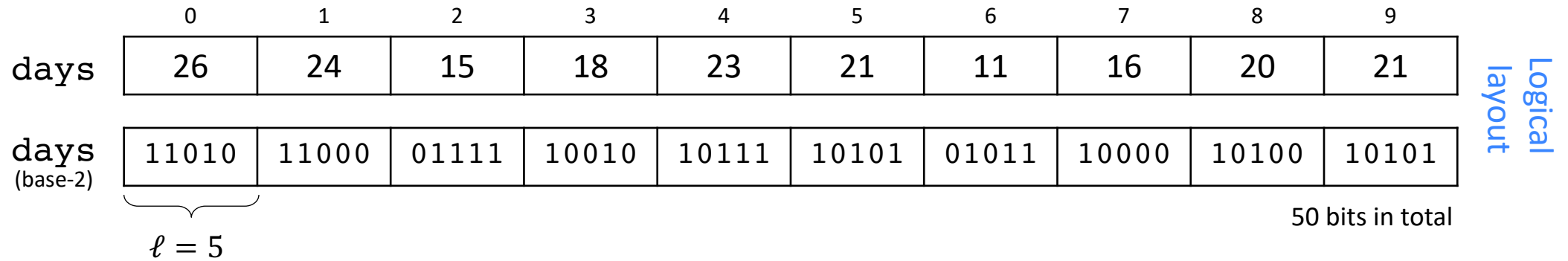
# The core of the library: `int_vector`

- Say you need to store the number of days per month worked by each employee of a company
- The `<cstdint>` header defines the (portable) types `intn_t` and `uintn_t` for  $n = 8, 16, 32, 64$
- So you would normally declare something like

```
std::vector<uint8_t> days(n_employees);
```
- What is the space? 1 byte per employee
- But only 5 bits are needed ( $31 = 11111_2$ )
- We are wasting  $8/5 - 1 = 60\%$  of extra space per employee!

# The core of the library: `int_vector`

- `int_vector<ℓ>` is a container of integers of fixed bit-width  $\ell$



$$j = i * \ell$$

$$\text{days}[i] = (\text{data}[j/w] \gg (j \% w)) \& \overbrace{((1 \ll \ell) - 1)}^{= 1^\ell}$$

Exercise: what is the formula when `days[i]` spans two words?



# Solution to the exercise

```
j = i * l
```

```
offset = j % w
```

```
if (offset + l <= w) // the integer spans one word
    return (data[j/w] >> (j % w)) & ((1 << l) - 1)
else // the integer spans two words
    return (data[j/w] >> (j % w)) |
           (data[j/w+1] & ((1 << ((offset+l) % w)) - 1)) << (w-offset)
```

## Example 1: `int_vector<ℓ>`, bit-width $\ell$ set at compile-time

```
#include <iostream>
#include <sdsl/vectors.hpp>

int main() {
    sds1::int_vector<5> v = {26, 24, 15, 18, 23, 21, 11, 16, 20, 21};
    std::cout << v << std::endl;
    v[3] = 10;
    v[1] = 194; // == 0b11000010
    std::cout << v << std::endl;
    std::cout << v.bit_size() << std::endl;
    std::cout << sds1::size_in_bytes(v) << std::endl;
    return 0;
}
```

26 24 15 18 23 21 11 16 20 21

26 2 15 10 23 21 11 16 20 21

50

16 ← 8 byte for the compressed v, 8 bytes for v.size()

## Example 2: `int_vector<>`, bit-width set at run-time

```
#include <iostream>
#include <sdsl/vectors.hpp>

int main() {
    sds1::int_vector<> v = {26, 24, 15, 18, 23, 21, 11, 16, 20, 21};
    std::cout << sds1::size_in_bytes(v) << std::endl;
    sds1::util::bit_compress(v);
    std::cout << v << std::endl;
    std::cout << (int) v.width() << std::endl;
    std::cout << sds1::size_in_bytes(v) << std::endl;
    return 0;
}
```

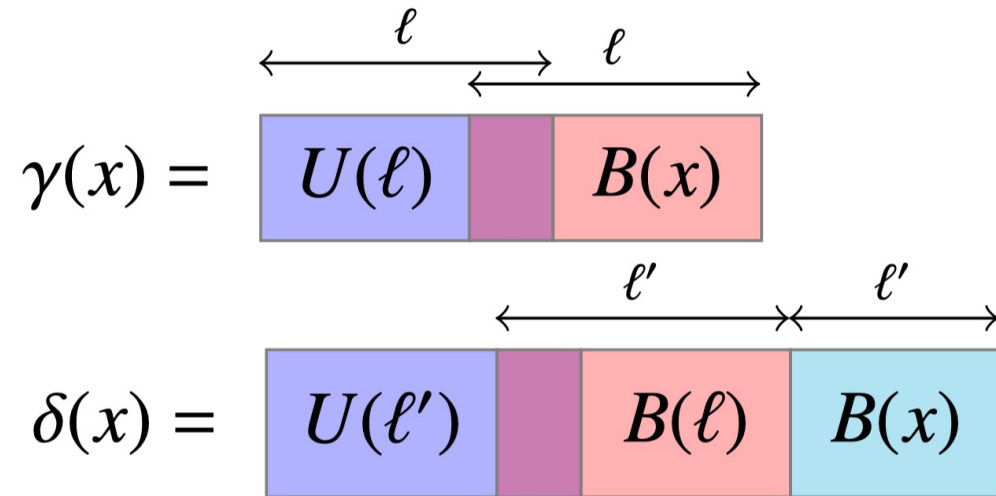
89 ← 80 bytes for the 10 ints, 8 bytes for `v.size()`, 1 byte for the bit-width

26 24 15 18 23 21 11 16 20 21

5

17 ← 8 byte for the compressed `v`, 8 bytes for `v.size()`, 1 byte for the bit-width

# Integer coders



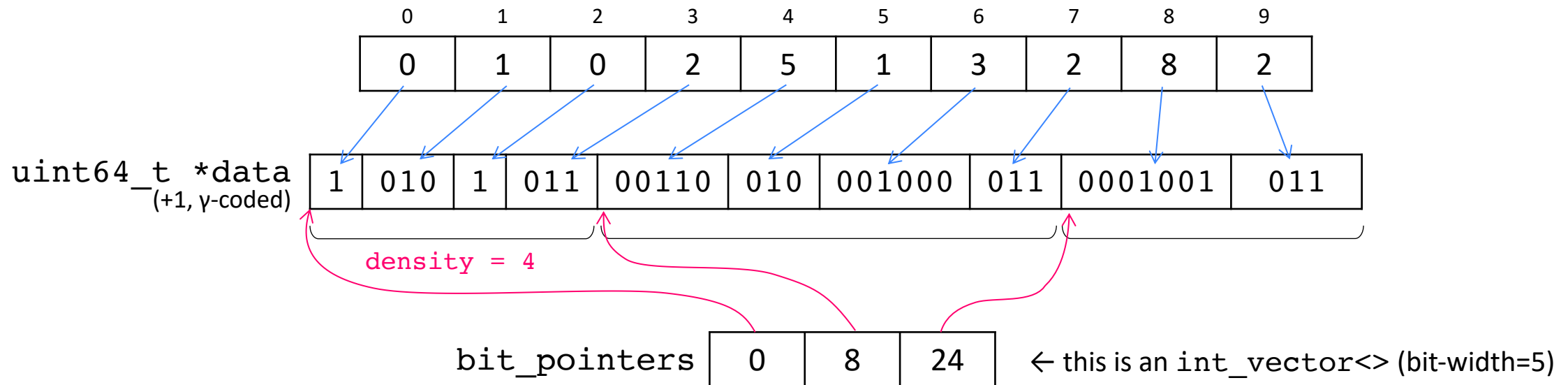
$$\gamma(9) = 0001001$$

$$\delta(14) = 0001001110$$

- `sdsl::coder::elias_gamma`
- `sdsl::coder::elias_delta`

# Compressed integer vectors: `vlc_vector`

- Stores Elias- $\gamma$  and - $\delta$  codes contiguously
- Zeros are permitted (internally, it uses  $x + 1$  instead of  $x$ )
- How to implement `vlc_vector::operator[]`?



`density` is a *space-time trade-off parameter*: decrease it for faster random access but higher space usage

## Example 3: `vlc_vector<coder, density>`

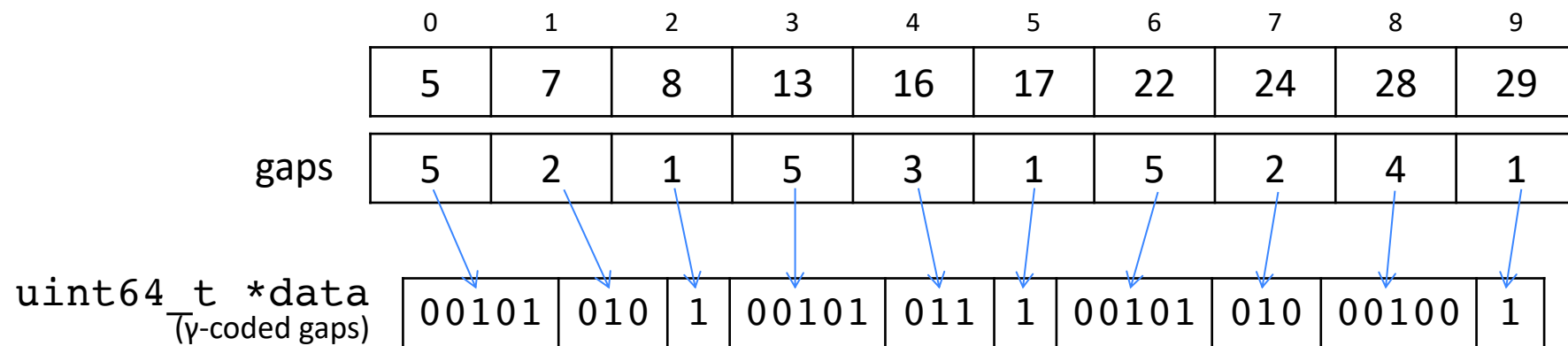
```
#include <iostream>
#include <sdsl/vectors.hpp>

int main() {
    sds1::int_vector<> v(10 * (1 << 20));
    v[100] = 1ULL << 63;
    sds1::util::bit_compress(v);
    std::cout << size_in_mega_bytes(v) << std::endl;
    sds1::vlc_vector<sds1::coder::elias_delta, 128> vlc(v);
    std::cout << size_in_mega_bytes(vlc) << std::endl;
    return 0;
}
```

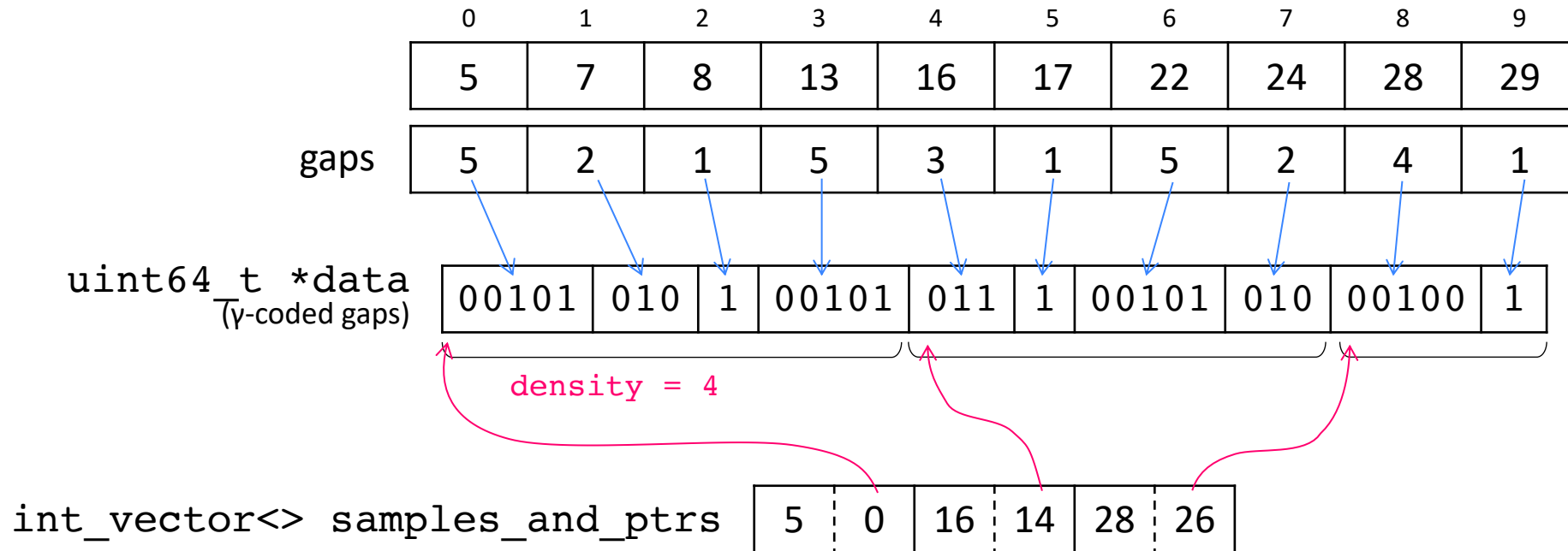
80  
1.48442

# Compressed integer vectors: `enc_vector`

- Say you have increasing integers  $x_1, x_2, x_3 \dots, x_n$   
A postings list (search engines), an adjacency list (graphs)
- $\gamma/\delta$ -code the gaps  $x_1, x_2 - x_1, x_3 - x_2, \dots, x_n - x_{n-1}$
- How to implement `vlc_vector::operator[]`?



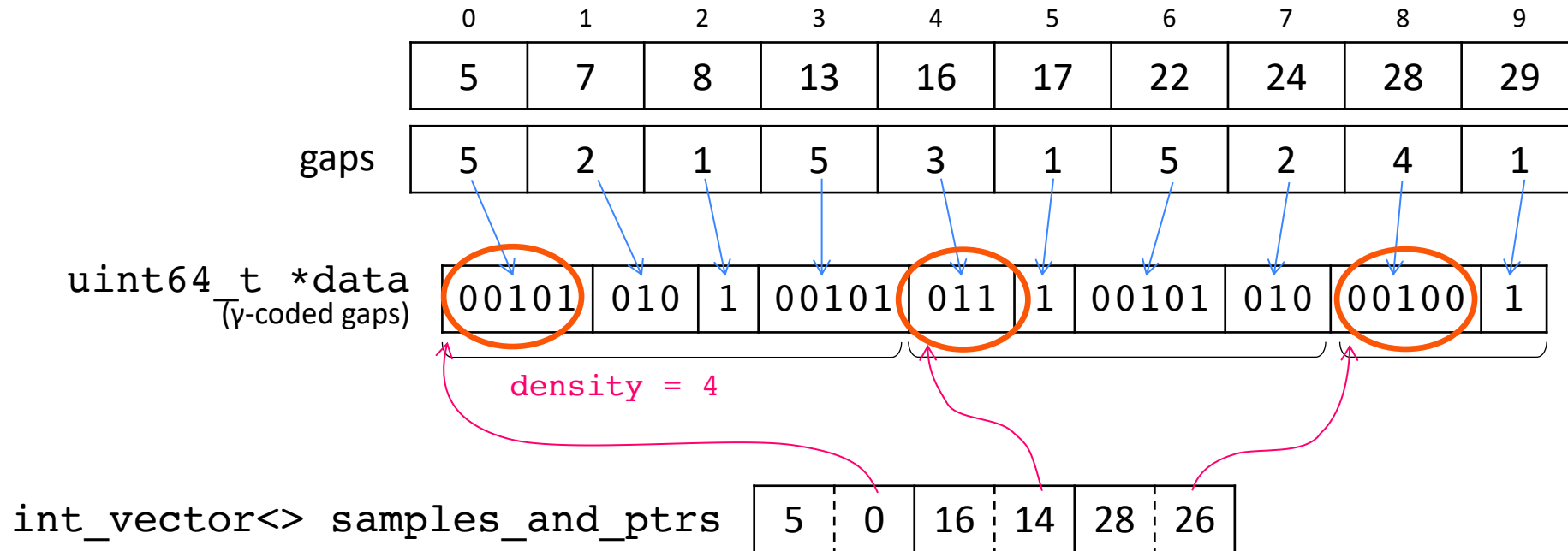
# Compressed integer vectors: `enc_vector`



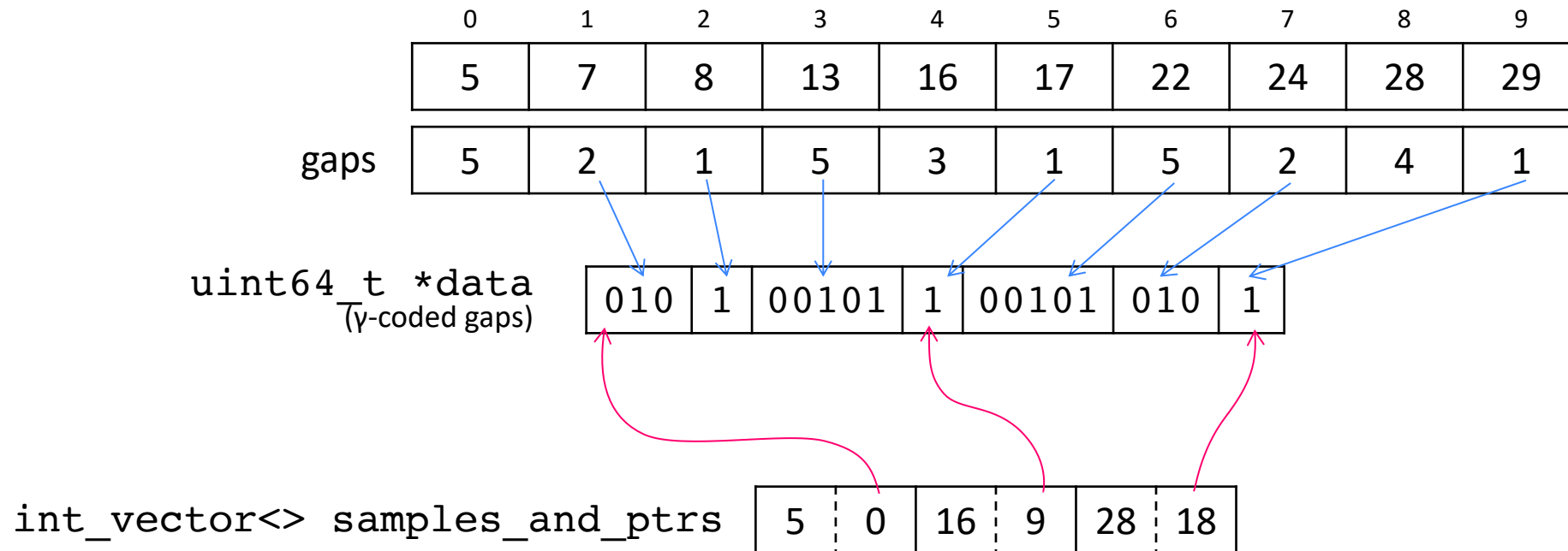


# Compressed integer vectors: `enc_vector`

- There is some **redundancy** here...



# Compressed integer vectors: `enc_vector`



**density** is a *space-time trade-off parameter*: decrease it for faster random access but higher space usage

# Example 4: `enc_vector<coder, density>`

```
#include <random>
#include <iostream>
#include <sdsl/vectors.hpp>

sdsl::int_vector<> random_data(size_t n, uint64_t u) { // Algorithm 3.4 of the notes
    if (n > u)
        throw std::invalid_argument("n > u");
    std::mt19937 rnd(std::random_device{}());
    sdsl::int_vector<> out(n);
    for (size_t j = 1, s = 0; j <= u && s < n; ++j)
        if (rnd() % (u - j + 1) < n - s)
            out[s++] = j;
    return out;
}

...
```

# Example 4: `enc_vector<coder, density>`

```
#include <random>
#include <iostream>
#include <sdsl/vectors.hpp>

sdsl::int_vector<> random_data(size_t n, uint64_t u) {...}

int main() {
    auto data = random_data(10 * (1 << 20), 1 << 24);
    std::cout << sdsl::size_in_mega_bytes(data) << std::endl;

    sdsl::vlc_vector<sdsl::coder::elias_delta, 128> vlc(data);
    std::cout << sdsl::size_in_mega_bytes(vlc) << std::endl;

    sdsl::enc_vector<sdsl::coder::elias_delta, 128> enc_delta(data);
    std::cout << sdsl::size_in_mega_bytes(enc_delta) << std::endl;

    return 0;
}
```

What if we use `sdsl::coder::elias_gamma`?

80  
39.0282  
3.19339

# Example 4(bis): `enc_vector<coder, density>`

```
#include <random>
#include <iostream>
#include <sdsl/vectors.hpp>

sdsl::int_vector<> random_data(size_t n, uint64_t u) {...}

int main() {
    auto data = random_data(10 * (1 << 20), 1 << 24);
    std::cout << sdsl::size_in_mega_bytes(data) << std::endl;

    sdsl::vlc_vector<sdsl::coder::elias_delta, 128> vlc(data);
    std::cout << sdsl::size_in_mega_bytes(vlc) << std::endl;

    sdsl::enc_vector<sdsl::coder::elias_delta, 128> enc_delta(data);
    std::cout << sdsl::size_in_mega_bytes(enc_delta) << std::endl;

    sdsl::enc_vector<sdsl::coder::elias_gamma, 128> enc_gamma(data);
    std::cout << sdsl::size_in_mega_bytes(enc_gamma) << std::endl;
    return 0;
}
```

```
80
39.0282
3.19339
2.79236
```

# Plain bitvectors (`bit_vector`)

- A specialised version of `int_vector<1>`
- Mutable
  - `b[i] = 1`
  - `b.flip()`
- Bitwise operations between bitvectors `b1 |= b2` (also `&=`, `^=`)
- Auxiliary data structures extend a (static) bitvector functionality

Class	+Bits	Time
<code>rank_support_v</code>	$0.25n$	$O(1)$
<code>rank_support_scan</code>	64	$O(n)$
<code>select_support_mcl</code>	$\leq 0.2n$	$O(1)$
<code>select_support_scan</code>	64	$O(n)$

# Example 5: bit\_vector

```
#include <iostream>
#include <sdsl/bit_vectors.hpp>

int main() {
    sdsl::bit_vector b1 = {1, 1, 0, 1, 0, 0, 1};
    sdsl::bit_vector b2 = {0, 0, 1, 1, 0, 1, 0};
    b1 |= b2;
    b1.flip();
    std::cout << b1 << std::endl;

    sdsl::bit_vector b(80 * (1 << 20), 0);
    for (size_t i = 0; i < b.size(); i += 100)
        b[i] = 1;
    std::cout << sdsl::size_in_mega_bytes(b) << std::endl;
}
```

0000100

10

# Recap

- **Integer vectors**

- `int_vector<>`, bit-width set at run-time, starts with 64 bits by default
- `int_vector< $\ell$ >`, bit-width  $\ell$  fixed at compile-time
- Implementation of random access

- **Compressed integer vectors**

- `vlc_vector<coder, density>`, vector of  $\gamma/\delta$ -coded integers
  - Implementation of random access via bit pointers
- `enc_vector<coder, density>`, vector of  $\gamma/\delta$ -coded gaps between ints
  - Implementation of random access via bit pointers and samples

- **Plain bitvectors**

- `bit_vector = int_vector<1>`
- Bitwise AND, OR, XOR
- Set individual bit, flip all bits



# Example 6: `bit_vector` with rank/select

```
#include <iostream>
#include <sdsl/bit_vectors.hpp>

int main() {
    sdsl::bit_vector b = {0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1};
    sdsl::rank_support_v<1> b_rank(&b);
    for (size_t i = 1; i <= b.size(); ++i)
        std::cout << b_rank(i) << " ";           counts #1s in b[0, i)
    std::cout << std::endl;

    sdsl::select_support_mcl<1> b_select(&b);
    size_t ones = b_rank(b.size());
    for (size_t i = 1; i <= ones; ++i)
        std::cout << b_select(i) << " ";
}
```

Output

```
b = 0 1 0 1 1 1 0 0 0 1 1
    0 1 1 2 3 4 4 4 4 5 6
    1 3 4 5 9 10
```

# Elias-Fano compressed bitvectors (`sd_vector`)

- Construct from a bitvector

```
sdsl::bit_vector b = {0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1};  
sdsl::sd_vector<> ef1(b);
```

- Construct from a vector of increasing integers

```
sdsl::int_vector<> v = {1, 3, 4, 5, 9, 10};  
sdsl::sd_vector<> ef2(v.begin(), v.end());
```

- `ef1 == ef2`

# Example 7: Elias-Fano

```
#include <iostream>
#include <sdsl/bit_vectors.hpp>

sdsl::int_vector<> random_data(size_t n, uint64_t u) {...}

int main() {
    auto data = random_data(10 * (1 << 20), 1 << 25);
    sdsl::sd_vector<> ef(data.begin(), data.end());
    std::cout << sdsl::size_in_mega_bytes(data) << std::endl;
    std::cout << sdsl::size_in_mega_bytes(ef) << std::endl;
    return 0;
}
```

80  
5.25673

# Example 8: Elias-Fano internals

```
#include <iostream>
#include <sdsl/bit_vectors.hpp>

int main() {
    sdsl::int_vector<> v = {1, 4, 7, 18, 24, 26, 30, 31};
    sdsl::sd_vector<> ef(v.begin(), v.end());
    uint64_t u = v[v.size() - 1] + 1;
    uint64_t bpi = std::ceil(std::log2(u));
    std::cout << "Data = " << v << std::endl
              << "Bitvector = " << ef << std::endl
              << "u = " << u << std::endl
              << "Bits per integer = " << bpi << std::endl
              << "Bits in the high part = " << bpi - ef.wl << std::endl
              << "Bits in the low part = " << (int) ef.wl << std::endl
              << "L = " << ef.low << " (in decimal)" << std::endl
              << "H = " << ef.high << std::endl;
    return 0;
}
```

Technical note: wrt the lecture notes, SDSL uses a different splitting point for the low/high parts (a difference of  $\pm 1$ ); but in this example the encodings match

```
Data = 1 4 7 18 24 26 30 31
Bitvector = 01001001000000000010000010100011
u = 32
Bits per integer = 5
Bits in the high part = 3
Bits in the low part = 2
L = 1 0 3 2 0 2 2 3 (in decimal)
H = 101100010011011000000000
```

1 =	000	01
4 =	001	00
7 =	001	11
18 =	100	10
24 =	110	00
26 =	110	10
30 =	111	10
31 =	111	11

$$h = \lceil \log n \rceil = 3$$
$$b = \lceil \log u \rceil = 5$$
$$\ell = b - h = 2$$

$L = 0100111000101011$

bucket 0 | 1 | 23 | 4 | 5 | 6 | 7  
 $H = 10111000100110110$

# Example 9: Elias-Fano access and query

```
#include <iostream>
#include <sdsl/bit_vectors.hpp>

int main() {
    sdsl::int_vector<> v = {1, 4, 7, 18, 24, 26, 30, 31};
    sdsl::sd_vector<> ef(v.begin(), v.end());
    std::cout << ef << std::endl;

    // Use as a bitvector with rank/select/operator[]
    sdsl::rank_support_sd<> ef_rank(&ef);
    sdsl::select_support_sd<> ef_select(&ef);
    std::cout << ef_rank(5) << std::endl;
    std::cout << ef_select(4) << std::endl;
    std::cout << ef[24] << std::endl; // access to the bit in pos 24

    // Use as an integer vector with access/nextGEQ given select/rank
    auto access = [&] (size_t i) { return ef_select(i + 1); };
    auto nextGEQ = [&] (uint64_t x) { return ef_select(ef_rank(x) + 1); };
    std::cout << access(5) << std::endl;
    std::cout << nextGEQ(4) << std::endl;
    std::cout << nextGEQ(27) << std::endl;
    return 0;
}
```

```
01001001000000000010000010100011
```

```
2
```

```
18
```

```
1
```

```
26
```

```
4
```

```
30
```