

**Indexing interlude:**  
**Bitvector Rank & Select:**  
**Primitives of succinct data  
structures**

# Thinking theoretically about data structure size

Assume that storing some data, in an information-theoretically optimal manner, requires  $Z$  bits

Representation of this data is:

*Implicit*:  $Z + O(1)$  bits

Only a constant size larger than the theoretical minimum

*Succinct*:  $Z + o(Z)$  bits

$Z$  bits, plus some term strictly smaller than  $Z$  bits

*Compact*:  $O(Z)$  bits

On the order of  $Z$  bits (grows linearly in  $Z$ )

# Thinking theoretically about data structure size

The idea of succinct data structures was first introduced by Jacobson in his thesis “Succinct static data structures”\*

In this thesis, among other things, he introduced succinct representations of trees and graphs that could be efficiently navigated.

As data sizes grow large, data structures that consume a lot of extra space become increasingly less feasible and so succinct data structures become increasingly important.

The *rank* and *select* operations become the basic building blocks of succinct data structures.

\*Jacobson, G. J (1988). Succinct static data structures (Ph.D.). Pittsburgh, PA: Carnegie Mellon University.

Slides for the following taken from:  
<https://www.cs.helsinki.fi/u/puglisi/dct2015/slides10.pdf>

**credit to Simon J. Puglisi, University of Helsinki**

# Succinct Data Structures

- Succinct data structure
  - = succinct representation of data + a succinct index
- (usually static)
- High-level goal: reduce space so the data structure might fit in RAM and therefore be faster to use
- Examples
  - Sets
  - Trees, graphs
  - Strings
  - Permutations, functions

# Succinct Representation

- A representation of data whose size (roughly) matches the information-theoretic lower bound
- If the input is taken from  $L$  distinct possible inputs, then its information-theoretic lower bound is  $\text{ceil}(\log L)$  bits
  - To be considered succinct a data structure must use:  
 $\text{ceil}(\log L) + o(\log L)$  bits
- Example: a lower bound for a set  $S$ , subset of  $\{1, 2, \dots, n\}$ 
  - $\log(2^n) = n$  bits
  - $n = 3$  we have 8 distinct sets... so d.s. will need at least 3 bits

$\emptyset$

$\{1\}$	$\{2\}$	$\{3\}$
$\{1, 2\}$	$\{1, 3\}$	$\{2, 3\}$
$\{1, 2, 3\}$		

# Succinct Index

- Auxiliary data structure to support queries on the succinct representation
- Size:  $o(\log L)$  bits
- The index should allow queries/operations on the succinct representation in (almost) the same time complexity as using a conventional data structure
  - This is the aim anyway
- Computational model is the word RAM
  - Assume word length  $w = \log \log L$
  - (this is the same pointer size as conventional data structures)
  - read/write  $w$  bits of memory in  $O(1)$  time
  - arithmetic/logical operations on  $w$  bit numbers take  $O(1)$  time
  - $+, -, *, /, \log, \&, |, !, >>, <<$

# Binary rank and select

- The ability to answer *rank* and *select* queries over bit vectors (binary strings, bit arrays) is essential for implementing succinct data structures
- Given a binary string  $B[1..n]$ 
  - $\text{rank}_B(i)$  returns the number of 1 bits in  $B[1..i]$
  - $\text{select}_B(i)$  returns the position of the  $i^{\text{th}}$  1 bit in  $B$



# Naïve rank

- To answer  $\text{rank}(i)$  scan  $B[1..i]$  and count 1-bits
- Simple but slow
  - $O(i)$  time =  $O(n)$  time in the worst case
- How can we do better?
  - After all, what are we?

## (Slightly) Less naïve rank

- Store an table  $A[1..n]$ , containing the rank answers

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
B	1	0	0	1	0	1	1	1	0	1	0	0	1	0	1	0
A	1	1	1	2	2	3	4	5	5	6	6	6	7	7	8	8

- $A[i] = \text{rank}(i)$ 
  - Now  $\text{rank}(i)$  takes constant time - just an array lookup!
- Drawback:
  - A requires  $n \log n$  bits -  $\log n$  times the size of B - not succinct!
  - We'd like a solution with  $O(1)$  queries and  $o(n)$  extra space...

# We want $O(1)$ queries with $o(n)$ extra bits...

- General approach will be to precompute some tables
- Each table stores part of the answer to every query
  - For any given query, we can extract needed parts in  $O(1)$  time
  - The total size of the tables is  $o(n)$  bits

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
B	1	0	0	1	0	1	1	1	0	1	0	0	1	0	1	0

- Premise:
  - Can read  $O(\log n)$  bits into an integer in range  $1..n$  in  $O(1)$  time
  - However, to inspect each of those bits take  $O(\log n)$  time

# Tables : Superblocks

- Divide B into superblocks of size  $s = \log^2 n / 2 = 4 * 4 / 2 = 8$
- Build a small table  $R_s$  containing ranks for only some positions

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
B	1	0	0	1	0	1	1	1	0	1	0	0	1	0	1	0
	0								1							
$R_s$	0								5							

- Store in  $R_s[j] = \text{rank}_B(j*s)$ , for all  $0 \leq j < n/s$

# Tables : Blocks

- Divide each superblock into blocks of size  $b = \log n / 2 = 2$
- Build a table  $R_b$  which contains the rank from the start of each block to the start of its superblock

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
B	1	0	0	1	0	1	1	1	0	1	0	0	1	0	1	0
	0								1							
R <sub>s</sub>	0								5							
	0	1	2	3	0	1	2	3								
R <sub>b</sub>	0	1	2	3	0	1	1	2								

- Store  $R_b[k/b] = \text{rank}_B(k*s) - \text{rank}_B(j*s)$ , for all  $0 \leq k < n/b$

# Intermission

- What we have so far (tables  $R_s$  and  $R_b$ ) almost gets us the answer we're after

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
B	1	0	0	1	0	1	1	1	0	1	0	0	1	0	1	0
	0								1							
R <sub>s</sub>	0								5							
	0	1	2	3	0	1	2	3								
R <sub>b</sub>	0	1	2	3	0	1	1	2								

- $\text{rank}_B(i) \approx R_s[i/s] + R_b[i/b]$ 
  - Just need to answer in-block queries in  $O(1)$  time

# Tables : Resolving in-block queries

- Solution? Use another table!
- Blocks have size  $b = \log_2 n / 2$ 
  - There are  $2^b$  such blocks possible
  - In each block there are  $b$  possible rank queries
  - Each answer (relative to the block) is in the range  $1..b$

$R_p$	Type	0	1	rank(0)	rank(1)
	0	0	0	0	0
	1	0	1	0	1
	2	1	0	1	1
	3	1	1	1	2

# Final Data Structure

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
B	1	0	0	1	0	1	1	1	0	1	0	0	1	0	1	0

	0								1							
R <sub>s</sub>	0								5							

	0	1	2	3	0	1	2	3
R <sub>b</sub>	0	1	2	3	0	1	1	2

	Type	0	1	rank(0)	rank(1)
R <sub>p</sub>	0	0	0	0	0
	1	0	1	0	1
	2	1	0	1	1
	3	1	1	1	2



# Size of table for within-block queries

- Blocks have size  $b = \log_2 n / 2$ 
  - There are  $2^b$  such blocks possible
  - In each block there are  $b$  possible rank queries
  - Each answer (relative to the block) is in the range  $1..b$

$R_p$	Type	0	1	rank(0)	rank(1)
	0	0	0	0	0
	1	0	1	0	1
	2	1	0	1	1
	3	1	1	1	2

- Therefore size of  $R_p$ , the in-block data structure is
  - $2^b * b * \log b = n^{1/2} * \log n * \log \log n / 2$  bits =  $o(n)$  bits

# Summing up sizes...

- The size of  $R_s$ , the superblock data structure is
  - $2n/\log^2 n$  superblocks, each of size  $\log n$  bits
  - $(n/\log^2 n) * \log n = 2n/\log n$  bits =  $o(n)$  bits
- The size of  $R_b$ , the block data structure is
  - $2n/\log n$  blocks, each of size  $\log \log n$  bits
  - $2n \log \log n / \log n$  bits =  $o(n)$  bits
- $R_s + R_b + R_p = o(n)$  extra bits for  $O(1)$  time rank queries
  - It is possible to construct this data structure in  $O(n)$  time

# Variations

- Just store  $R_s$  + use manual counting within superblocks
  - Saves space for  $R_b$  and  $R_p$ , takes time  $O(\log^2 n)$  per query
- Store  $R_s$  and  $R_b$  + use manual counting within blocks
  - Saves only space for  $R_p$ , takes time  $O(\log n)$  per query
- Use different superblock & block sizes
  - No more theoretical guarantees, but...
  - Perhaps faster in practice: blocks that are multiples of word sizes (32-bits) can be faster to handle

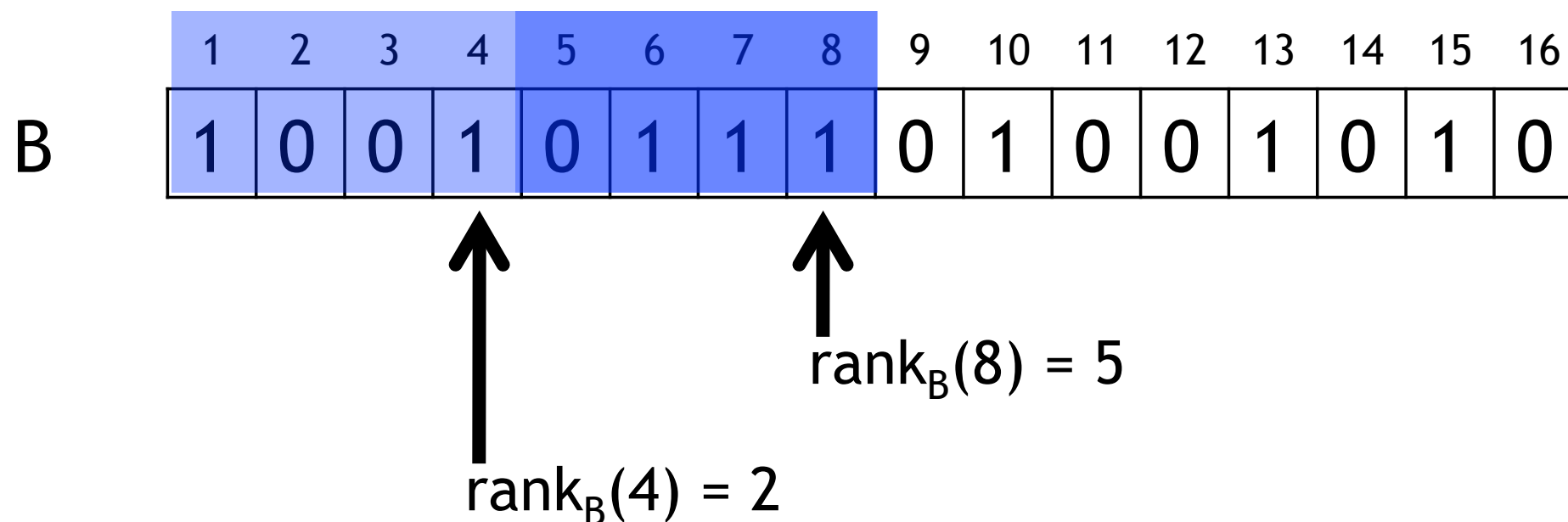
# Summary of rank

- Rank index takes  $O(n \log \log n / \log n) = o(n)$  bits so we use  $n + o(n)$  overall and can answer queries in  $O(1)$  time
- While it is sublinear, we'd still like the  $o(n)$  term to be small
  - Best is by Patrascu:  $O(n / \log^k n)$  bits,  $O(k)$  time queries
- Dynamic solutions exist
  - Queries no longer constant:  $O(\log n / \log \log n)$  time (Raman et al.)

# Relationship to select(i)

- We can use our solution to rank to get a (fairly) efficient solution to select(i), with this observation:
- If  $\text{rank}(n/2) > i$ , then the  $i^{\text{th}}$  1-bit is in  $B[1..n/2]$ 
  - Otherwise it is in  $B[n/2+1..n]$

$\text{select}_B(3)$



# Relationship to select(i)

- Applying this idea recursively to arrive at select(i)
  - $O(\log_2 n)$  time,  $o(n)$  space
- $O(1)$  time,  $o(n)$  space solutions for select also exist
  - Slightly more complicated than  $O(1)$  rank
  - (Munro and Clark)
- Similar variations as we discussed with rank (trading space for query time) are also possible

Information Systems 73 (2018) 25–34



Contents lists available at [ScienceDirect](#)

Information Systems

journal homepage: [www.elsevier.com/locate/is](http://www.elsevier.com/locate/is)



Rank and select: Another lesson learned

[Szymon Grabowski\\*](#), [Marcin Raniszewski](#)

Lodz University of Technology, Institute of Applied Computer Science, Al. Politechniki 11, Łódź 90–924, Poland



# These rank & select operations work over a binary alphabet — can be extended

## High-Order Entropy-Compressed Text Indexes

Roberto Grossi\*      Ankur Gupta†      Jeffrey Scott Vitter‡

Introduces the idea of the wavelet tree, a versatile index that can be extended to arbitrary alphabets. We'll discuss the simplest of variants according to the exposition of:

## Wavelet Trees for All \*

Gonzalo Navarro

Dept. of Computer Science, University of Chile. [gnavarro@dcc.uchile.cl](mailto:gnavarro@dcc.uchile.cl)

# Preliminaries

$S[1,n] = s_1 s_2 \dots s_n$  is a sequence of symbols where  $s_i$  in  $\Sigma$

$\Sigma = [1 \dots \sigma]$  is an alphabet of symbols

Representing  $S$  requires  $n * \lceil \lg \sigma \rceil = n * \lg \sigma + O(n)$  bits.

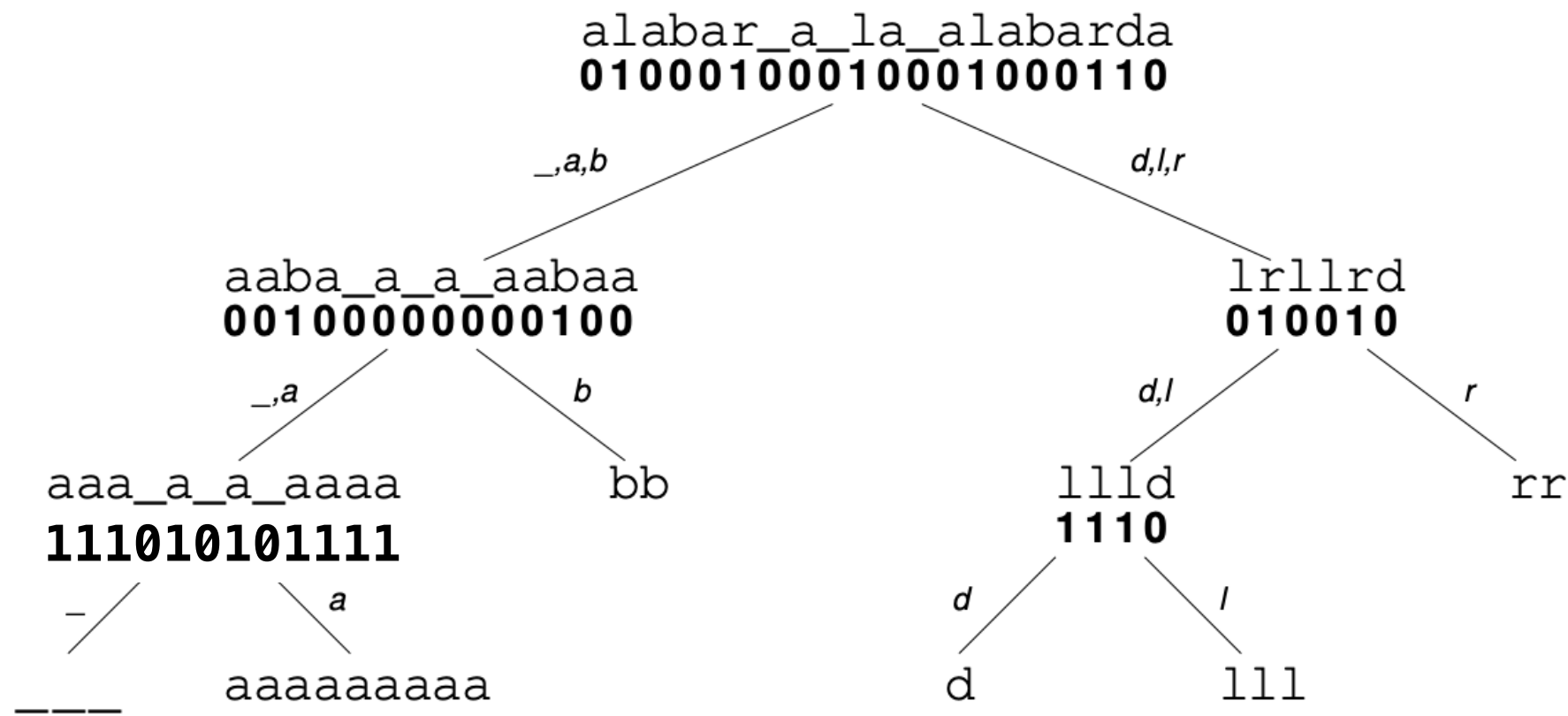
Wavelet tree: balanced binary tree with  $\sigma$  nodes, where each subtree is also a wavelet tree (i.e. it is recursive)



# Preliminaries

**Structure.** A wavelet tree [54] for sequence  $S[1, n]$  over alphabet  $[1..\sigma]$  can be described recursively, over a sub-alphabet range  $[a..b] \subseteq [1..\sigma]$ . A wavelet tree over alphabet  $[a..b]$  is a binary balanced tree with  $b - a + 1$  leaves. If  $a = b$ , the tree is just a leaf labeled  $a$ . Else it has an internal root node,  $v_{root}$ , that represents  $S[1, n]$ . This root stores a bitmap  $B_{v_{root}}[1, n]$  defined as follows: if  $S[i] \leq (a + b)/2$  then  $B_{v_{root}}[i] = 0$ , else  $B_{v_{root}}[i] = 1$ . We define  $S_0[1, n_0]$  as the subsequence of  $S[1, n]$  formed by the symbols  $c \leq (a + b)/2$ , and  $S_1[1, n_1]$  as the subsequence of  $S[1, n]$  formed by the symbols  $c > (a + b)/2$ . Then, the left child of  $v_{root}$  is a wavelet tree for  $S_0[1, n_0]$  over alphabet  $[a..\lfloor (a + b)/2 \rfloor]$  and the right child of  $v_{root}$  is a wavelet tree for  $S_1[1, n_1]$  over alphabet  $[1 + \lfloor (a + b)/2 \rfloor..b]$ .

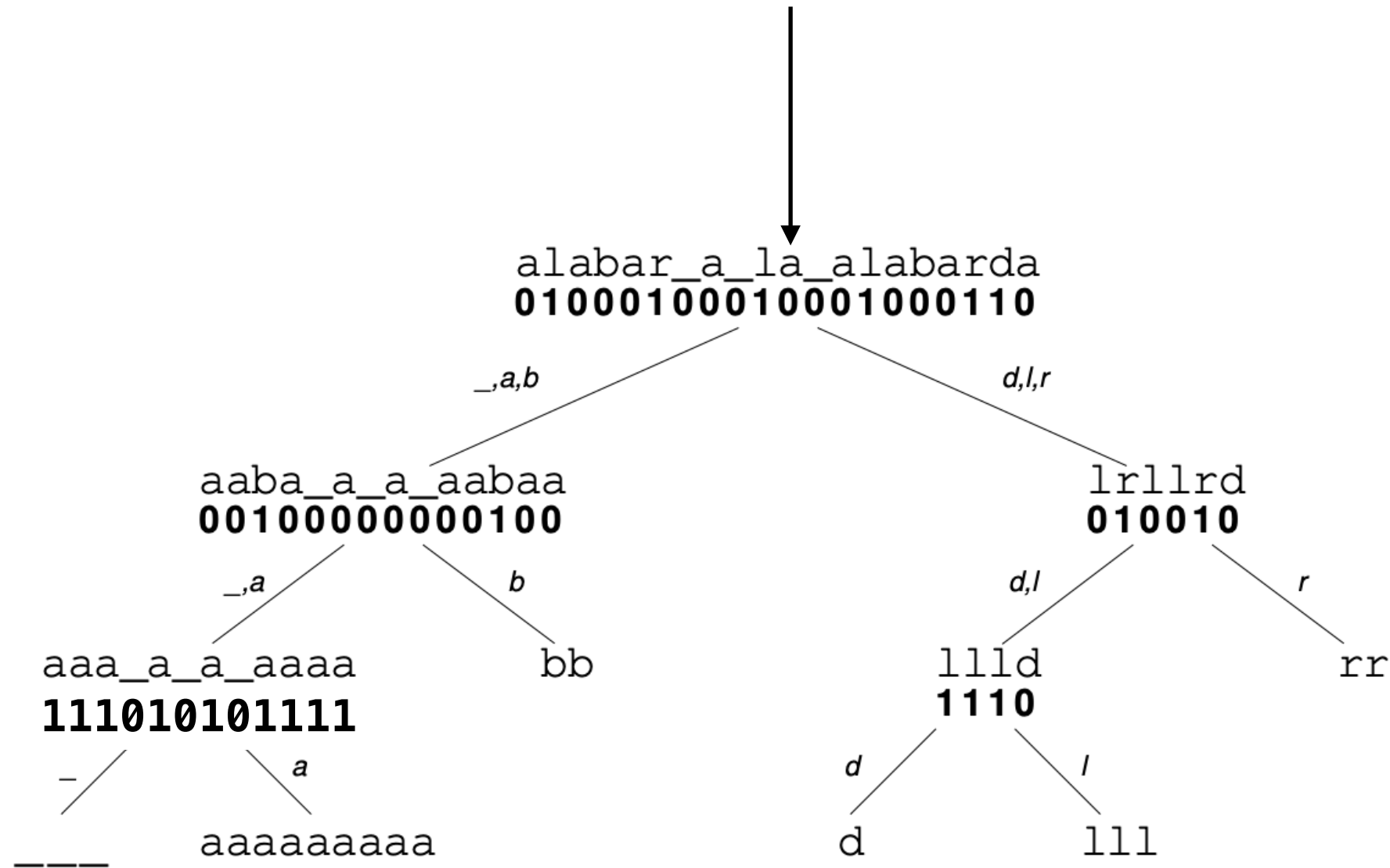
# Preliminaries



**Fig. 1.** A wavelet tree on string  $S = \text{"alabar a la alabarda"}$ . We draw the spaces as underscores. The subsequences of  $S$  and the subsets of  $\Sigma$  labeling the edges are drawn for illustration purposes; the tree stores only the topology and the bitmaps.

# Example: Rank

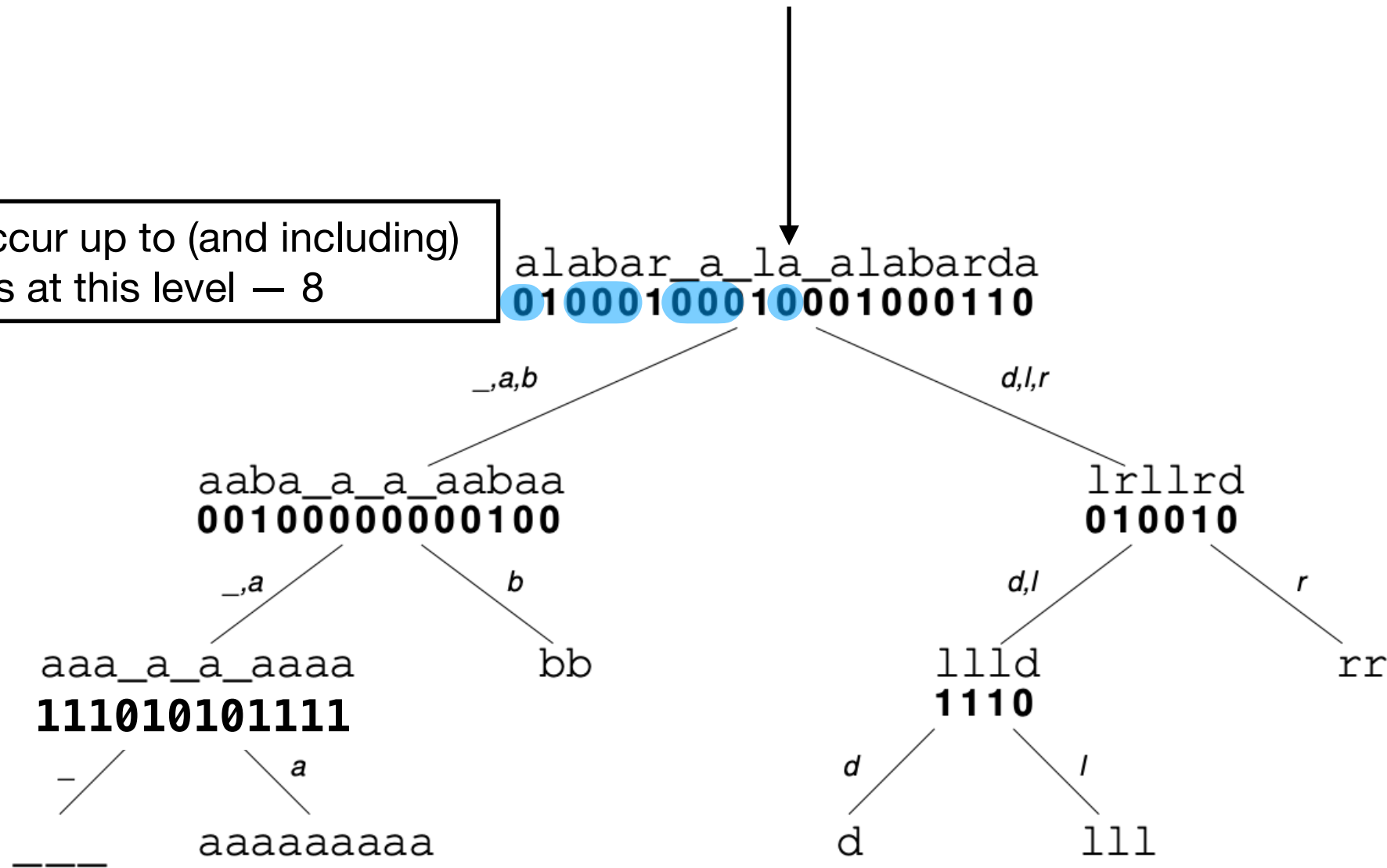
Consider asking for the rank of this “a”



# Example: Rank

Consider asking for the rank of this “a”

how many \_,a,b occur up to (and including) this one? Count 0's at this level — 8

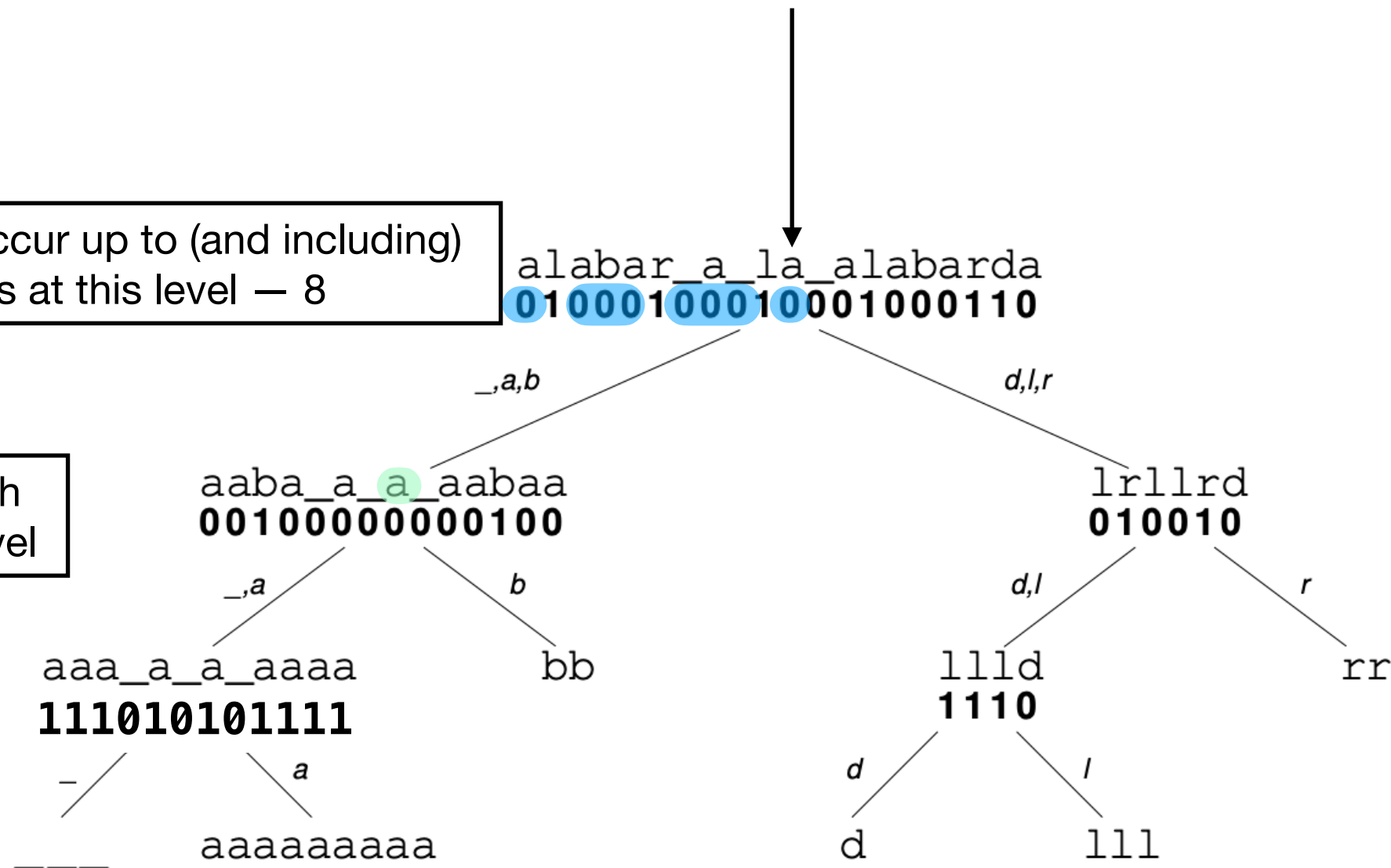


# Example: Rank

Consider asking for the rank of this “a”

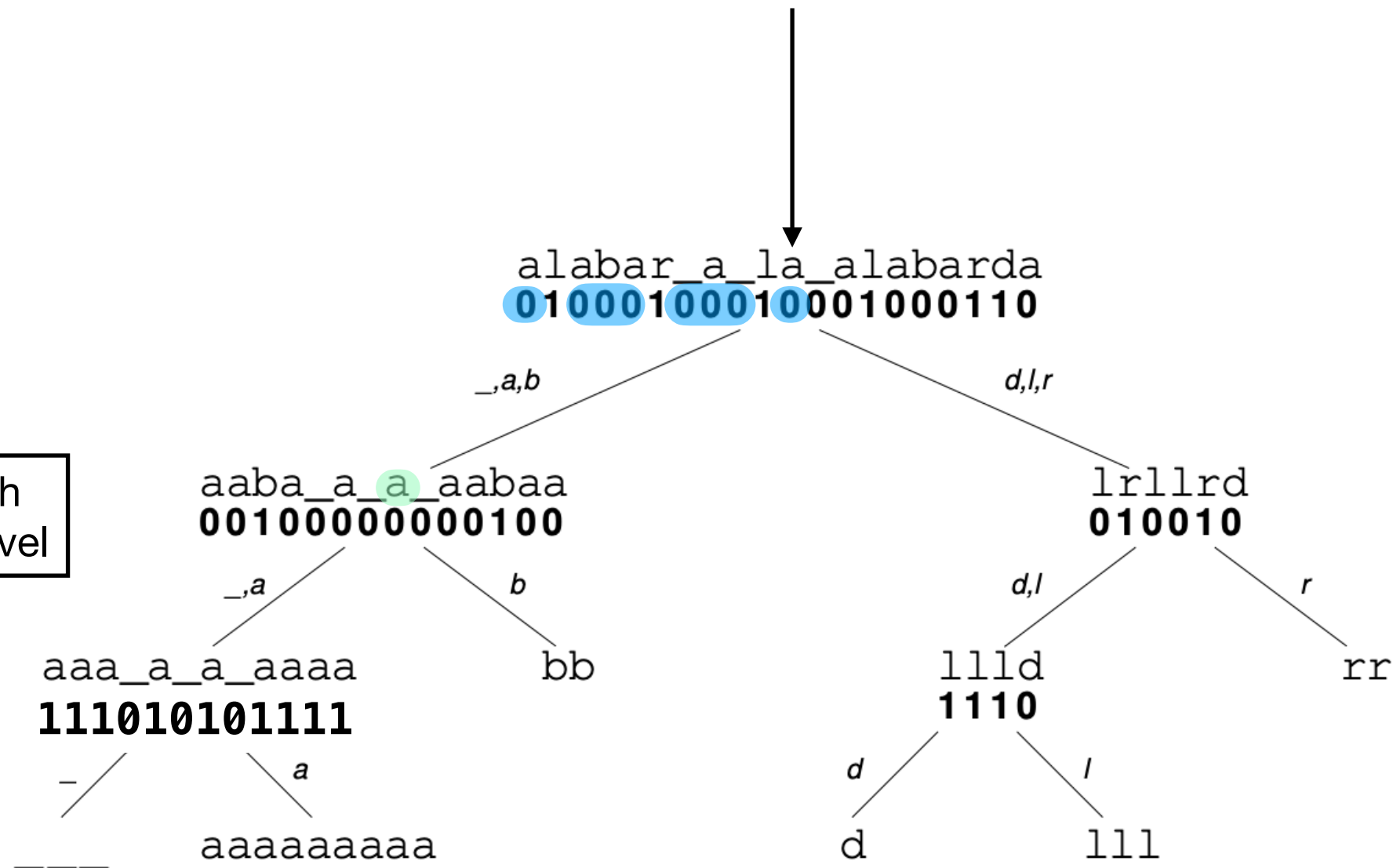
how many \_,a,b occur up to (and including) this one? Count 0's at this level — 8

So that maps to 8th character at this level



# Example: Rank

Consider asking for the rank of this “a”

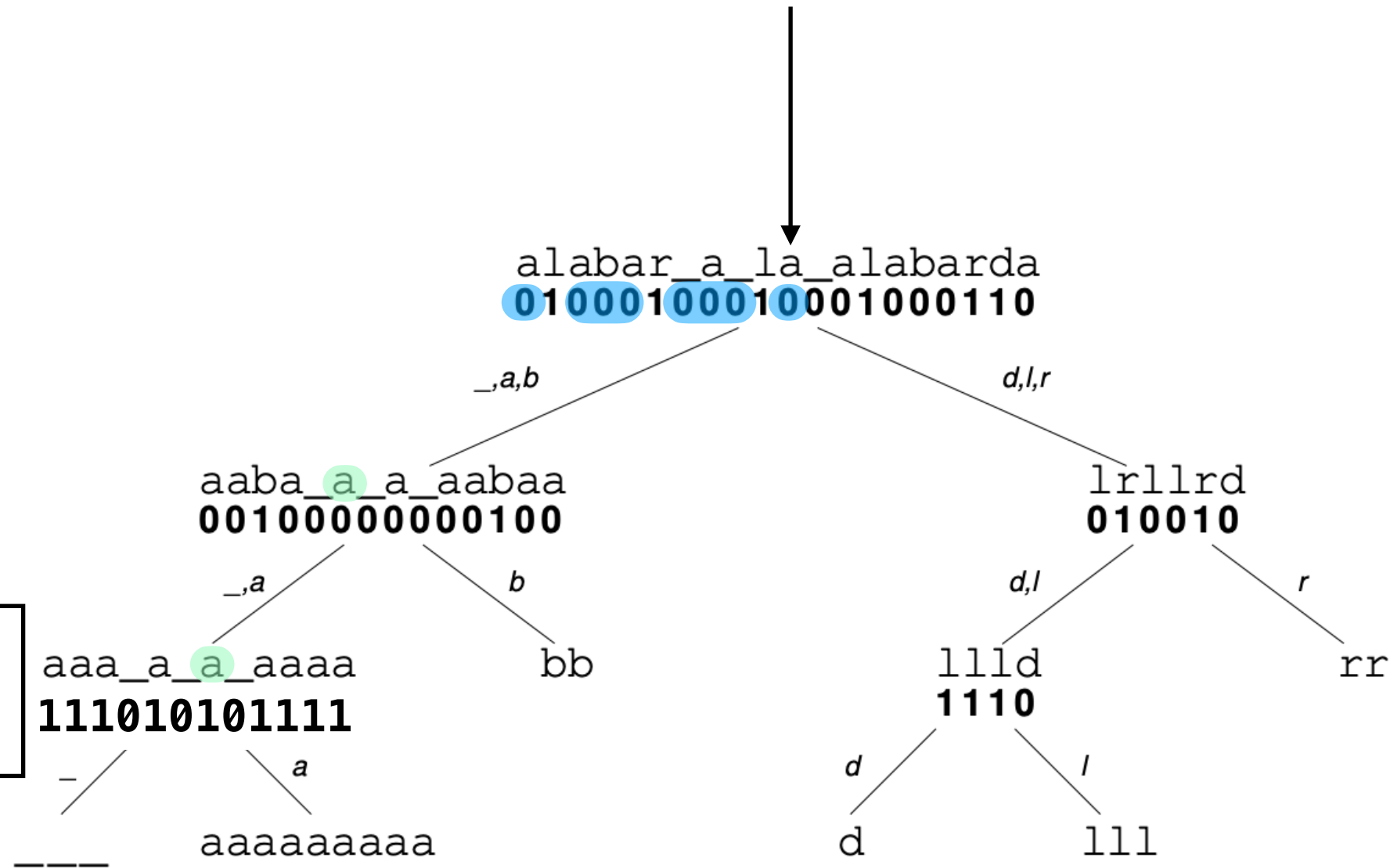


So that maps to 8th character at this level

Count 0's at this level — 7

# Example: Rank

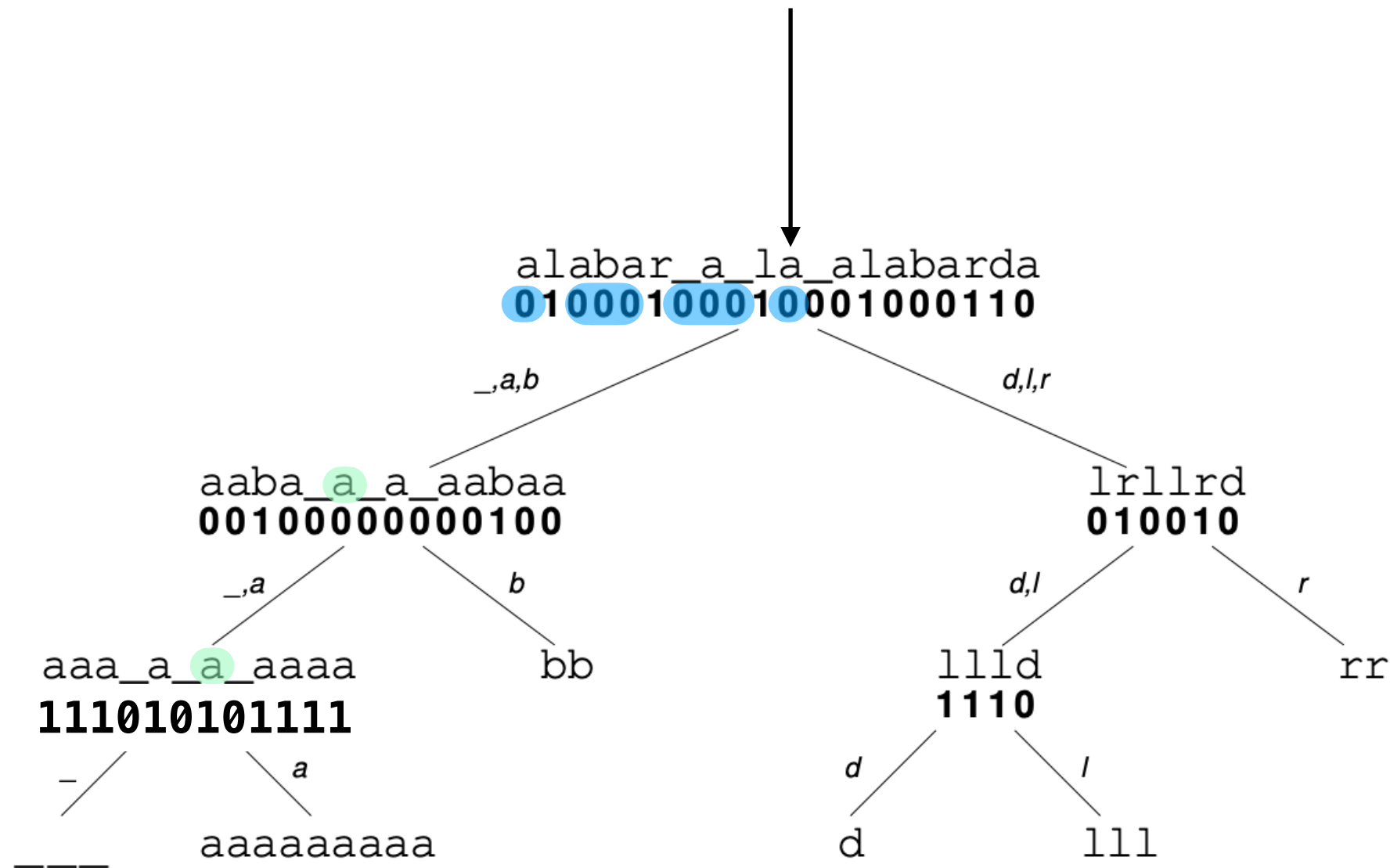
Consider asking for the rank of this “a”



So that maps to  
7th character at  
this level

# Example: Rank

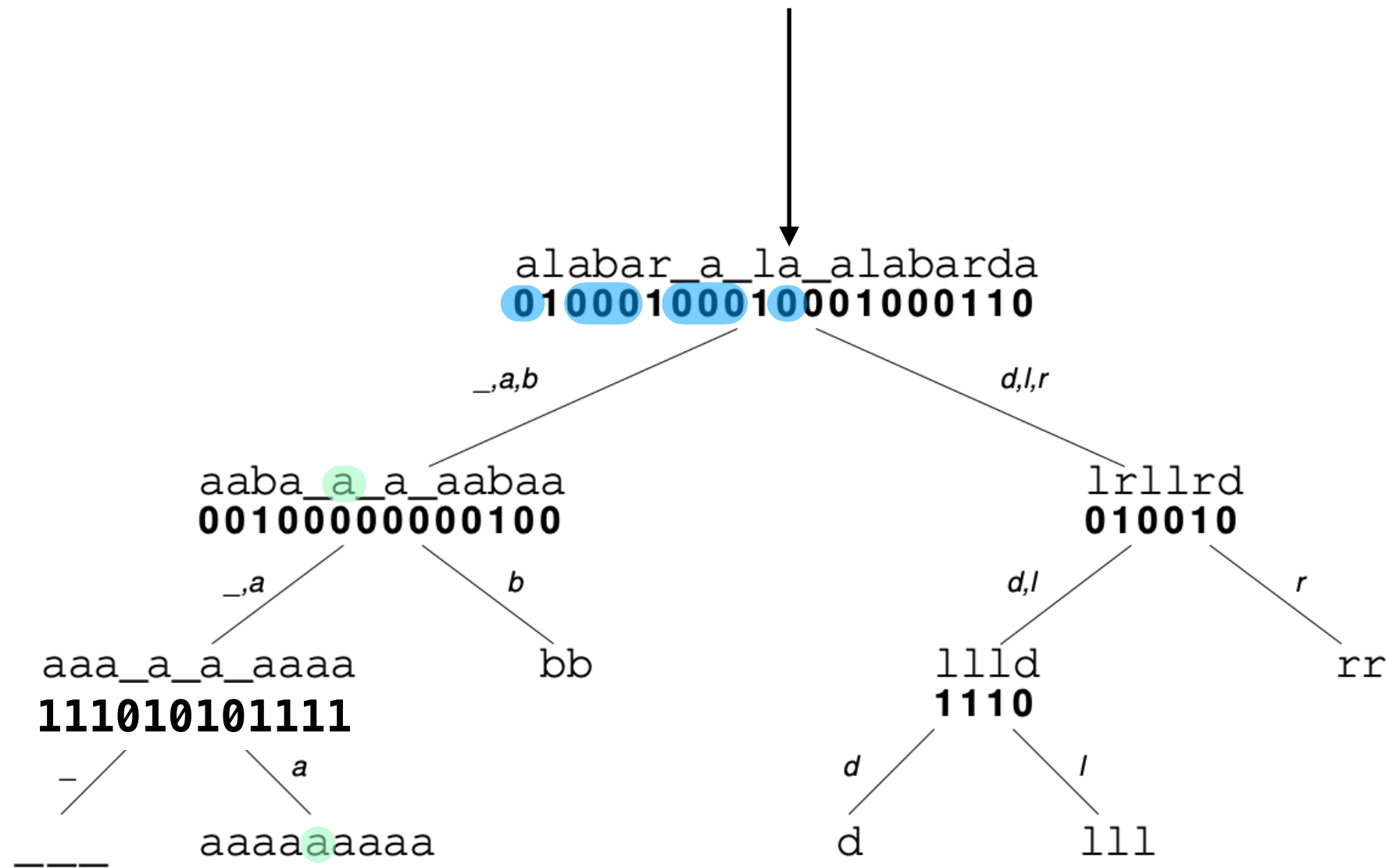
Consider asking for the rank of this “a”





# Example: Rank

Consider asking for the rank of this “a”

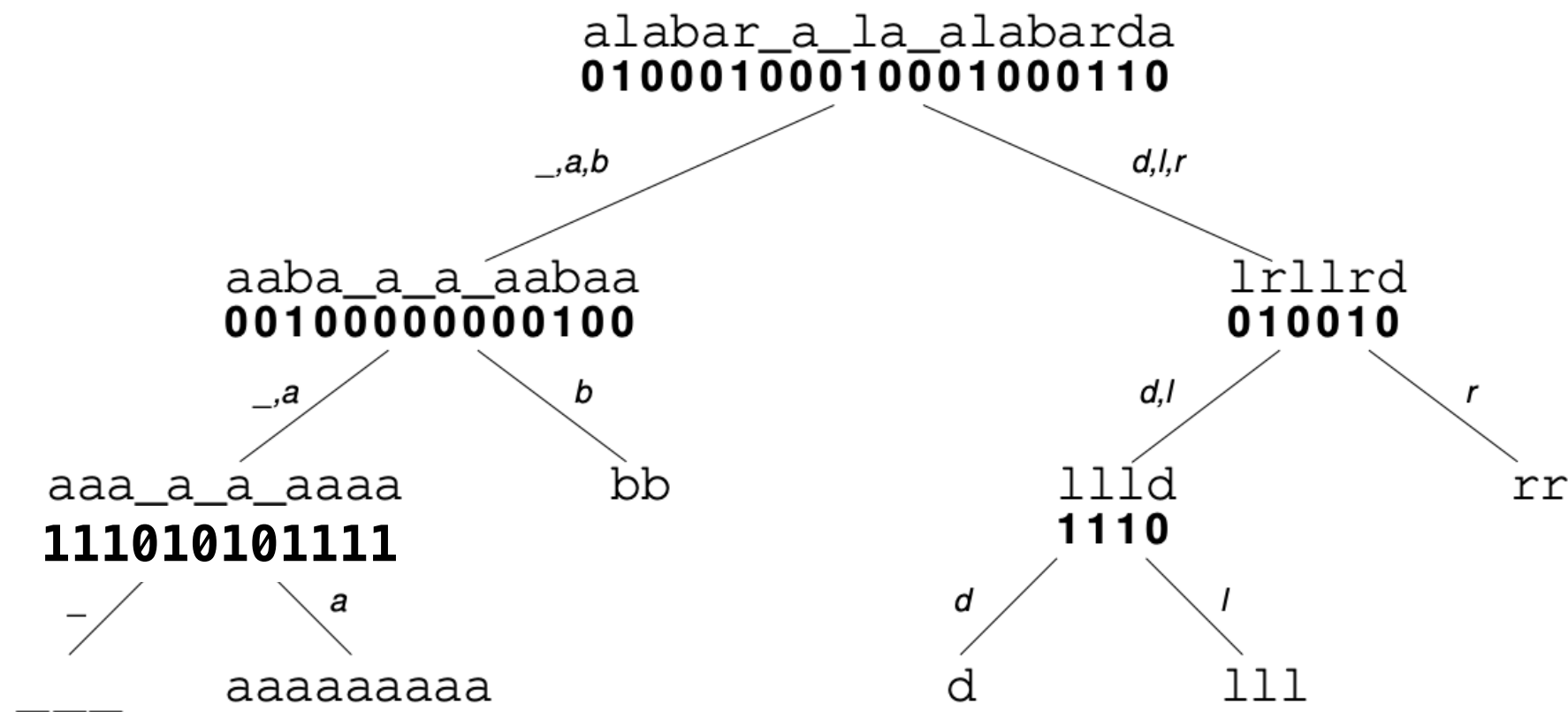


This is the 5th 'a'  
—  $\text{Rank}_a(S, 11) = 5$

If we are 1-indexing

# Example: Rank

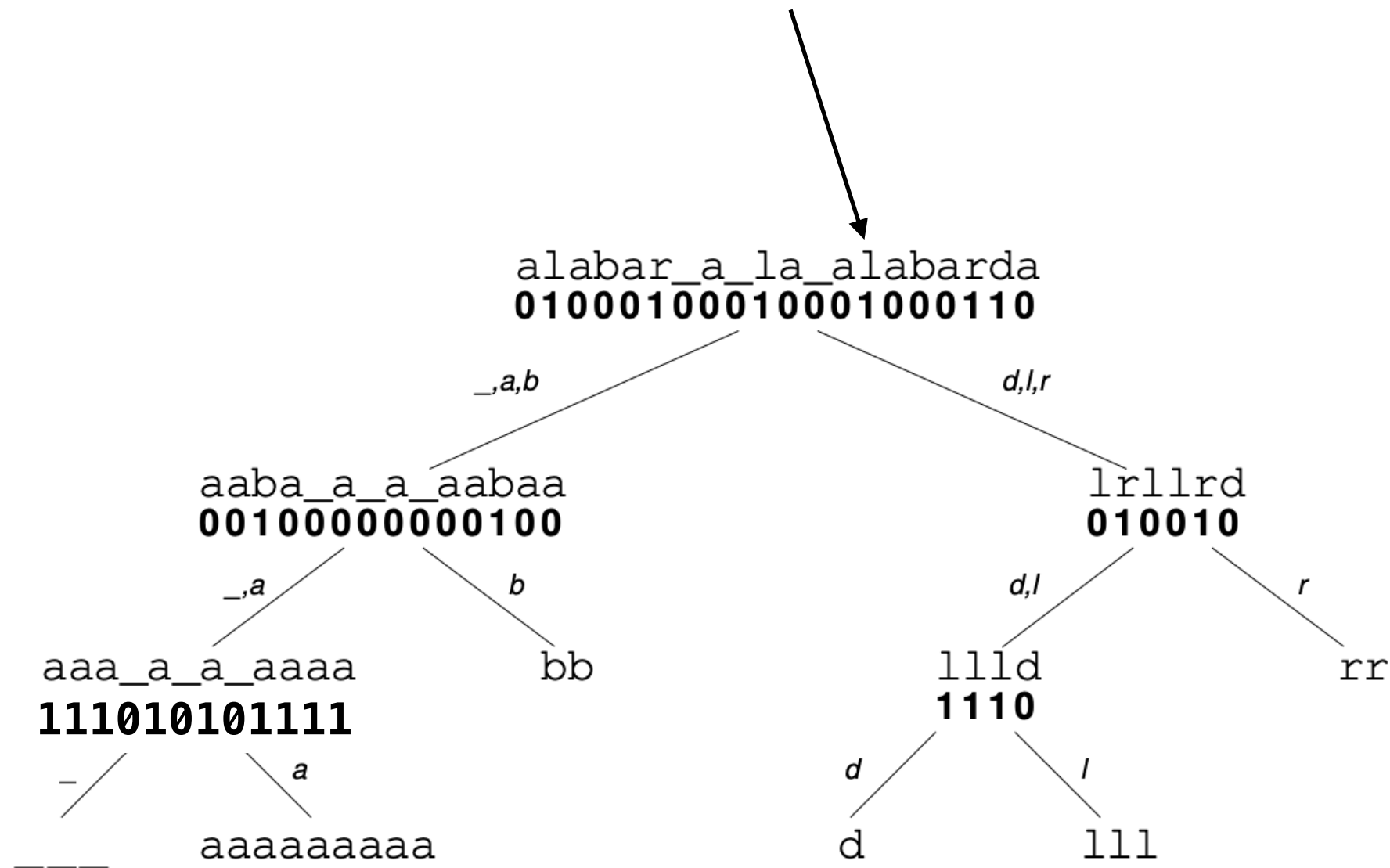
This procedure turns rank for any character in the alphabet into  $\lg \sigma$  rank calculations over bitvectors.



We can answer rank queries for an arbitrary character in  $\lg \sigma * O(1) = O(\lg \sigma)$  time. For small, constant alphabets, through the magic of Big-O, this is *constant time*. :)

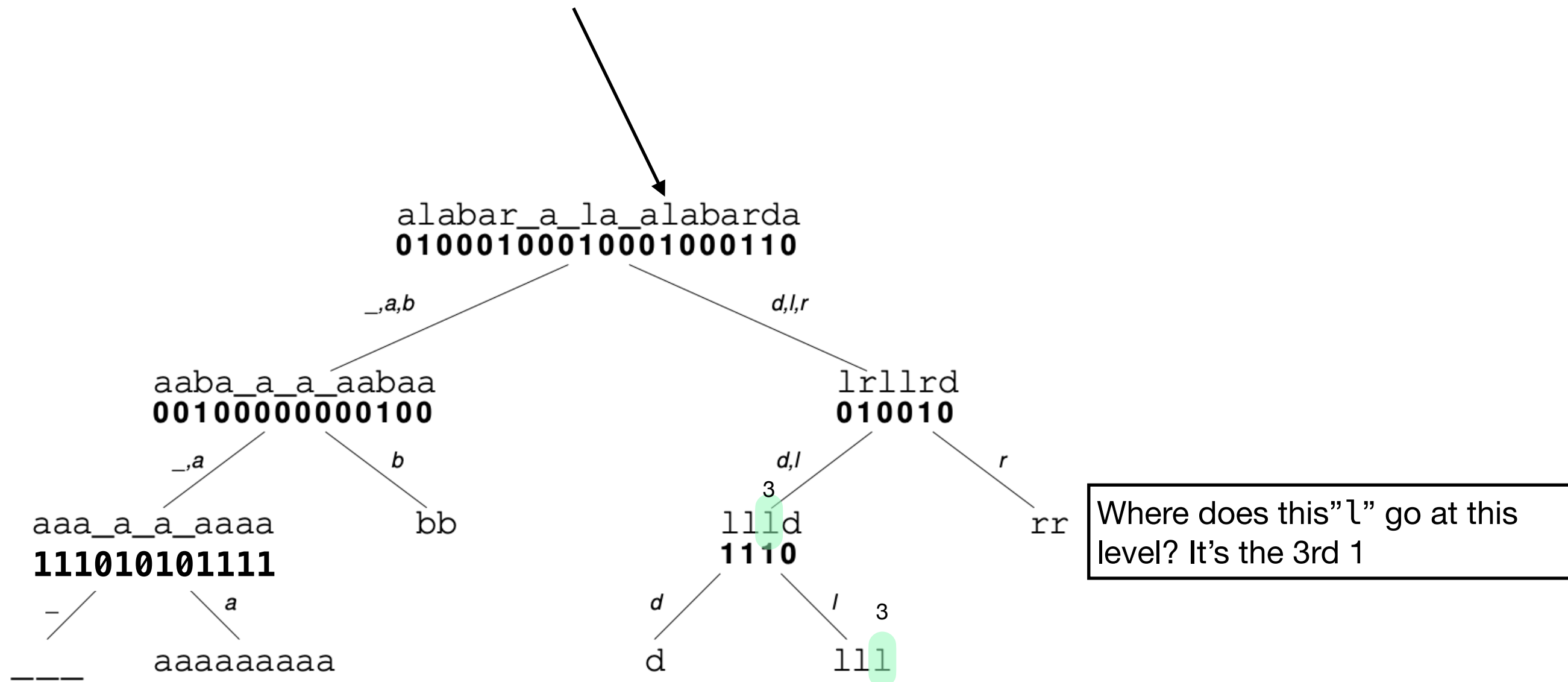
# Example: Select

Select the 3rd “l” (at what index does it occur?)



# Example: Select

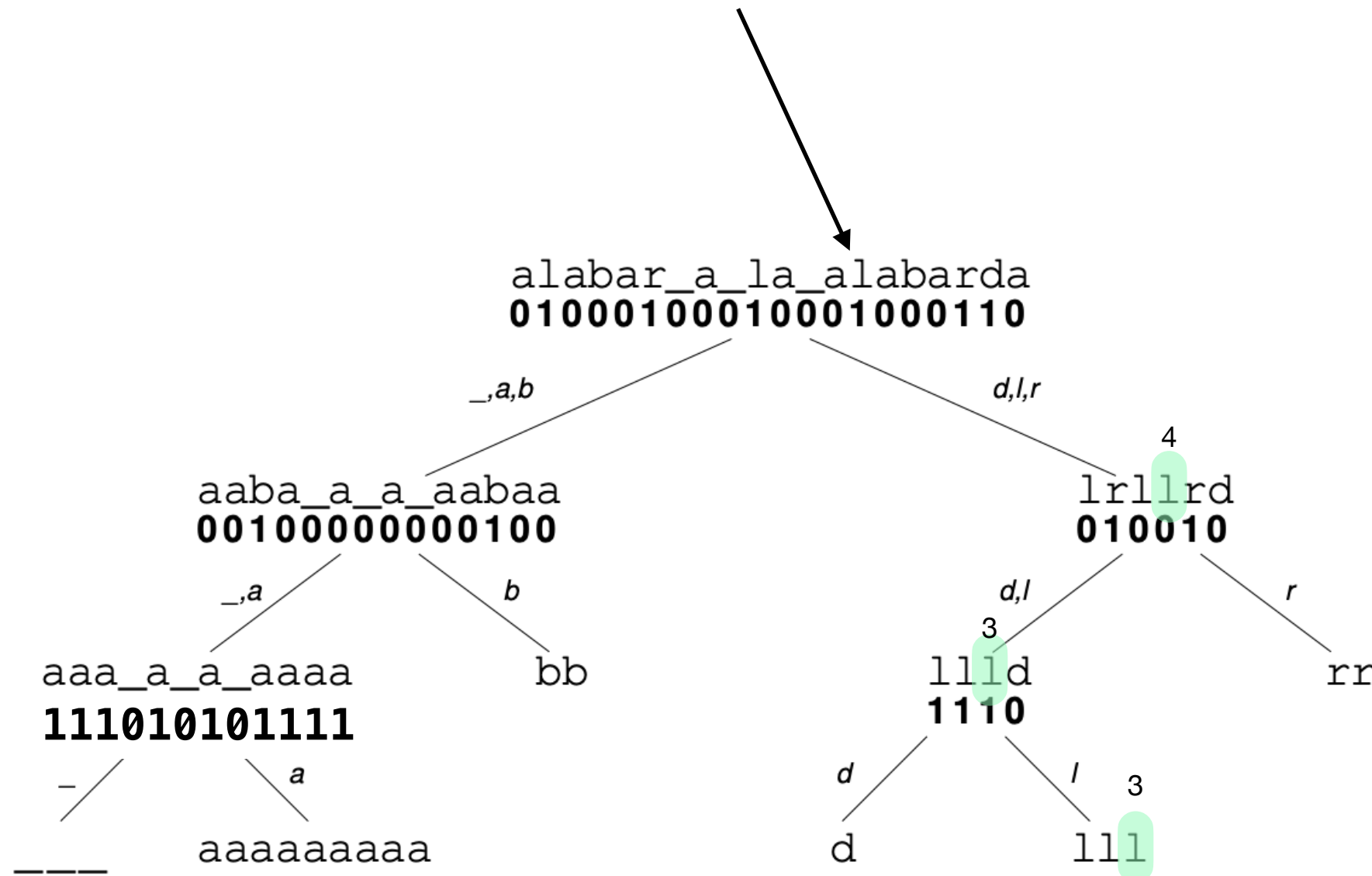
Select the 3rd “l” (at what index does it occur?)



Here, we start at the bottom of the tree and work up.

# Example: Select

Select the 3rd “1” (at what index does it occur?)

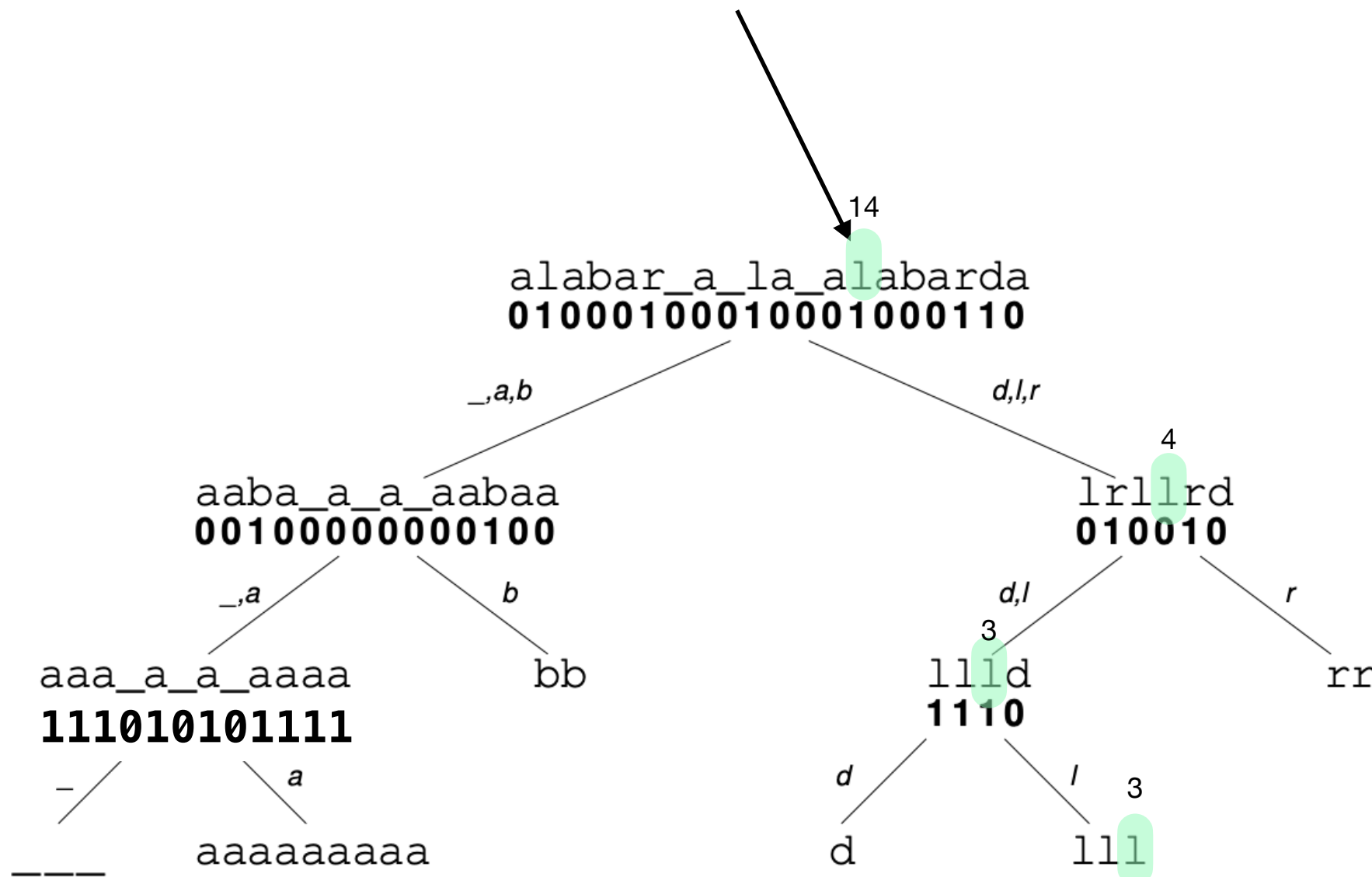


Where does this "l" go at this level? It's the 3rd 0

Here, we start at the bottom of the tree and work up.

# Example: Select

Select the 3rd “1” (at what index does it occur?)

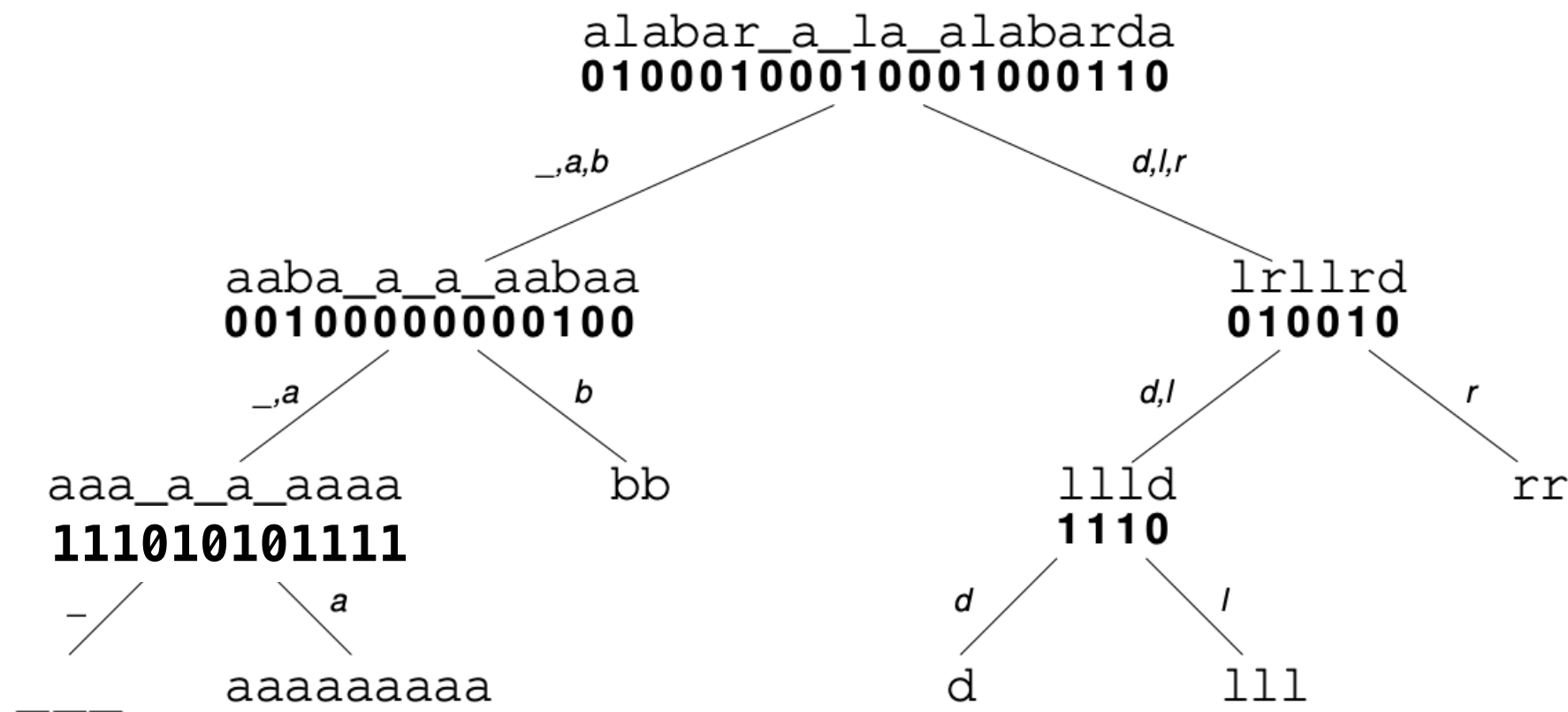


Where does this “l” go at this level? It’s the 4th 1

Here, we start at the bottom of the tree and work up.

# Example: Select

This procedure turns select for any character in the alphabet into  $\lg \sigma$  select calculations over bitvectors.



We can answer select queries for an arbitrary character in  $\lg \sigma * O(1) = O(\lg \sigma)$  time. For small, constant alphabets, through the magic of Big-O, this is *constant time*. :)

# Succinct Data Structures

We have only scratched the surface on what is possible with rank & select and succinct data structures in general.

However, we'll assume familiarity with rank and select moving forward as we talk about data structures in Comp Bio that use them.

Gonzalo Navarro alone publishes 14-24 papers / year in this field :):

<b>2019 (21)</b>	<b>2018 (14)</b>
<b>2017 (21)</b>	<b>2016 (19)</b>

A google search on Gonzalo, and succinct data structures will send you down a wonderful rabbit-hole; I recommend you try it!



# Some practical advice

Succinct data structure papers tend to be quite theoretical (go figure!).

Luckily, there is a *go-to* library for implementation of these ideas.

 [simongog](#) / [sdsl-lite](#)

 Watch

119

 Star

1,582

 Fork

243

 Code

 Issues

39



Pull requests

18



Projects

0



Wiki



Security



Insights

## All your code in one place

Over 40 million developers use GitHub together to host and review code, project manage, and build software together across more than 100 million projects.

[Sign up for free](#)

[See pricing for teams and enterprises](#)

Dismiss

Succinct Data Structure Library 2.0

# Some practical advice

Succinct data structure papers tend to be quite theoretical (go figure!).

Luckily, there is a *go-to* library for implementation of these ideas.

 [simongog](#) / [sdsl-lite](#)

 Watch

119

 Star

1,582

 Fork

243

 Code

 Issues

39

 Pull requests

18

 Projects

0

 Wiki

Security

 Insights

## All your code in one place

Over 40 million developers use GitHub together to host and review code, project manage, and build software together across more than 100 million projects.

[Sign up for free](#)

[See pricing for teams and enterprises](#)

Dismiss

Succinct Data Structure Library 2.0

Provides a modern, modular C++ implementation of many different succinct data structures.

# Some practical advice

Create an FM-index over some text?

```
#include <sdsl/suffix_arrays.hpp>
#include <string>
#include <iostream>

string index_file = string(argv[1])+index_suffix;
csa_wt<wt_huff<rrr_vector<127> >, 512, 1024> fm_index;

if (!load_from_file(fm_index, index_file)) {
    ifstream in(argv[1]);
    if (!in) {
        cout << "ERROR: File " << argv[1] << " does not exist. Exit." << endl;
        return 1;
    }
    cout << "No index "<<index_file<< " located. Building index now." << endl;
    construct(fm_index, argv[1], 1); // generate index
    store_to_file(fm_index, index_file); // save it
}
```

# Some practical advice

Perform rank queries over a bit vector?

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

using namespace std;
using namespace sdsl;

int main()
{
    bit_vector b = bit_vector(8000, 0);
    for (size_t i=0; i < b.size(); i+=100)
        b[i] = 1;
    rank_support_v<1> b_rank(&b);
    for (size_t i=0; i<=b.size(); i+= b.size()/4)
        cout << "(" << i << ", " << b_rank(i) << ") ";
    cout << endl;
}
```

# Some practical advice

How about select<sub>0</sub>

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

using namespace std;
using namespace sdsl;

int main()
{
    bit_vector b = {0,1,0,1,1,1,0,0,0,1,1};
    size_t zeros = rank_support_v<0>(&b)(b.size());
    bit_vector::select_0_type b_sel(&b);

    for (size_t i=1; i <= zeros; ++i) {
        cout << b_sel(i) << " ";
    }
    cout << endl;
}
```



# Some practical advice

You get the idea! An incredibly powerful library, at your fingertips.

## sdsl Cheat Sheet

### Data structures

The library code is in the `sdsl` namespace. Either import the namespace in your program (`using namespace sdsl;`) or qualify all identifiers by a `sdsl::`-prefix.

Each section corresponds to a header file. The file is hyperlinked as part of the section heading.

We have two types of data structures in `sdsl`. *Self-contained* and *support* structures. A support object `s` can extend a self-contained object `o` (e.g. add functionality), but requires access to `o`. Support structures contain the substring support in their class names.

### Integer Vectors (IV)

The core of the library is the class `int_vector<w>`. Parameter  $w$  corresponds to the fixed length of each element in bits. For  $w = 8, 16, 32, 64, 1$  the length is fixed during compile time and the vectors correspond to `std::vector<uintw_t>` resp. `std::vector<bool>`. If  $w = 0$  (default) the length can be set during runtime. *Constructor:* `int_vector<>(n, x, l)`, with  $n$  equals size,  $x$  default integer value,  $l$  width of integer (has no effect for  $w > 0$ ).

*Public methods:* `operator[i]`, `size()`, `width()`, `data()`.

### Manipulating `int_vector<w> v`

Method	Description
<code>v[i]=x</code>	Set entry <code>v[i]</code> to $x$ .
<code>v.width(l)</code>	Set width to $l$ , if $w = 0$ .
<code>v.resize(n)</code>	Resize <code>v</code> to $n$ elements.
Useful methods in namespace <code>sdsl::util</code> :	
<code>set_to_value(v, k)</code>	Set <code>v[i]=k</code> for each $i$ .
<code>set_to_id(v)</code>	Set <code>v[i]=i</code> for each $i$ .
<code>set_random_bits(v)</code>	Set elements to random bits.
<code>mod(v, m)</code>	Set <code>v[i]=v[i] mod m</code> for each $i$ .
<code>bit_compress(v)</code>	Gets $x = \max_i v[i]$ and $l = \lceil \log(x-1) \rceil + 1$ and packs the entries in $l$ -bit integers.
<code>expand_width(v, l)</code>	Expands the width of each integer to $l$ bits, if $l \geq v.width()$ .

### Compressed Integer Vectors (CIV)

For a vector `v`, `enc_vector` stores the self-delimiting coded deltas (`v[i+1]-v[i]`). Fast random access is achieved by sampling values of `v` at rate `t_dens`. Available coder are `coder::elias_delta`, `coder::elias_gamma`, and `coder::fibonacci`.

Class `vlc_vector` stores each `v[i]` as self-delimiting codeword. Samples at rate `t_dens` are inserted for fast random access.

Class `dac_vector` stores for each value  $x$  the least  $(t_b - 1)$  significant bits plus a bit which is set if  $x \geq 2^{b-1}$ . In the latter case, the process is repeated with  $x' = x/2^{b-1}$ .

### Bitvectors (BV)

Representations for a bitvector of length  $n$  with  $m$  set bits.

Class	Description	Space
<code>bit_vector</code>	plain bitvector	$64 \lceil n/64 \rceil + 1$
<code>bit_vector_il</code>	interleaved bitvector	$\approx n(1 + 64/K)$
<code>rrr_vector</code>	$H_0$ -compressed bitvector	$\approx \lceil \log \binom{n}{m} \rceil$
<code>sd_vector</code>	sparse bitvector	$\approx m \cdot (2 + \log \frac{n}{m})$
<code>hyb_vector</code>	hybrid bitvector	

`bit_vector` equals `int_vector<1>` and is therefore dynamic.  
*Public Methods:* `operator[i]`, `size()`, `begin()`, `end()`  
*Public Types:* `rank_1_type`, `select_1_type`, `select_0_type`<sup>1</sup>.  
Each bitvector can be constructed out of a `bit_vector` object.

### Rank Supports (RS)

RSs add rank functionality to BV. Methods `rank(i)` and `operator(i)` return the number of set bits<sup>2</sup> in the prefix  $[0..i]$  of the supported BV for  $i \in [0, n]$ .

Class	Compatible BV	+Bits	Time
<code>rank_support_v</code>	<code>bit_vector</code>	$0.25n$	$O(1)$
<code>rank_support_v5</code>	<code>bit_vector</code>	$0.0625n$	$O(1)$
<code>rank_support_scan</code>	<code>bit_vector</code>	64	$O(n)$
<code>rank_support_il</code>	<code>bit_vector_il</code>	128	$O(1)$
<code>rank_support_rrr</code>	<code>rrr_vector</code>	80	$O(k)$
<code>rank_support_sd</code>	<code>sd_vector</code>	64	$O(\log \frac{n}{m})$
<code>rank_support_hyb</code>	<code>hyb_vector</code>	64	-

Call `util::init_support(rs, bv)` to initialize rank structure `rs` to bitvector `bv`. Call `rs(i)` to get `rank(i) = \sum_{k=0}^{i-1} bv[k]`

### Select Supports (SLS)

SLSs add select functionality to BV. Let  $m$  be the number of set bits in BV. Methods `select(i)` and `operator(i)` return the position of the  $i$ -th set bit<sup>3</sup> in BV for  $i \in [1..m]$ .

Class	Compatible BV	+Bits	Time
<code>select_support_mcl</code>	<code>bit_vector</code>	$\leq 0.2n$	$O(1)$
<code>select_support_scan</code>	<code>bit_vector</code>	64	$O(n)$
<code>select_support_il</code>	<code>bit_vector_il</code>	64	$O(\log n)$
<code>select_support_rrr</code>	<code>rrr_vector</code>	64	$O(\log n)$
<code>select_support_sd</code>	<code>sd_vector</code>	64	$O(1)$

Call `util::init_support(sls, bv)` to initialize `sls` to bitvector `bv`. Call `sls(i)` to get `select(i) = \min\{j \mid \text{rank}(j+1) = i\}`.

### Wavelet Trees (WT=BV+RS+SLS)

Wavelet trees represent sequences over byte or integer alphabets of size  $\sigma$  and consist of a tree of BVs. Rank and select on the sequences is reduced to rank and select on BVs, and the runtime is multiplied by a factor in  $[H_0, \log \sigma]$ .

Class	Shape	lex_ordered	Default alphabet	Traversable
<code>wt_rlmm</code>	underlying WT dependent			×
<code>wt_gmr</code>	none	×	integer	×
<code>wt_ap</code>	none	×	integer	×
<code>wt_huff</code>	Huffman	×	byte	✓
<code>wm_int</code>	Balanced	×	integer	✓
<code>wt_blcd</code>	Balanced	✓	byte	✓
<code>wt_hutu</code>	Hu-Tucker	✓	byte	✓
<code>wt_int</code>	Balanced	✓	integer	✓

*Public types:* `value_type`, `size_type`, and `node_type` (if WT is

traversable). In the following let  $c$  be a symbol,  $i, j, k$ , and  $q$  integers,  $v$  a node, and  $r$  a range.

*Public methods:* `size()`, `operator[i]`, `rank(i, c)`, `select(i, c)`, `inverse_select(i)`, `begin()`, `end()`.  
Traversable WTs provide also: `root()`, `is_leaf(v)`, `empty(v)`, `size(v)`, `sym(v)`, `expand(v)`, `expand(v, r)`, `expand(v, std::vector<r>)`, `bit_vec(v)`, `seq(v)`.  
lex\_ordered WTs provide also: `lex_count(i, j, c)` and `lex_smaller_count(i, c)`. `wt_int` provides: `range_search_2d`.  
`wt_algorithm.hpp` contains the following generic WT method (let `wt` be a WT object): `intersect(wt, vector<r>)`, `quantile_freq(wt, i, j, q)`, `interval_symbols(wt, i, j, k, ...)`, `symbol_lte(wt, c)`, `symbol_gte(wt, c)`, `restricted_unique_range_values(wt, x_i, x_j, y_i, y_j)`.

### Suffix Arrays (CSA=IV+WT)

Compressed suffix arrays use CIVs or WTs to represent the suffix arrays (SA), its inverse (ISA), BWT,  $\Psi$ , and LF. CSAs can be built over byte and integer alphabets.

Class	Description
<code>csa_bitcompressed</code>	Based on SA and ISA stored in a IV.
<code>csa_sada</code>	Based on $\Psi$ stored in a CIV.
<code>csa_wt</code>	Based on the BWT stored in a WT.

*Public methods:* `operator[i]`, `size()`, `begin()`, `end()`.  
*Public members:* `isa`, `bwt`, `lf`, `psi`, `text`, `L`, `F`, `C`, `char2comp`, `comp2char`, `sigma`.  
*Policy classes:* alphabet strategy (e.g. `byte_alphabet`, `succinct_byte_alphabet`, `int_alphabet`) and SA sampling strategy (e.g. `sa_order_sa_sampling`, `text_order_sa_sampling`)

### Longest Common Prefix (LCP) Arrays

Class	Description
<code>lcp_bitcompressed</code>	Values in a <code>int_vector&lt;&gt;</code> .
<code>lcp_dac</code>	Direct accessible codes used.
<code>lcp_byte</code>	Small values in a byte; 2 words per large.
<code>lcp_wt</code>	Small values in a WT; 1 word per large.
<code>lcp_vlc</code>	Values in a <code>vlc_vector</code> .
<code>lcp_support_sada</code>	Values stored permuted. CSA needed.
<code>lcp_support_tree</code>	Only depths of CST inner nodes stored.
<code>lcp_support_tree2</code>	+ large values are sampled using LF.

*Public methods:* `operator[i]`, `size()`, `begin()`, `end()`

### Balanced Parentheses Supports (BPS)

We represent a sequence of parentheses as a `bit_vector`. An opening/closing parenthesis corresponds to 1/0.

Class	Description
<code>bp_support_g</code>	Two-level pioneer structure.
<code>bp_support_gg</code>	Multi-level pioneer structure.
<code>bp_support_sada</code>	Min-max-tree over excess sequence.

*Public methods:* `find_open(i)`, `find_close(i)`, `enclose(i)`, `double_enclose(i, j)`, `excess(i)`, `rr_enclose(i, j)`, `rank(i)`<sup>4</sup>, `select(i)`.  
Call `util::init_support(bps, bv)` to initialize a BPS `bps` to `bit_vector bv`.



# Some practical advice

You get the idea! An incredibly powerful library, at your fingertips.

## Suffix Trees (**CST**=**CSA**+**LCP**+**BPS**)

A CST can be parametrized by any combination of CSA, LCP, and BPS. The operation of each part can still be accessed through member variables. The additional operations are listed below. CSTs can be built for byte or integer alphabets.

Class	Description
<code>cst_sada</code>	Represents a node as position in BPS. Navigational operations are fast (they are directly translated in BPS operations on the DFS-BPS). Space: $4n + o(n) +  CSA  +  LCP $ bits.
<code>cst_sct3</code>	Represents nodes as intervals. Fast construction, but slower navigational operations. Space: $3n + o(n) +  CSA  +  LCP $

**Public types:** `node_type`. In the following let  $v$  and  $w$  be nodes and  $i, d, lb, rb$  integers.

**Public methods:** `size()`, `nodes()`, `root()`, `begin()`, `end()`, `begin_bottom_up()`, `end_bottom_up`, `size(v)`, `is_leaf(v)`, `degree(v)`, `depth(v)`, `node_depth(v)`, `edge(v, d)`, `lb(v)`, `rb(v)`, `id(v)`, `inv_id(i)`, `sn(v)`, `select_leaf(i)`, `node(lb, rb)`, `parent(v)`, `sibling(v)`, `lca(v, w)`, `select_child(v, i)`, `child(v, c)`, `children(v)`, `sl(v)`, `wl(v, c)`, `leftmost_leaf(v)`, `rightmost_leaf(v)`

**Public members:** `csa`, `lcp`.

The `traversal example` shows how to use the DFS-iterator.

## Range Min/Max Query (**RMQ**)

A RMQ `rmq` can be used to determine the position of the minimum value<sup>5</sup> in an arbitrary subrange  $[i, j]$  of an preprocessed vector  $v$ . Operator `operator(i, j)` returns  $x = \min\{r \mid r \in [i, j] \wedge v[r] \leq v[k] \ \forall k \in [i, j]\}$

Class	Space	Time
<code>rmq_support_sparse_table</code>	$n \log^2 n$	$O(1)$
<code>rmq_succint_sada</code>	$4n + o(n)$	$O(1)$
<code>rmq_succint_sct</code>	$2n + o(n)$	$O(1)$

## Constructing data structures

Let  $o$  be a WT-, CSA-, or CST-object. Object  $o$  is built with `construct(o, file, num_bytes=0)` from a sequence stored in `file`. File is interpreted dependent on the value of `num_bytes`:

Value	File interpreted as
<code>num_bytes=0</code>	serialized <code>int_vector&lt;&gt;</code> .
<code>num_bytes=1</code>	byte sequence of length <code>util::file_size(file)</code> .
<code>num_bytes=2</code>	16-bit word sequence.
<code>num_bytes=4</code>	32-bit word sequence.
<code>num_bytes=8</code>	64-bit word sequence.
<code>num_bytes=d</code>	Parse decimal numbers.

Note: `construct` writes/reads data to/from disk during construction. Accessing disk for small instances is a considerable overhead. `construct_im(o, data, num_bytes=0)` will build  $o$  using only main memory. Have a look at [this handy tool for an example](#).

## Configuring construction

The locations and names of the intermediate files can be configured by a `cache_config` object. It is constructed by `cache_config(del, tmp_dir, id, map)` where `del` is a boolean variable which specifies if the intermediate files should be deleted after construction, `tmp_dir` is a path to the directory

where the intermediate files should be stored, `id` is used as part of the file names, and `map` contains a mapping of keys (e.g. `conf::KEY_BWT`, `conf::KEY_SA`, ...) to file paths.

The `cache_config` parameter extends the construction method to: `construct(o, file, config, num_bytes)`.

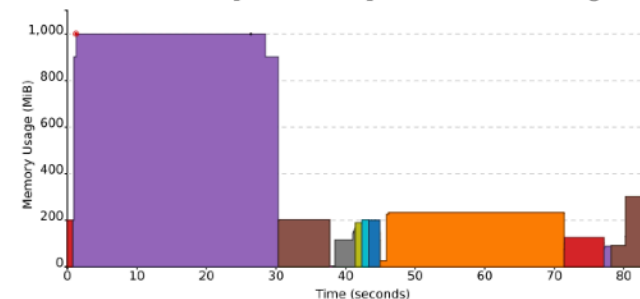
The following methods (`key` is a key string, `config` represent a `cache_config` object, and  $o$  a `sds` object) should be handy in customized construction processes:

`cache_file_name(key, config)`  
`cache_file_exists(key, config)`  
`register_cache_file(key, config)`  
`load_from_cache(o, key, config)`  
`store_to_cache(o, key, config)`

## Resource requirements

**Memory:** The memory peak of CSA and CST construction occurs during the SA construction, which is 5 times the texts size for byte-alphabets and inputs < 2 GiB (see the Figure below for a 200 MB text) and 9 times for larger inputs. For integer alphabets the construction takes about twice the space of the resulting output.

**Time:** A CST construction processes at about 2 MB/s. The Figure below shows the resource consumption during the construction of a `cst_sct3<>` CST for [200 MB English text](#). For a detailed description of the phases click on the figure.



This diagram was generated using the sample program [memory-visualization.cpp](#).

## Reading and writing data

### Importing data into `sds` structures

`load_vector_from_file(v, file, num_bytes)`  
Load file into an `int_vector v`. Interpretation of file depends on `num_bytes`; see method `construct`.

### Store `sds` structures

Use `store_to_file(o, file)` to store an `sds` object  $o$  to file. Object  $o$  can also be serialized into a `std::ostream`-object out by the call `o.serialize(out)`.

### Load `sds` structures

Use `load_from_file(o, file)` to load an `sds` object  $o$ , which is stored in file. Call `o.load(in)` reads  $o$  from `std::istream`-object in.

## Utility methods

More useful methods in the `sds1::util` namespace:

Method	Description
<code>pid()</code>	Id of current process.
<code>id()</code>	Get unique id inside the process.
<code>basename(p)</code>	Get filename part of a path $p$ .
<code>dirname(p)</code>	Get directory part of a path $p$ .
<code>demangle(o)</code>	Demangles output of <code>typeid(o).name()</code> .
<code>demangle2(o)</code>	Simplifies output of <code>demangle</code> . E.g. removes <code>sds1::</code> -prefixes, ...
<code>to_string(o)</code>	Transform object $o$ to a string.
<code>assign(o1, o2)</code>	Assign $o1$ to $o2$ , or swap $o1$ and $o2$ if the objects are of the same type.
<code>clear(o)</code>	Set $o$ to the empty object.

## Measuring and Visualizing Space

`size_in_bytes(o)` returns the space used by an `sds` object  $o$ . Call `write_structure<JSON_FORMAT>(o, out)` to get a detailed space breakdown written in JSON format to stream out. `<HTML_FORMAT>` will write a HTML page ([like this](#)), which includes an interactive SVG-figure.

## Methods on words

Class `bits` contains various fast methods on a 64-bit word  $x$ . Here the most important ones.

Method	Description
<code>bits::cnt(x)</code>	Number of set bits in $x$ .
<code>bits::sel(x, i)</code>	Position of $i$ -th set bit, $i \in [0, \text{cnt}(x) - 1]$ .
<code>bits::lo(x)</code>	Position of least significant set bit.
<code>bits::hi(x)</code>	Position of most significant set bit.

**Note:** Positions in  $x$  start at 0. `lo` and `hi` return 0 for  $x = 0$ .

## Tests

A `make test` call in the `test` directory, downloads test inputs, compiles tests, and executes them.

## Benchmarks

Directory `benchmark` contains configurable benchmarks for various data structure, like WTs, CSAs/FM-indexes (measuring time and space for operations `count`, `locate`, and `extract`).

## Debugging

You get the gdb command `pv <int_vector> <idx1> <idx2>`, which displays the elements of an `int_vector` in the range  $[\text{idx1}, \text{idx2}]$  by appending the file `sds1.gdb` to your `.gdbinit`.

© Simon Gog  
Cheatsheet template provided by Winston Chang  
<http://www.stdout.org/~winston/latex/>

## Notes

- <sup>1</sup> `select_0_type` not defined for `sd_vector`.
- <sup>2</sup> It is also possible to rank 0 or the patterns 10 and 01.
- <sup>3</sup> It is also possible to select 0 or the patterns 10 and 01.
- <sup>4</sup> For PBS the bits are counted in the prefix  $[0..i]$ .
- <sup>5</sup> Or maximum value; can be set by a template parameter.