

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 $\ell = \lceil \log(x-1) \rceil + 1$ and packs the entries in ℓ -bit integers.
<code>expand_width(v, ℓ)</code>	Expands the width of each integer to ℓ bits, if $\ell \geq v.width()$.

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` 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.

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	$\approx n(1 + K/64)$
<code>rrr_vector</code>	H_0 -compressed bitvector	$\approx \lceil \log \binom{m}{n} \rceil$
<code>sd_vector</code>	sparse bitvector	$\approx 2m \cdot \log \frac{n}{m}$

`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`¹. 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 bits² 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>	$0.25n$	$\mathcal{O}(1)$
<code>rank_support_v5</code>	<code>bit_vector</code>	$0.0625n$	$\mathcal{O}(1)$
<code>rank_support_scan</code>	<code>bit_vector</code>	64	$\mathcal{O}(n)$
<code>rank_support_il</code>	<code>bit_vector_il</code>	128	$\mathcal{O}(1)$
<code>rank_support_rrr</code>	<code>rrr_vector</code>	80	$\mathcal{O}(k)$
<code>rank_support_sd</code>	<code>sd_vector</code>	64	$\mathcal{O}(1)$

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 bit³ in BV for $i \in [1..m]$.

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

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 σ 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
<code>wt_blcd</code>	Balanced	✓	byte
<code>wt_huff</code>	Huffman	×	byte
<code>wt_hutu</code>	Hu-Tucker	✓	integer
<code>wt_int</code>	Balanced	✓	byte
<code>wt_rlmn</code>	underlying WT dependent		

Public methods: `operator[i]`, `rank(i, c)`, `select(i, c)`,

`inverse_select(i, c)`, `size()`, `begin()`, `end()`

Order preserving WTs (`lex_ordered`) provide more methods: `lex_count` and `lex_smaller_count`. In addition `wt_int` provides: `range_search_2d`, `quantil_freq`, `topk_greedy`, `topk_qprobing`, `intersect`.

Suffix Arrays (CSA=IV+WT)

Compressed suffix arrays use CIVs or WTs to represent the suffix arrays (SA), its inverse (ISA), BWT, Ψ , 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 Ψ 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`, `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<></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.
<i>Public methods</i> : <code>operator[i]</code> , <code>size()</code> , <code>begin()</code> , <code>end()</code>	

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

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: `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)`, `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 value⁵ 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<></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)
```

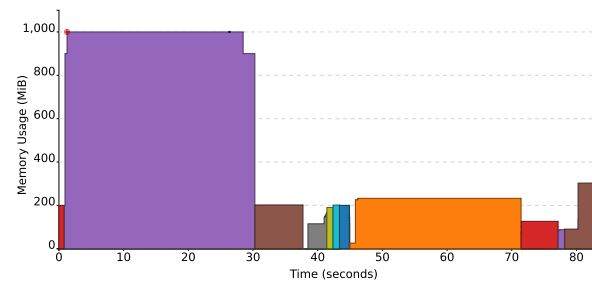
Resource requirements

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 Phase #2 in 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](#).

Phases description: #1=copy TEXT, #2=SA, #3=BWT, #4=wt_huff<>, #7=LCP, #8=BP, #9=bp_support_sada<>, #10=lcp_dac<>.



This diagram was generated from the output of this [example program](#), which print the memory log to `stdout`.

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 <code>p</code> .
<code>dirname(p)</code>	Get directory part of a path <code>p</code> .
<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 <code>o</code> to a string.
<code>assign(o1, o2)</code>	Assign <code>o1</code> to <code>o2</code> , or swap <code>o1</code> and <code>o2</code> if the objects are of the same type.
<code>clear(o)</code>	Set <code>o</code> 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 <code>x</code> .
<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 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

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