# Counting with Prediction: Rank and Select Queries with Adjusted Anchoring
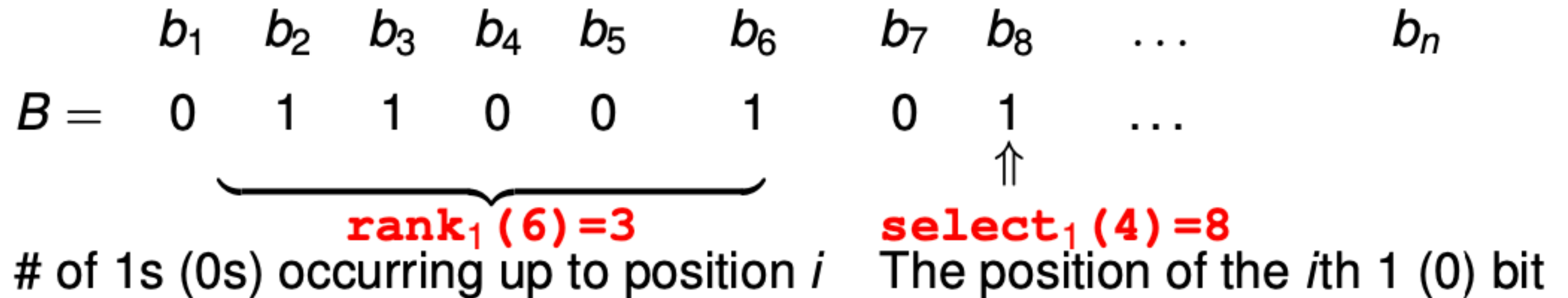
Oguzhan Kulekci
Indiana University Bloomington
okulekci@iu.edu

# Rank & Select Queries

$$b_1 \quad b_2 \quad b_3 \quad b_4 \quad b_5 \quad b_6 \quad b_7 \quad b_8 \quad \ldots \quad b_n$$

$$B = \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 0 \quad 1 \quad \ldots$$

$$\mathtt{rank_1(6)=3} \qquad \mathtt{select_1(4)=8}$$

\# of 1s (0s) occurring up to position $i$     The position of the $i$th 1 (0) bit
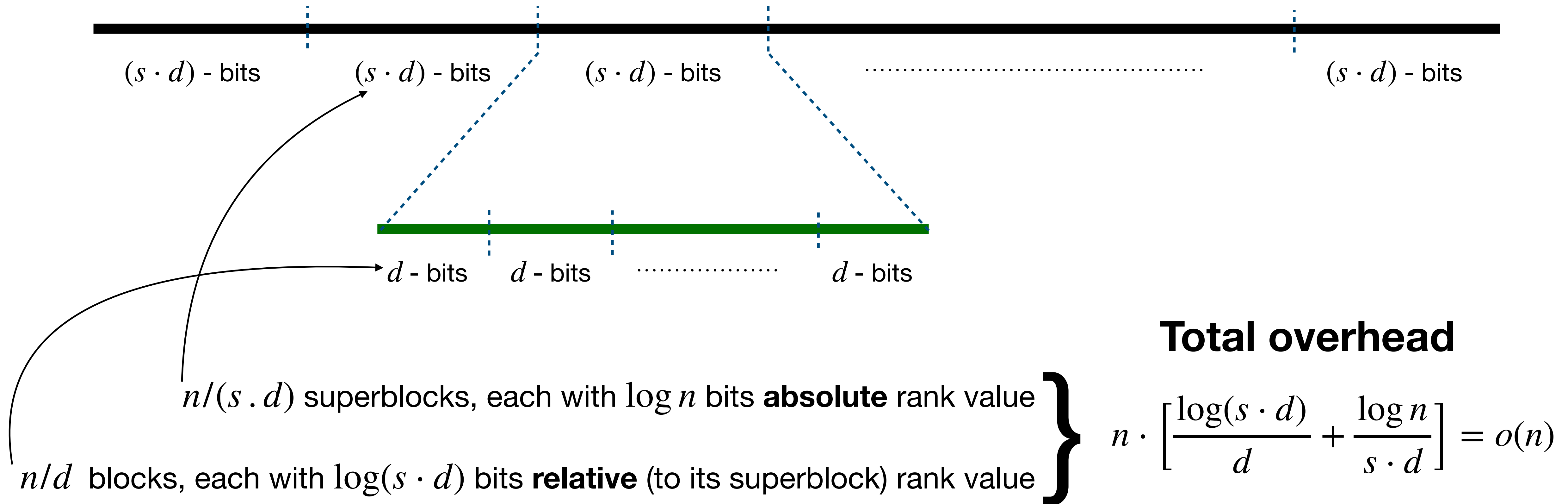
- The fundamental building block in compressed data structures.

- Deeply studied for more than 30 years (please see the references in the paper)

# Rank&Select Dictionaries

**Maintain a dictionary of rank values for some positions and use it to answer queries efficiently.**

$(s \cdot d)$ - bits   $(s \cdot d)$ - bits   $(s \cdot d)$ - bits   $(s \cdot d)$ - bits

$d$ - bits   $d$ - bits   $d$ - bits

$n/(s \cdot d)$ superblocks, each with $\log n$ bits **absolute** rank value

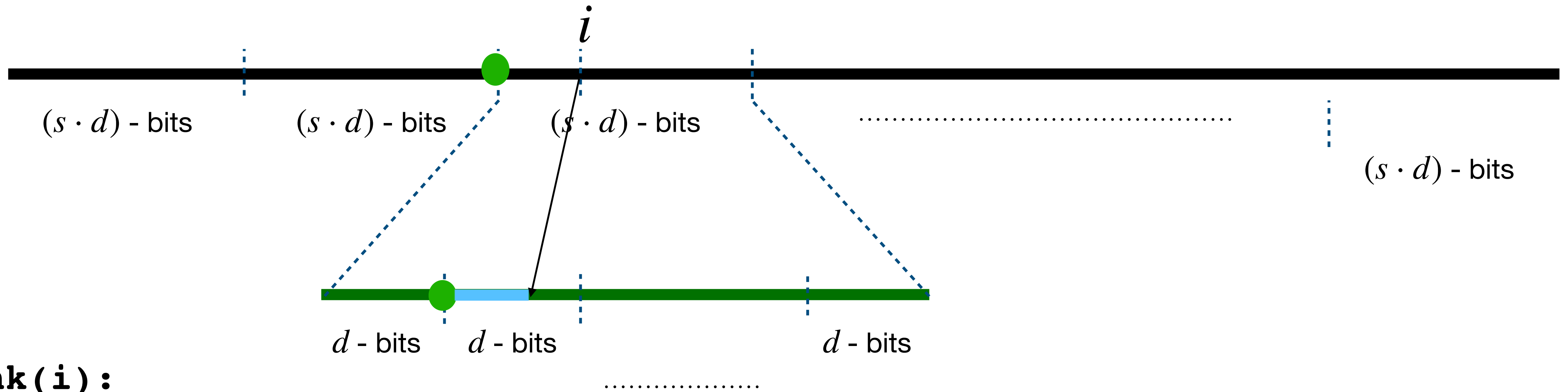$n/d$ blocks, each with $\log(s \cdot d)$ bits **relative** (to its superblock) rank value

**Total overhead**

$$n \cdot \left[ \frac{\log(s \cdot d)}{d} + \frac{\log n}{s \cdot d} \right] = o(n)$$

*,which becomes $o(n)$ with proper selection of $s$ and $d$*

# Rank&Select Dictionaries

$i$

$(s \cdot d)$ - bits      $(s \cdot d)$ - bits      $(s \cdot d)$ - bits      ............................................

$(s \cdot d)$ - bits

$d$ - bits      $d$ - bits      $d$ - bits

....................

`Rank(i):`

- Sum the corresponding superblock and block rank values ● from the maintained dictionary

- Add the number of set bits detected inside the inner-block up to the queried position ▬

  - **SIMD instructions are used to compute this value fast in constant time.**

Overall, $O(1)$-time solution for rank with $o(n)$ overhead.
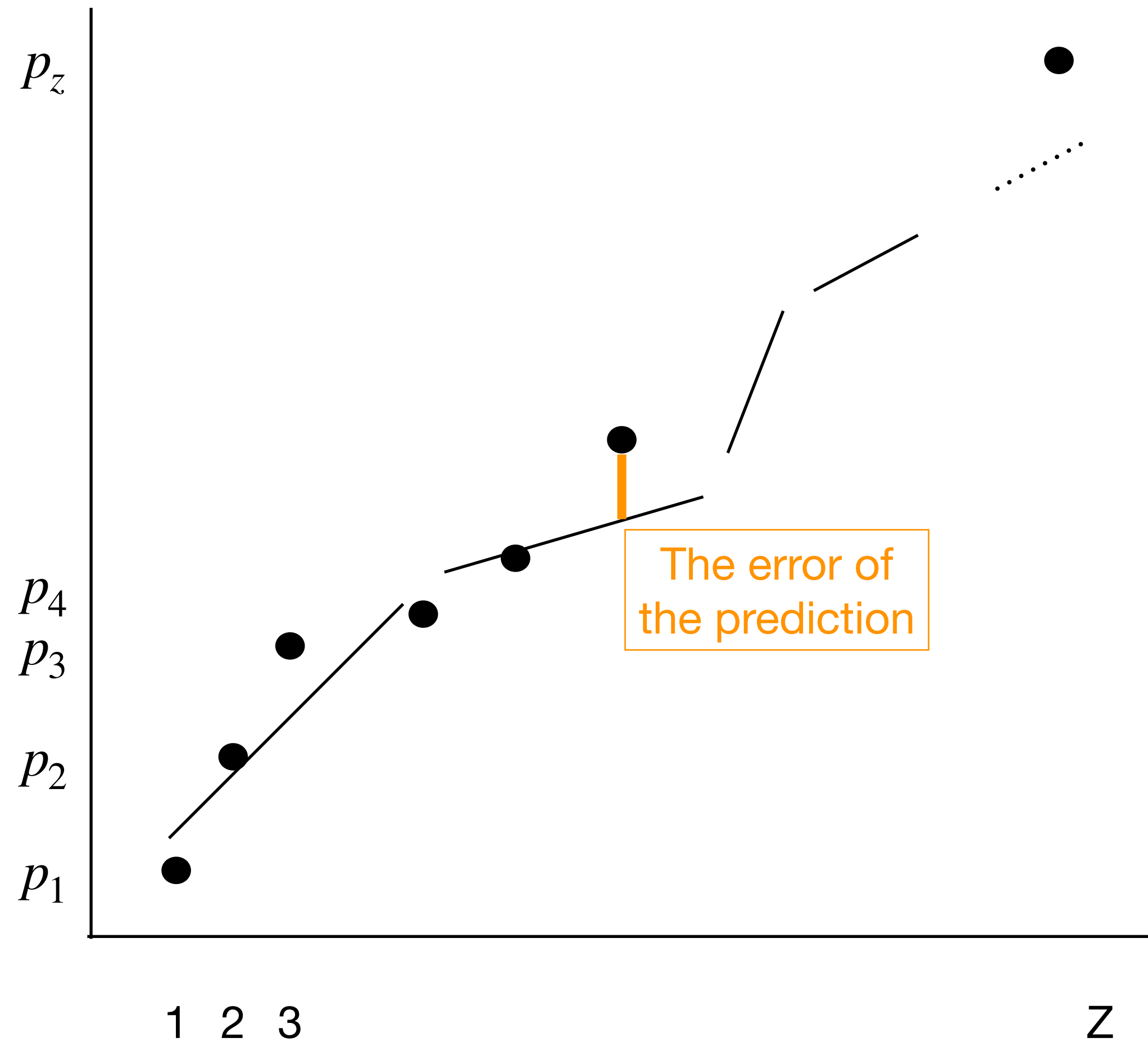
`Select(i):`

- Yes, it needs a different data structure with variable size blocks to answer in constant time.

# Rank&Select Dictionaries

- In practice, the data structures **favor either rank or select** operations, but not both ! Usually, the dictionaries constructed with fixed-size blocks favor rank, and variable-size ones select.

- If the bitmap is sparse (the polarity is far from 0.5) then keeping the bitmap compressed make sense, and leads to **compressed R&S solutions**. (Although they are a bit slow in practice still)

- New research direction by **using machine learning techniques in data structure design**, *Learned data structures, Ferragina, Vinviguerra, 2020.*

- ***Boffa et al, ALENEX'21:*** *A learned-approach to quicken and compress R&S Dictionaries*
  - *Targets **compressed** bitmap, favors* `SELECT` *with variable size blocks*

- ***This study***
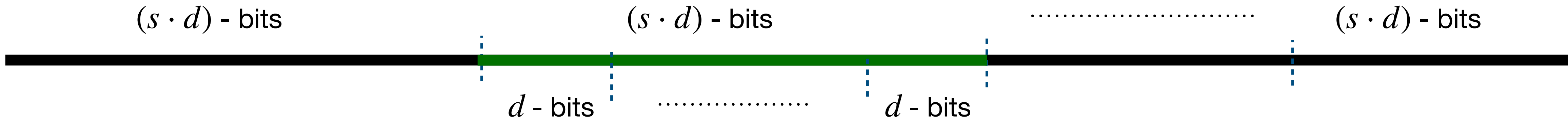  - *Targets **uncompressed** bitmaps, and favors* `RANK` *with fixed size blocks*

# Previous Work: R&S with Learning Approach



The error of the prediction

- The **positions of the set bits** on a given bitmap is a sequence of **increasing** integers
  $$P = \langle p_1, p_2, \ldots, p_z \rangle$$

- Fit $\ell$ lines, where each covers variable number of positions such that the error between the prediction and actual value is denotable by $c$ bits.

- Maintain the parameters for each $ax + b$ line, and also the $c$ bit correction values for each position

- Also some metadata for the number of positions covered per each line is stored

- **The $\ell$ depends on the regularity and number of the positions**

- **Select(i):** Simply go the the line corresponding line, get the prediction and correct it.

- **Rank(i):** Search which line includes I and search the closest previous position on that line.

# Proposed Data Structure

$(s \cdot d)$ - bits           $(s \cdot d)$ - bits        ..........................    $(s \cdot d)$ - bits

$d$ - bits    ...................    $d$ - bits

- The block rank values in each super-block is an increasing sequence

  - $B_1 = \langle b_1^1, b_2^1, \ldots b_{s-1}^1 \rangle, \quad \ldots \ldots ,$
    $B_{n/sd} = \langle b_1^{n/sd}, b_2^{n/d}, \ldots b_{s-1}^{n/sd} \rangle$

  - Linear regression ($ax + b$) per each $B_i$ and **store** the $(a, b)$ values that occupies $\dfrac{n}{s \cdot d} \cdot (32 + 32)$ bits.

- Per each block maintain **a single validity** bit and a $\log m$ bit **correction value** that consumes $\dfrac{n}{d} \cdot (1 + \log m)$ bits.
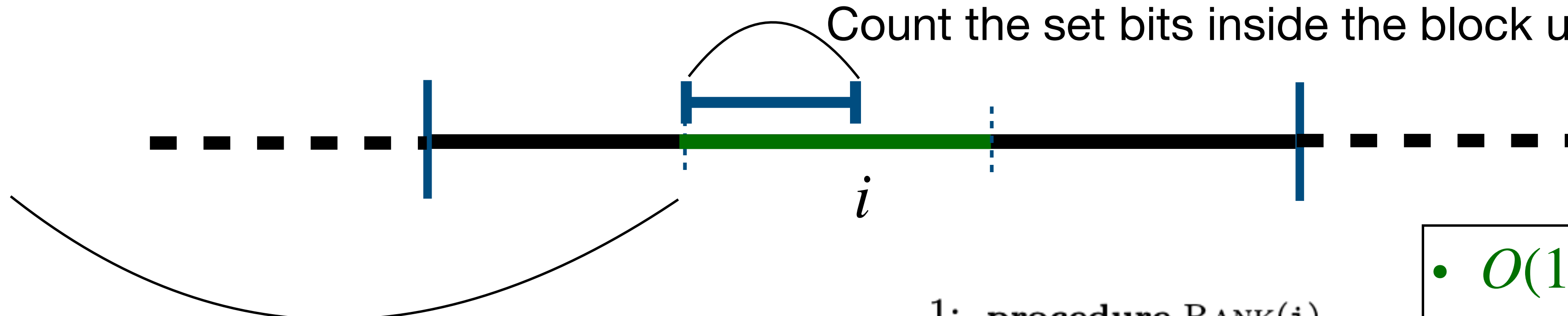
**Total space usage in bits**

$$\frac{n}{d} \cdot \left( \log 2m + \frac{64}{s} \right)$$

- **Nothing is stored for the $n/sd$ super-block rank sequence !**

- **They are _almost_ encoded with the $b$ parameters of the regression lines $ax + b$, when $x = 0$**

# Rank with Adjusted Anchoring

Count the set bits inside the block until position $i$ via popcount.

$i$

- **Predict** the count of the set bits until the beginning of the block by using the corresponding regression line
- **Adjust** predicted value by using the corresponding correction value.

- If the error is larger than the adjustable range, then **scan towards left** until hitting a correctly predictable block.

- $O(1)$-time, without scan
- $O(s)$-time, if scan is required

```
1: procedure RANK(i)
2:     bID ← ⌊(i + d − 1)/d⌋
3:     excess ← popcnt(B[i + 1 ... bID · d])
4:     rnk ← 0
5:     while (V[bID]=1)&(R[bID]=0)&(bID>0) do
6:         rnk ← rnk+popcnt(B[(bID−1)·d+1..bID·d])
7:         bID ← bID − 1
8:     if (bID > 0) then
9:         sID ← ⌊(bID + s − 1)/s⌋
10:        innerbID ← i − (bID − 1) · d
11:        rnk ← rnk + ⌊β_0^{sID} + β_1^{sID} · innerbID⌋
12:        rnk ← rnk − m/2 + R[bID]
13:        if V[bID] = 1 then
14:            if R[bID] ≥ m/2 then
15:                rnk ← rnk + m/2
16:            else
17:                rnk ← rnk − m/2
    return (rnk − excess)
```

# Select with Adjusted Anchoring

```
 1: procedure SELECT(i)
 2:     q ← ⌊i/setBitRatio⌋
 3:     sID ← ⌊(q + s·d − 1)/(s·d)⌋
 4:     while (β₀^sID ≥ i)&(sID > 0) do sID ← sID−1
 5:     while (β₀^sID < i) do sID ← sID+1
 6:     innerbID ← ⌊(i−β₀^sID)/β₁^sID⌋
 7:     if innerbID < 1 then innerbID = 1
 8:     bID ← innerbID + sID · s
 9:     if bID > sbc then
10:         bID←sbc
11:         innerbID←(bID mod s) + 1
12:     rank ← RANK(bID · d)
13:     c ← popcnt(B[(bID−1)·d+1..bID·d])
14:     while rank < i do
15:         bID ← bID + 1
16:         c ← popcnt(B[(bID−1)·d+1..bID·d])
17:         rank ← rank + c
18:     rank ← rank − c
19:     while rank ≥ i do
20:         bID ← bID − 1
21:         c ← popcnt(B[(bID−1)·d+1..bID·d])
22:         rank ← rank − c
23:     p ← KTHSETBIT_SIMD(i − rank)
24:     return (bID−1)·d + p
```

- Start with a rough **prediction** of the super-block according to *average set bit ratio*
- Perform a **linear scan** to locate it explicitly

- By using the corresponding linear regression, **predict** the block position inside the super block
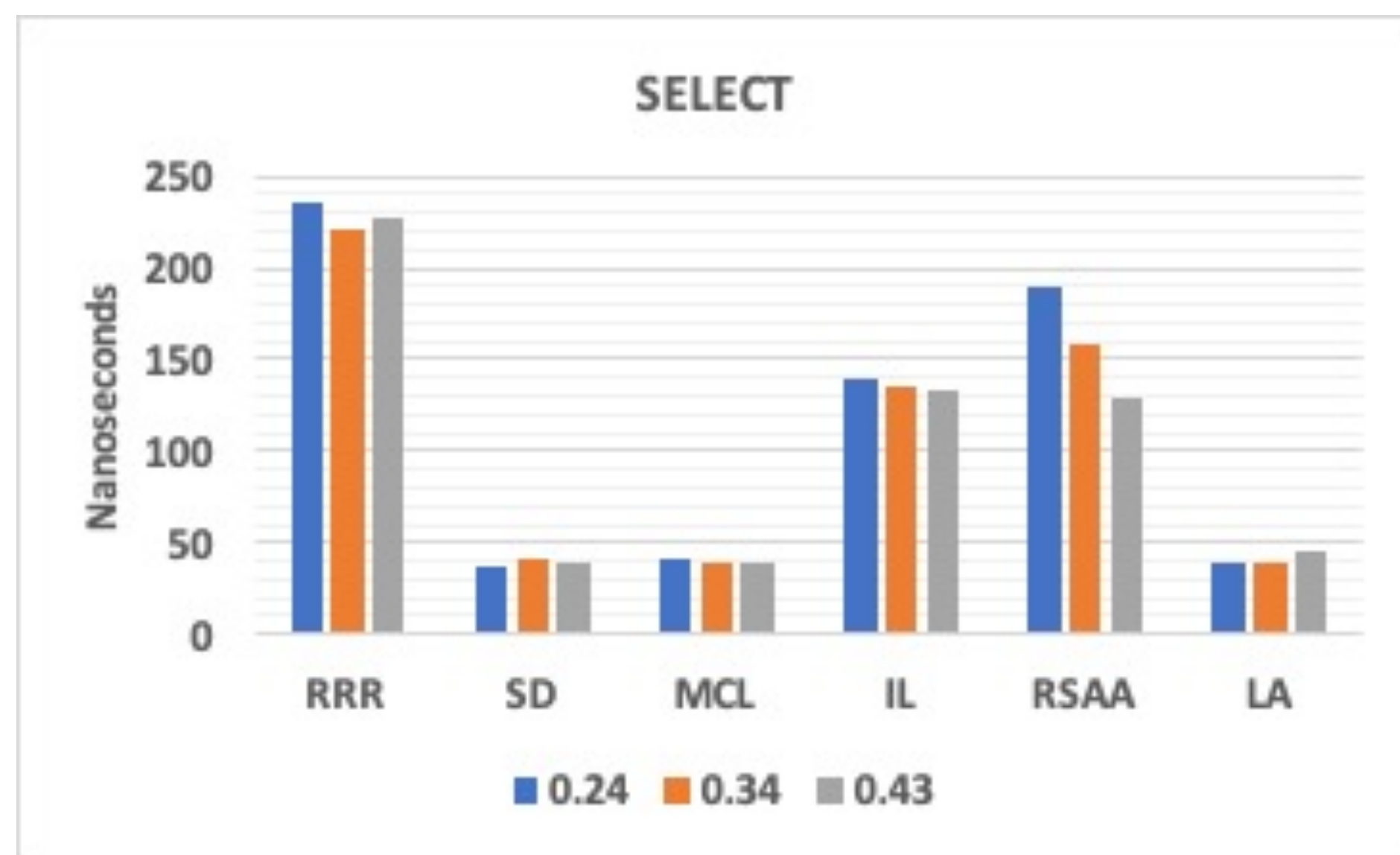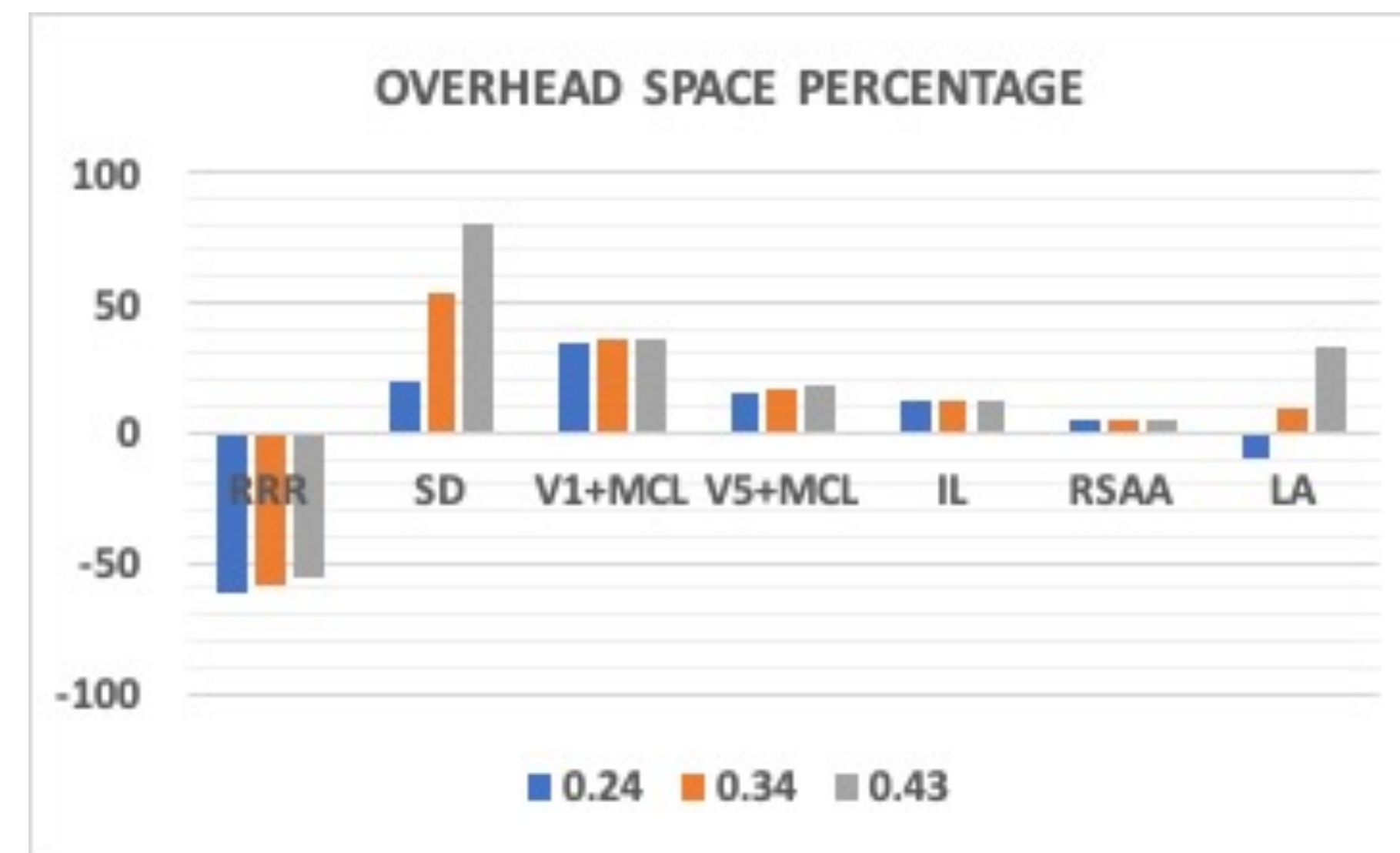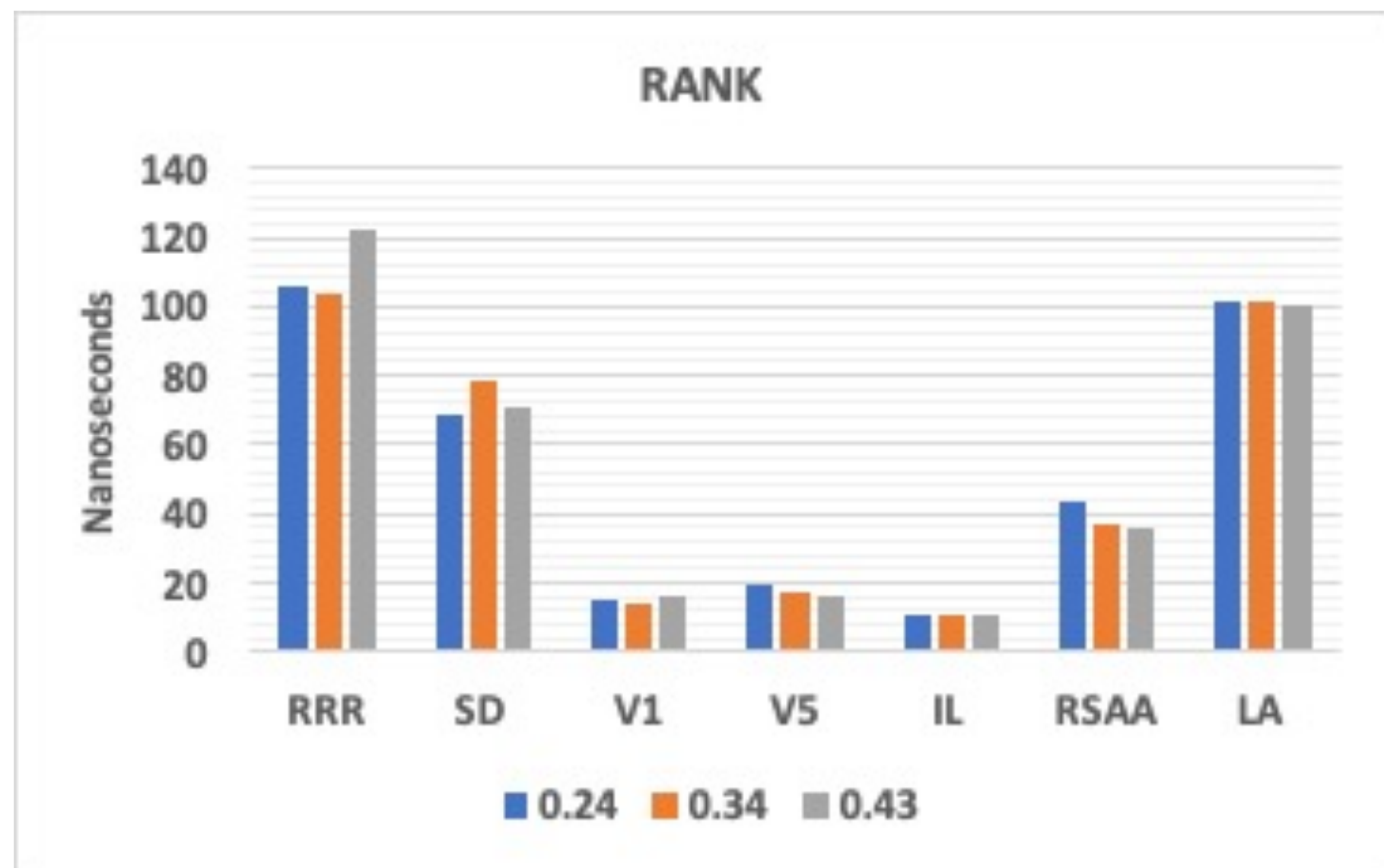
- The predicted block ID is adjusted by checking its rank value
- In case it is needed, again the neighboring blocks are scanned towards left or right until the correct block is located.

- Last but not the least, the kth set bit is determined with **SIMD instructions**.
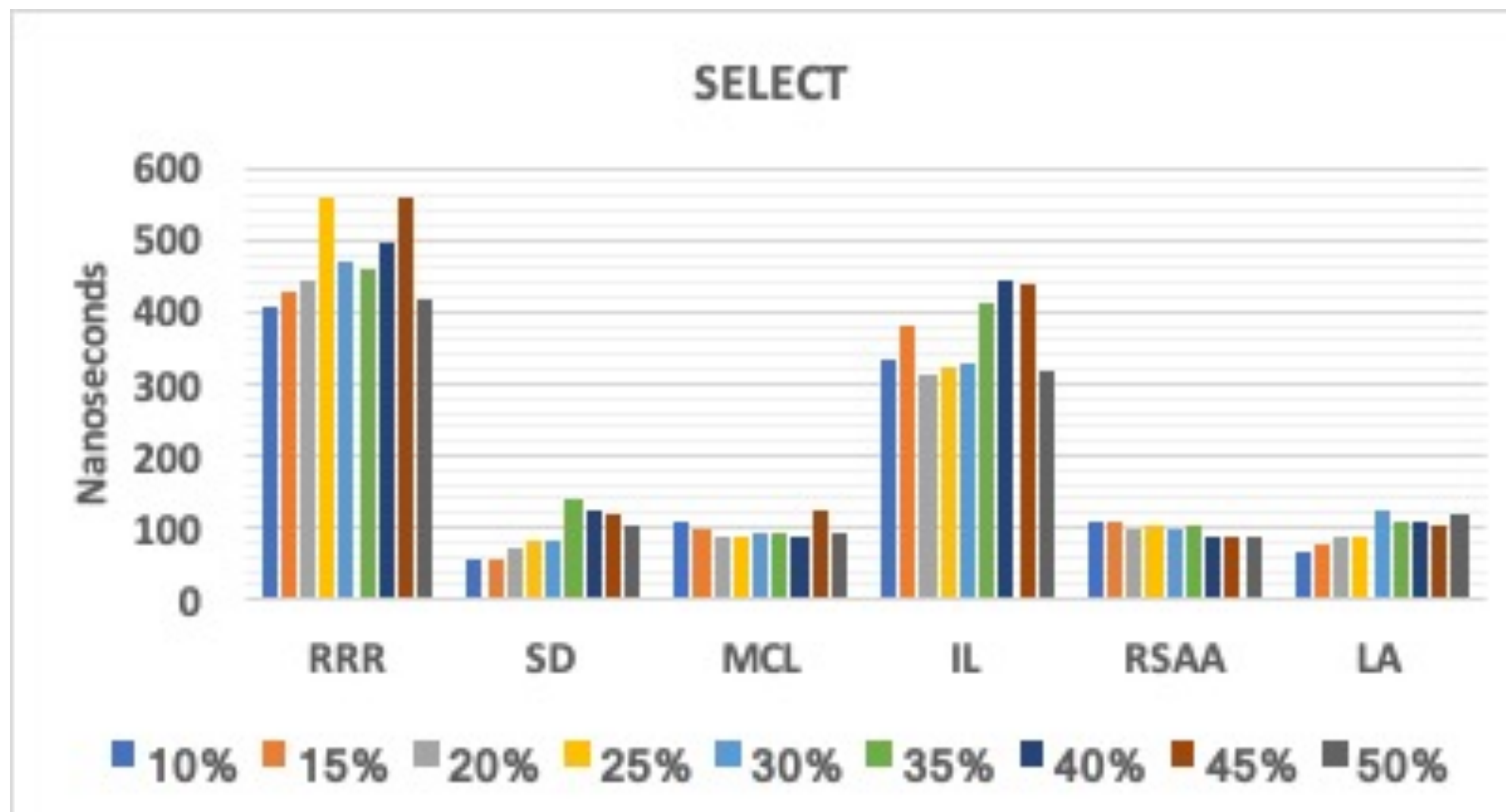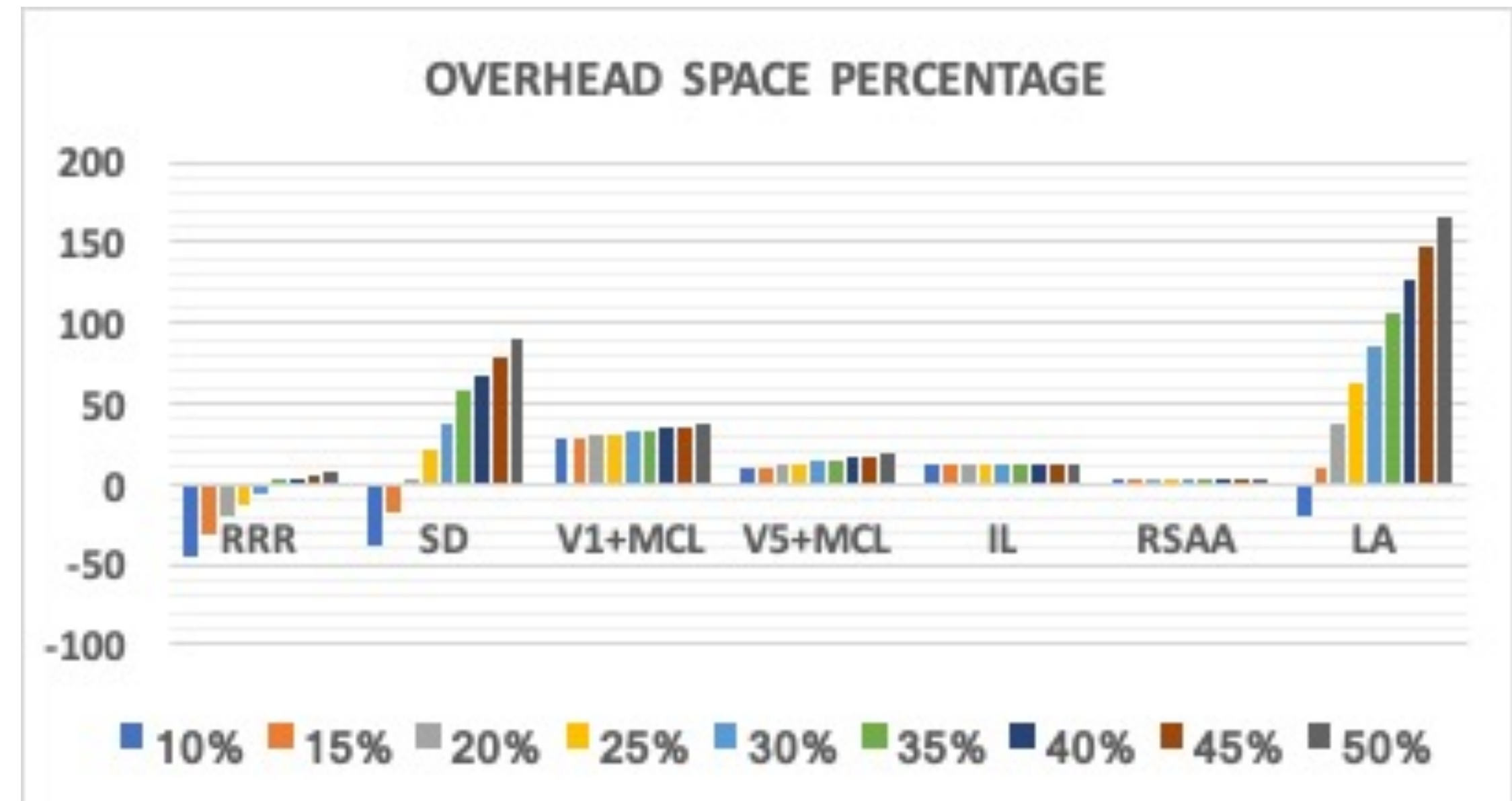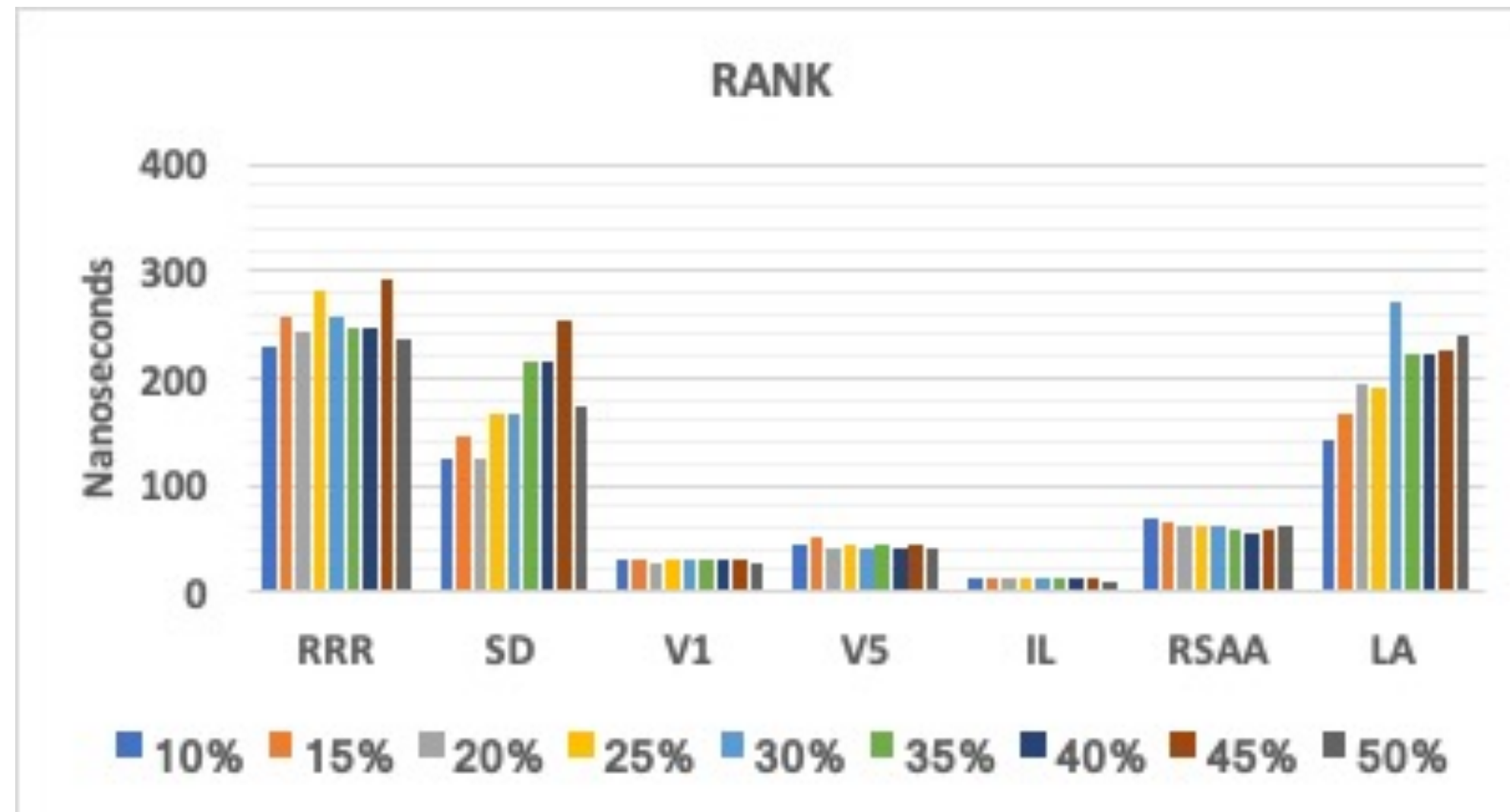
# Tuning the parameters

- **The overhead space is less than 5 percent, and correct prediction of the block rank values is more than 99% .**

- The performance of the RSAA scheme highly depends on accurate predictions, where the parameters $s, d, m$ are central in prediction performance.

- Setting $m = d = 256$, $s = 16$ has been observed to provide most reasonable results empirically.

- Larger $m$, which denotes the recoverable error threshold, results in improved prediction success, but increases space usage (*and vice versa for sure*). For fast processing with small space consumption, $m = 256$ has been the best value.

- $d$, the block size in bits, is set to 256 as well to keep space consumption less than $5\%$ while not hurting the prediction performance with $m = 256$.

- To keep the irregular block count less than $1\%$, $s$, the number of blocks in a super block is set to 16.

# Experimental Results on Real-Data



- RRR,SD, LA are compressed, V1,V5,MCL are uncompressed schemes
- Data sets are the real sequences used in previous study, which are averaged according to 0-1 ratios as 24, 34, and 43 percent.
- RSAA has the least overhead space
- LA compressed only the sequences with less than 30 percent set bits!

# Experimental Results on Synthetic-Data



- On randomly generated bit sequences with different densities

- RSAA has the least overhead space

- Randomly distributed, balanced 0-1 distribution favors RSAA.

# Conclusions

- Rank and select on **uncompressed bitmaps** with the machine learning support

- Studies appeared targeting sparse bitmaps previously, where, in contrast, RSAA targets **balanced density** bitmaps

- Overhead around 3-5% of the input

- **Would it be possible to have better time-space resuşts with other learning paradigms ?**

- **More generally, can ML techniques help in basic combinatorial tasks ?**