

CSE 410 – Advanced Data Structures

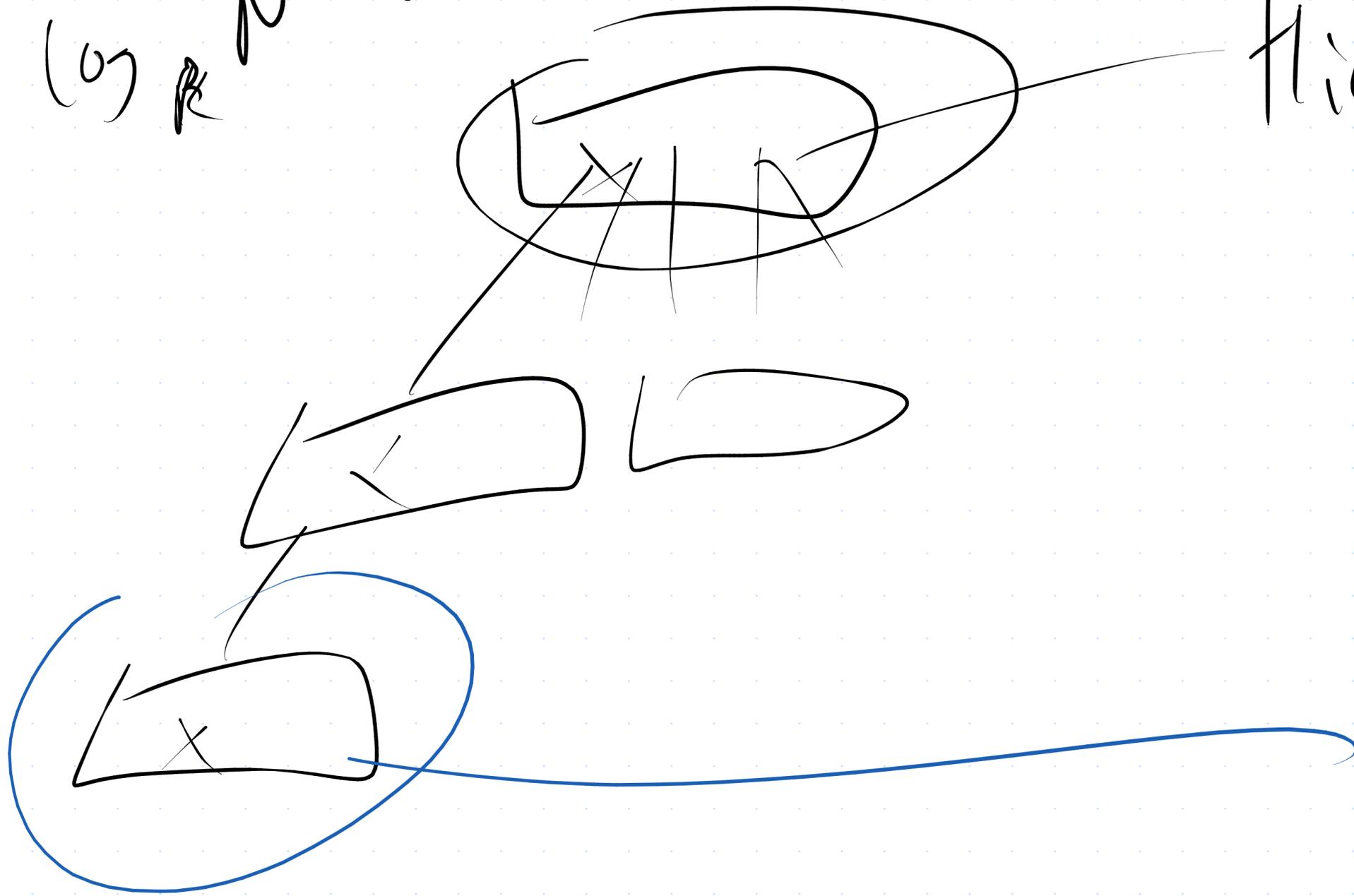
Topic 07: Write-Optimized Structures

Oliver Kennedy

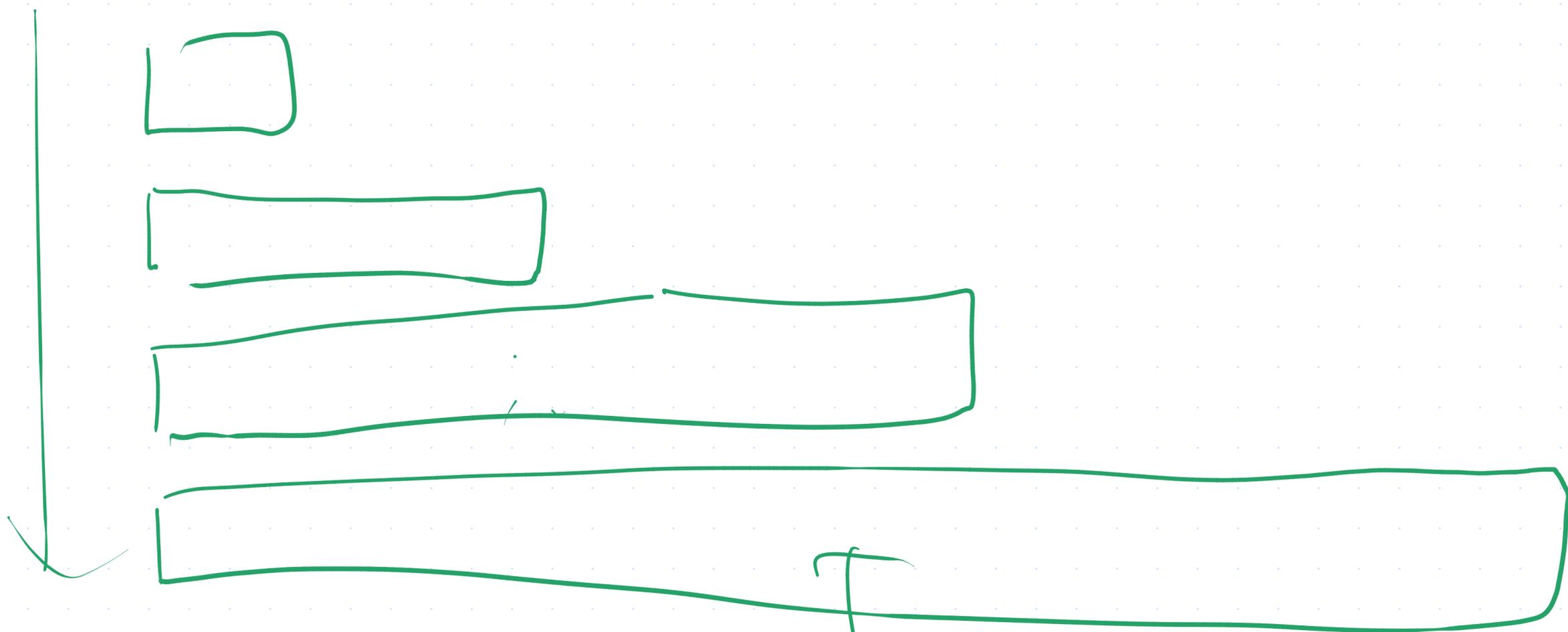


$\log R^N$ writes

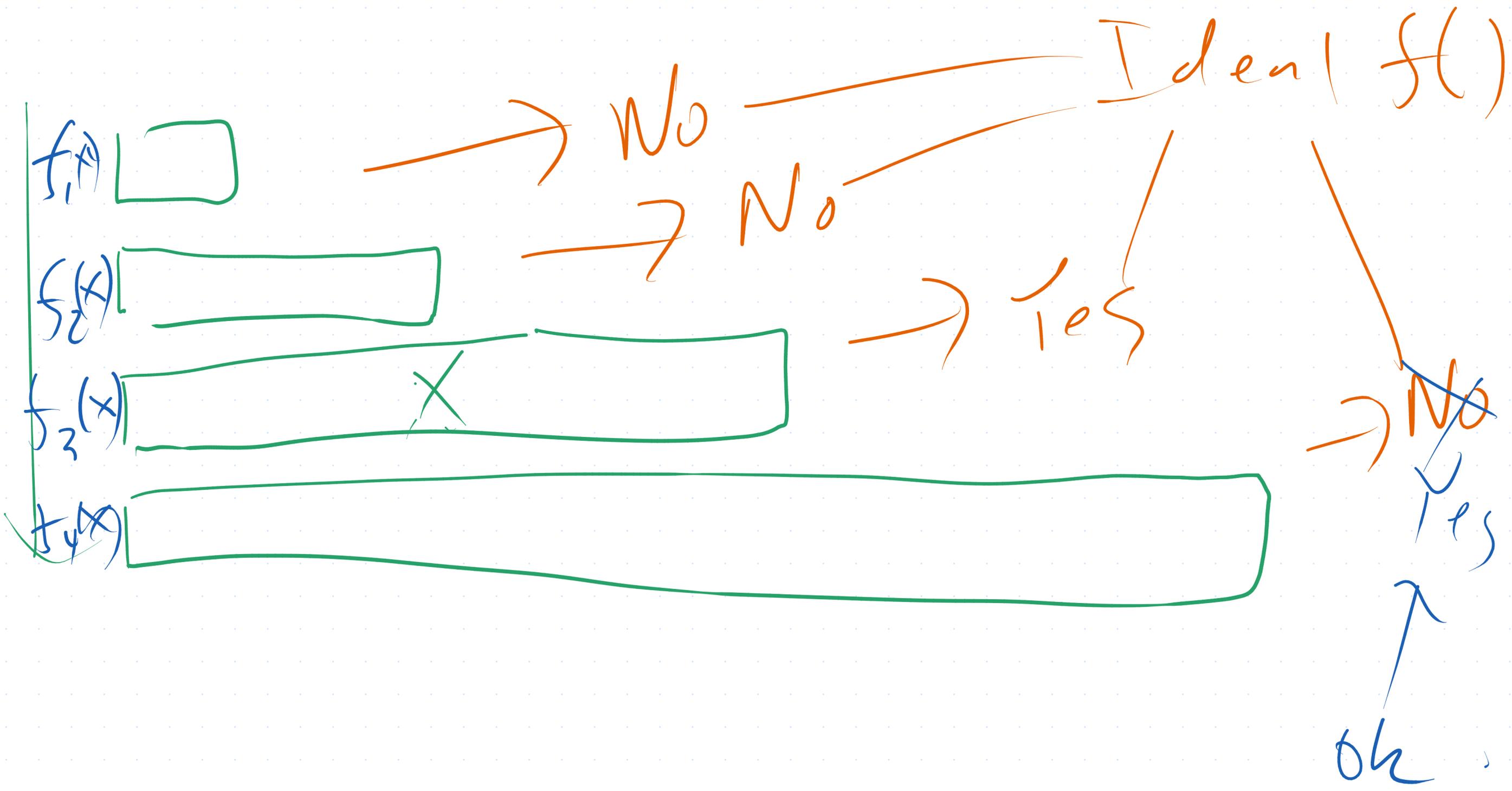
High contention



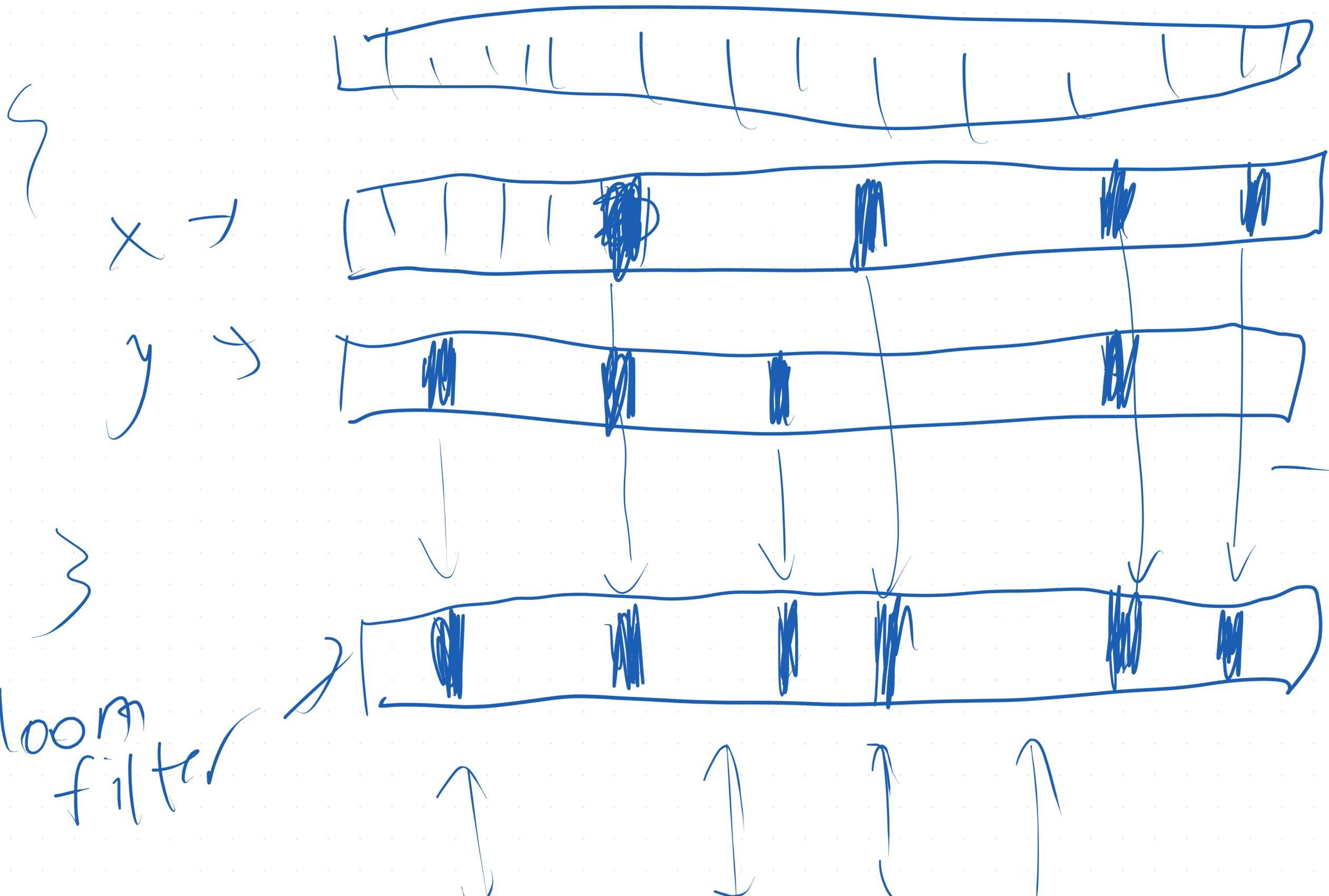
Coordination
required

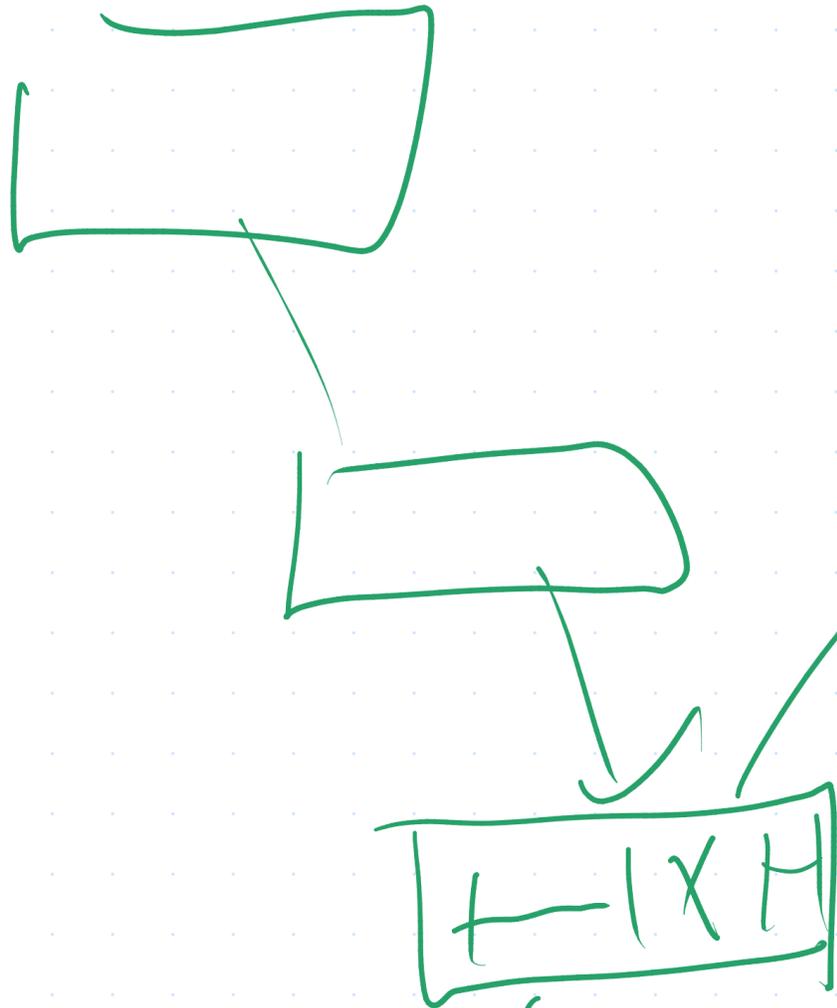


sorted
array
↑



Bloom filter

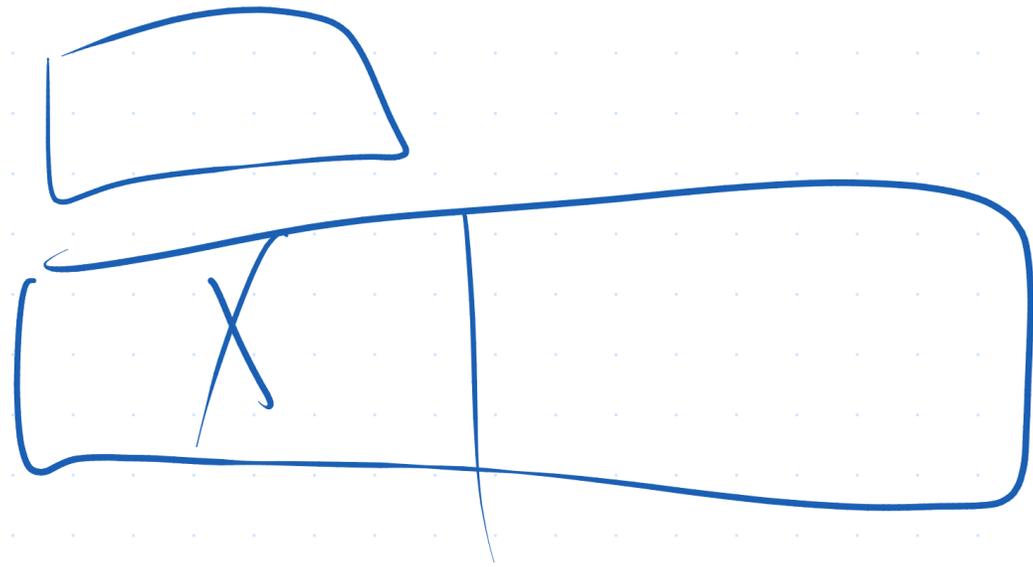




change
one
record

one change
all records
need to be
written
out

Deletions & updates



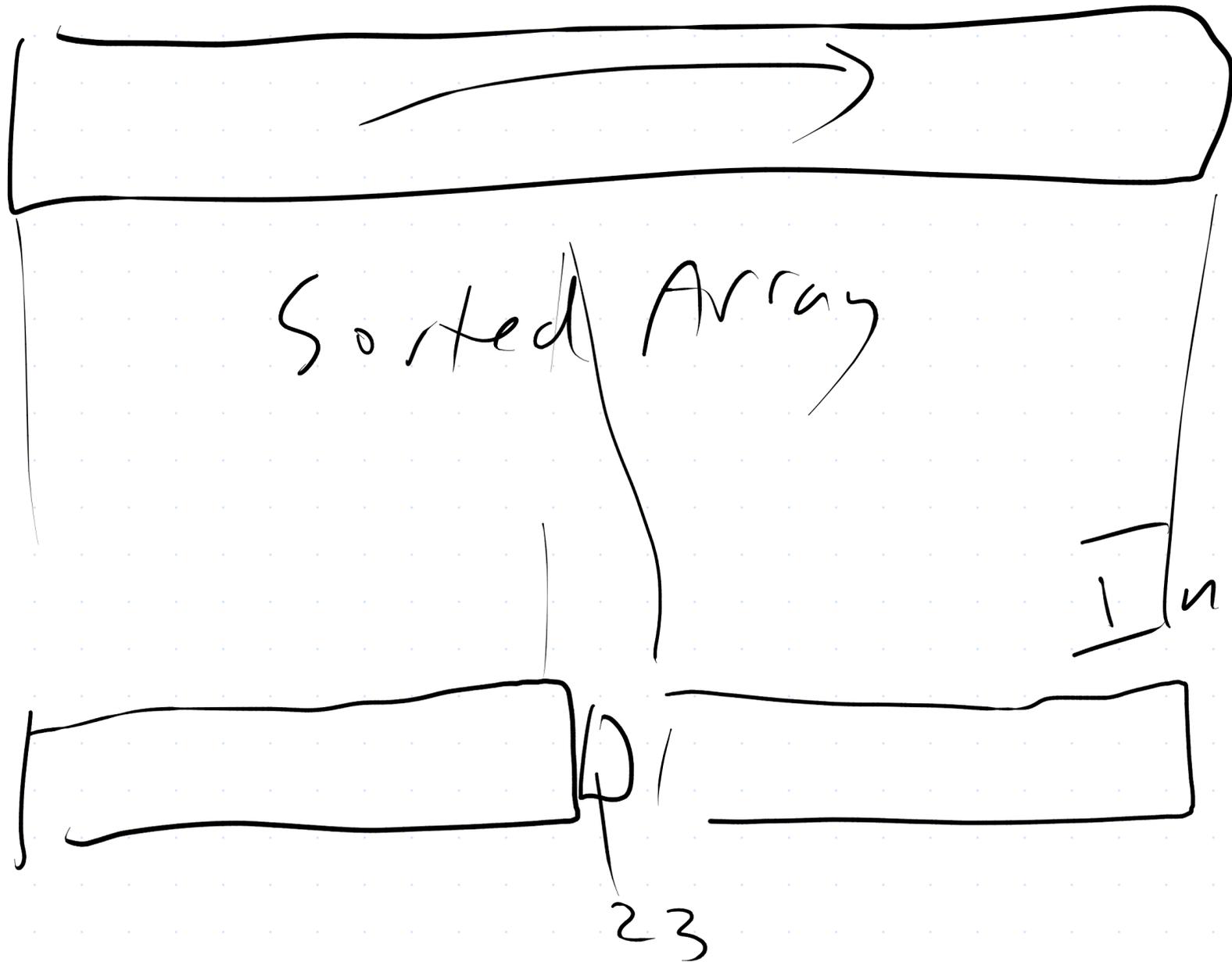
Deletion = update to
Null



exactly one sorted run



Level 4



$O(N)$
IOs
for all inserts

How is a sorted run stored?

- Array (sorted)

- ISAM Index / Btree (original)

- (DF)-Based Index

↳ or Interpolation Search

- Fence pointer table

Idea 1 : store min, max

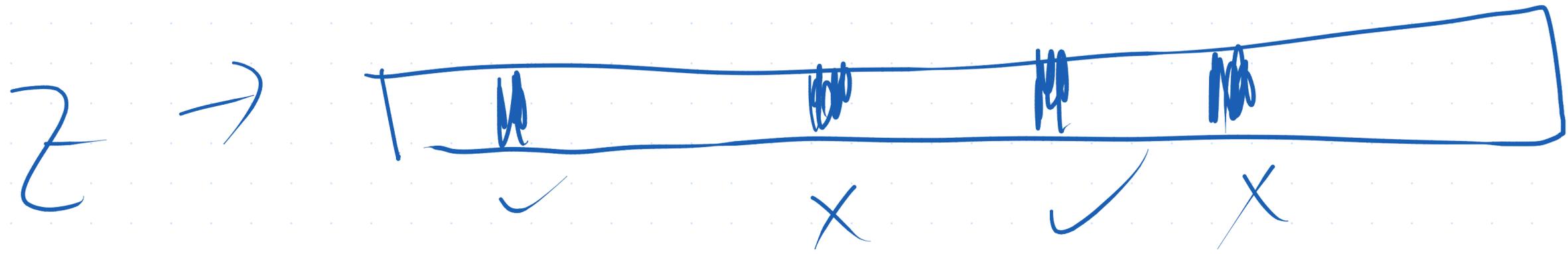
↳ discard requests

that are not in range

Req: min, max

↳ Fence pointer table gives us min
max

↳ No overhead!

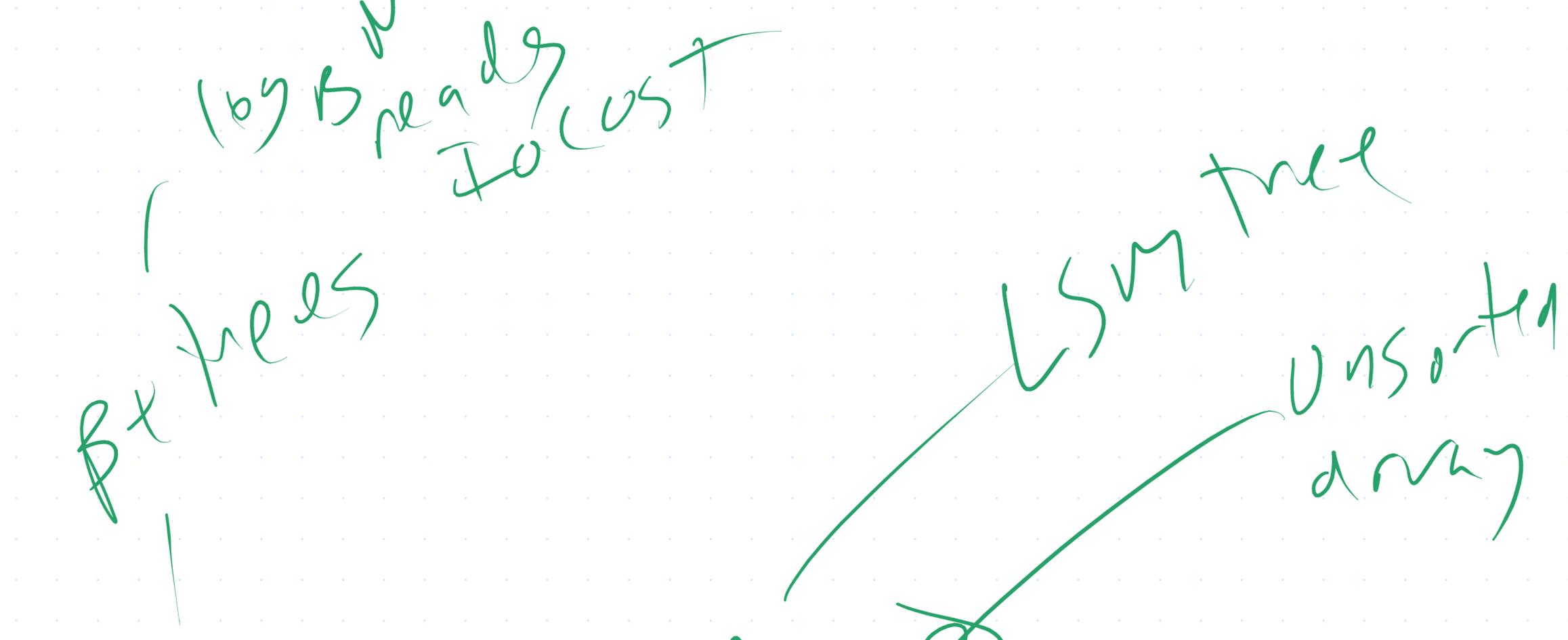


Bloom filter parameters

→ size of the bit vector

→ # of bits that each element sets

→ size of bit vector & number of elements stored



read optimized

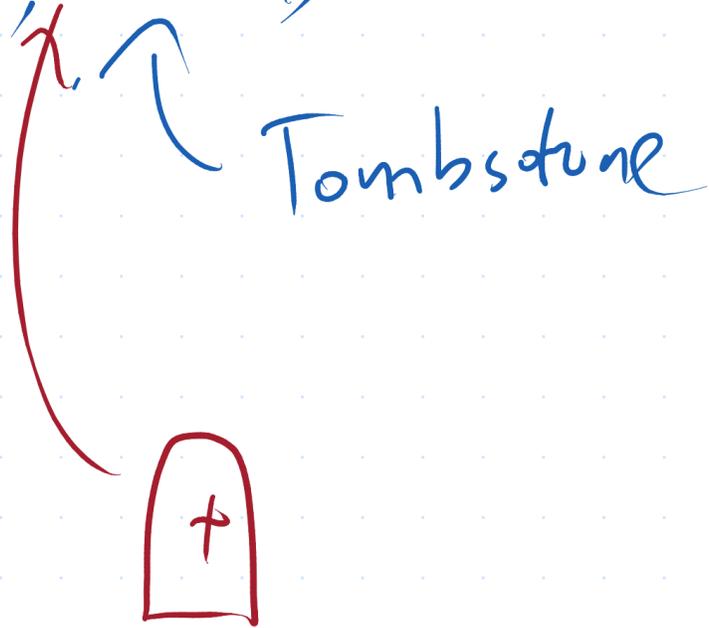
write optimized

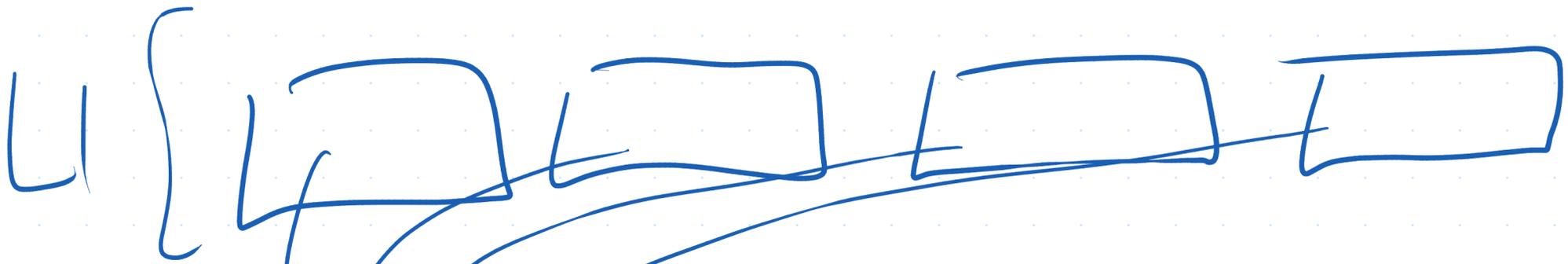
$\log^2(W)$

goal
 same read complexity
 AS B+ tree
 better
 read cost \uparrow

Delete (k) = Update (k, NULL)

So how do we implement
Update (k, 42)







23 48 =

← B →

Sorted Array

$$O(N)$$

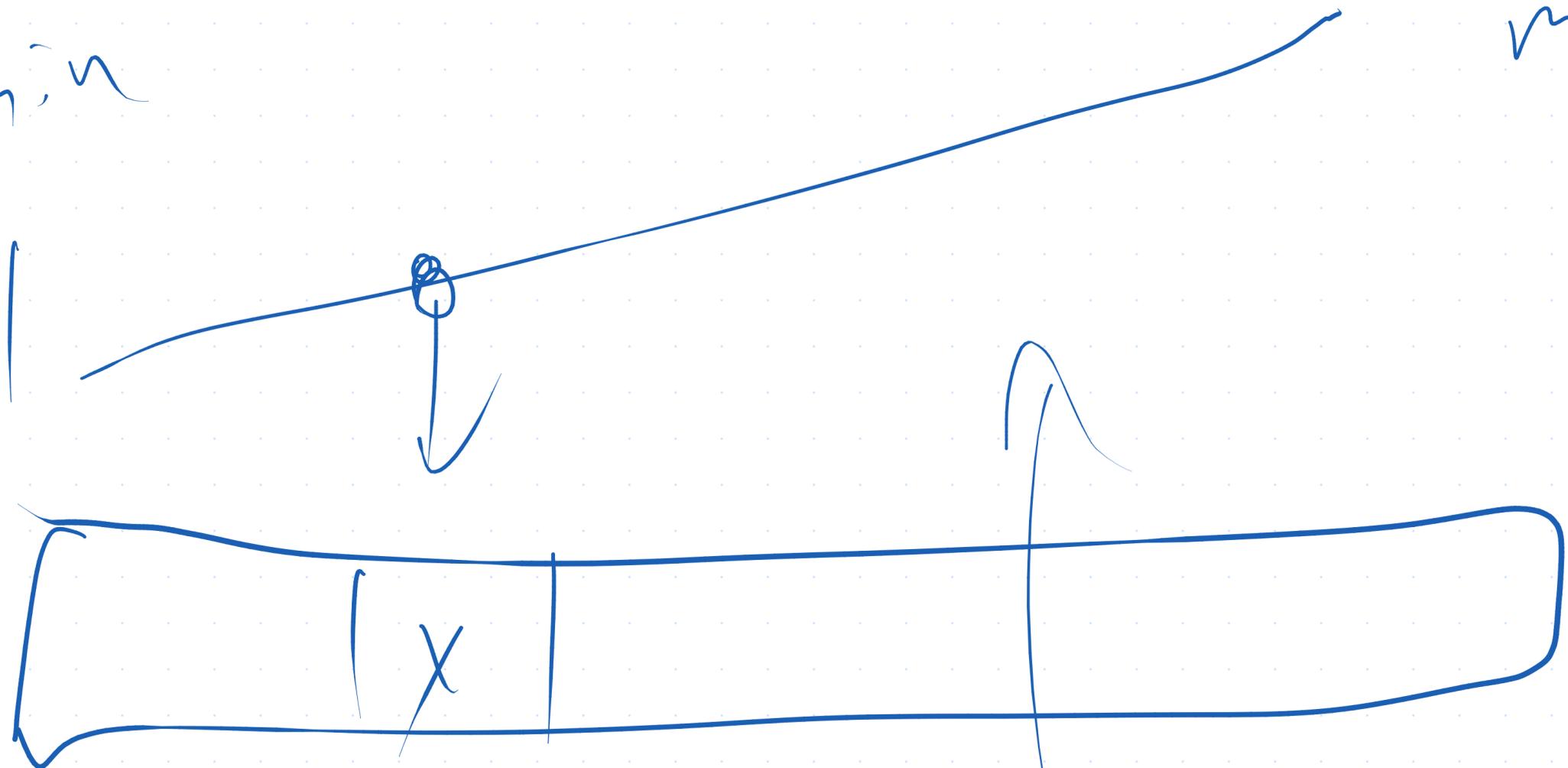
every B insertions

vs
 $O(N)$ per insertion

$$O\left(\frac{N}{B}\right)$$

min

max



CDF-Indexes
Learned Indexes

$$f_i(x) \Rightarrow \begin{cases} \text{True if } x \text{ is a key} \\ \text{on that level} \\ \text{false otherwise} \end{cases}$$

↑
Ideal

$x \in$ the set of keys
at that level

or
↓

$$f_i(x) = \begin{cases} \text{false if } x \notin \text{set of keys} \\ \text{true if } x \text{ can be in the set} \\ \text{of keys} \end{cases}$$

More bits per element

↳ More chances to avoid collisions on test

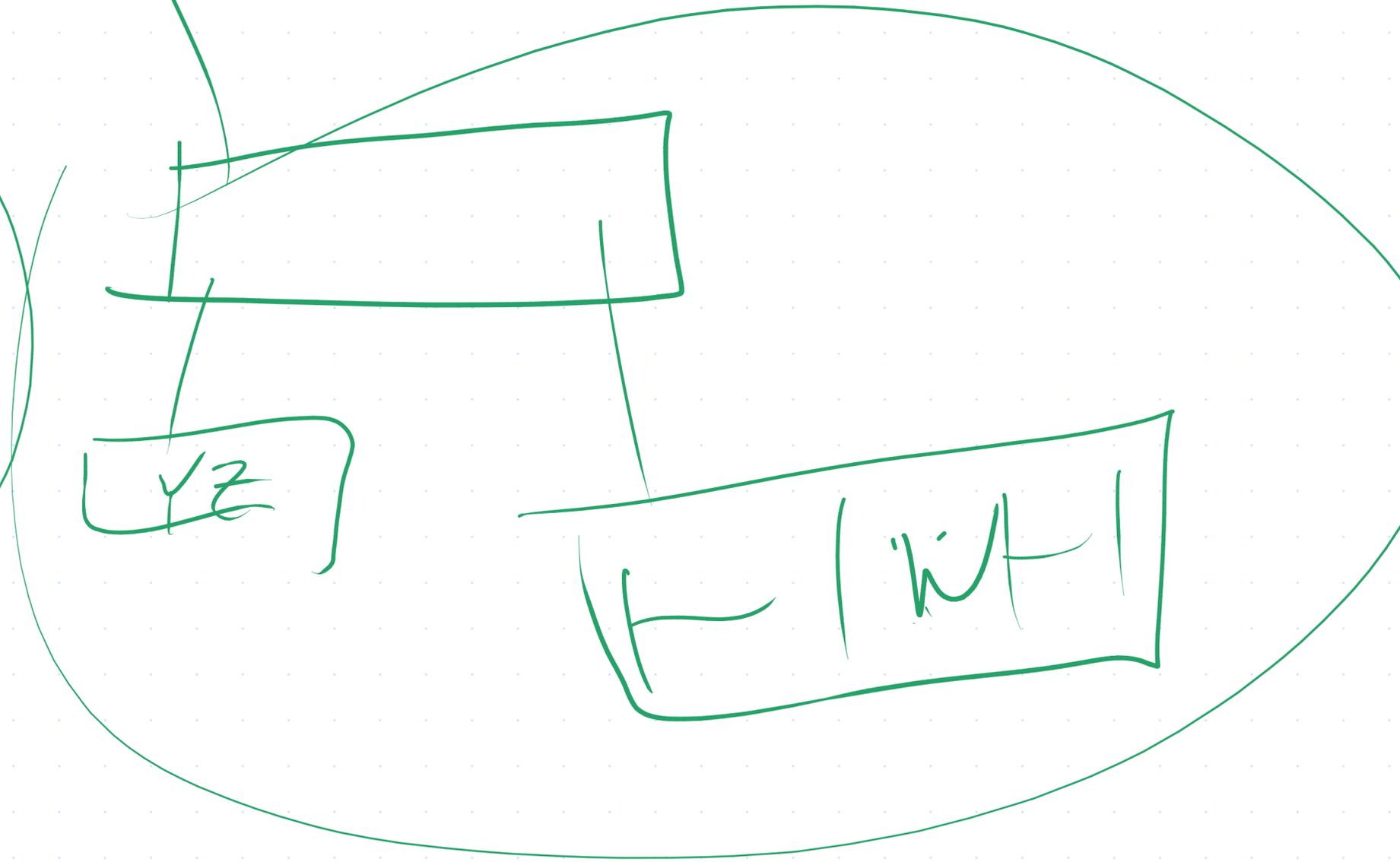
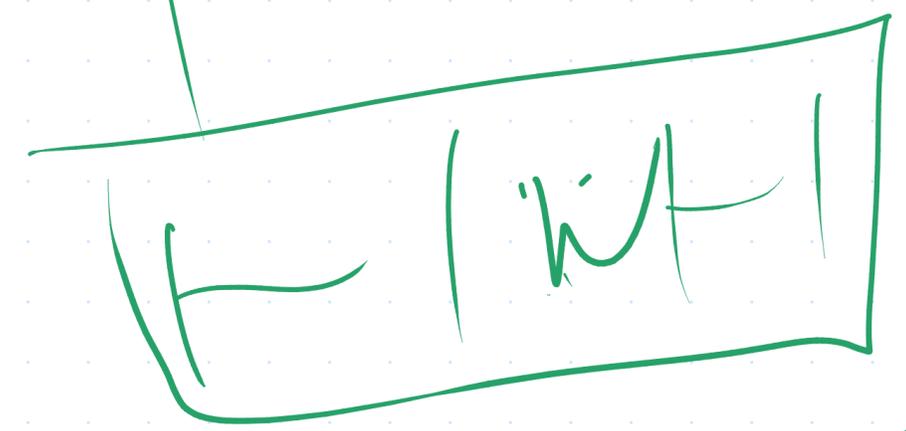
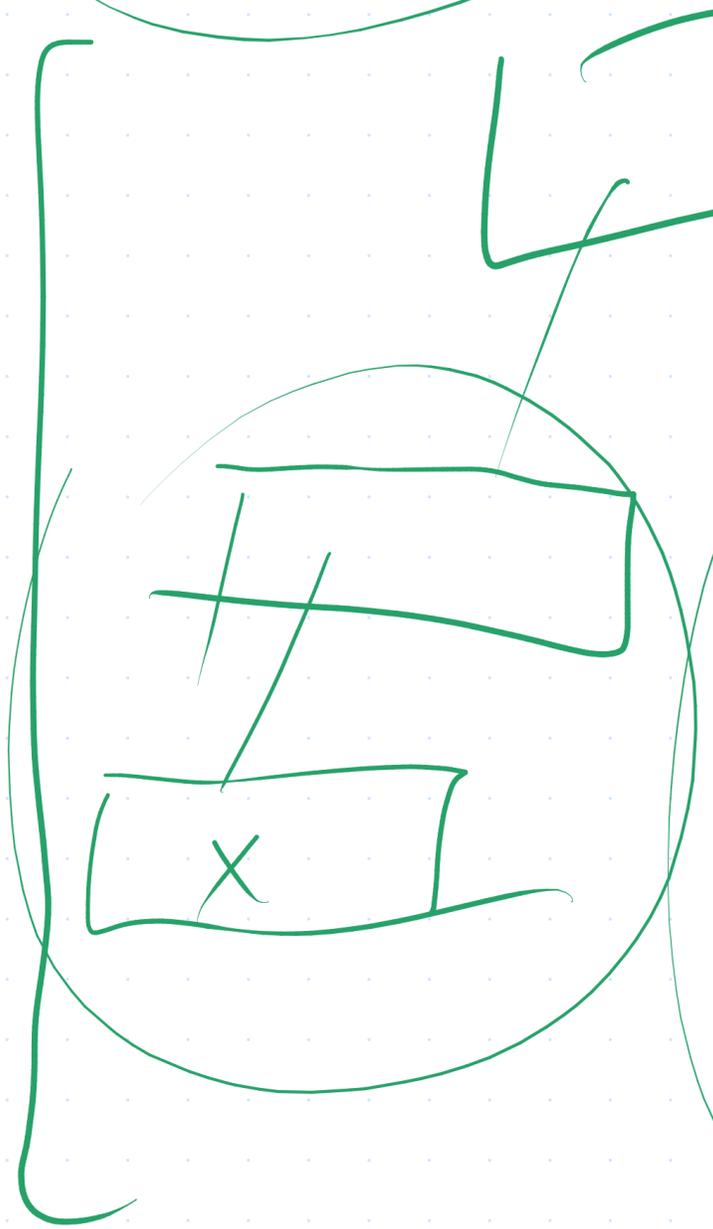
↳ More bits get set, so more collisions

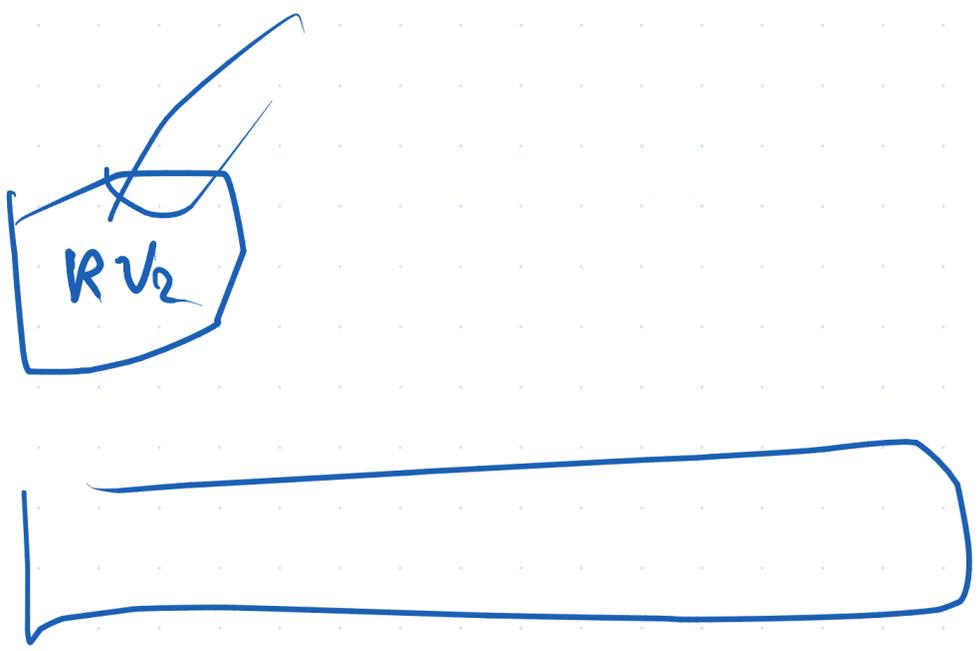
but with performance

mem



disk



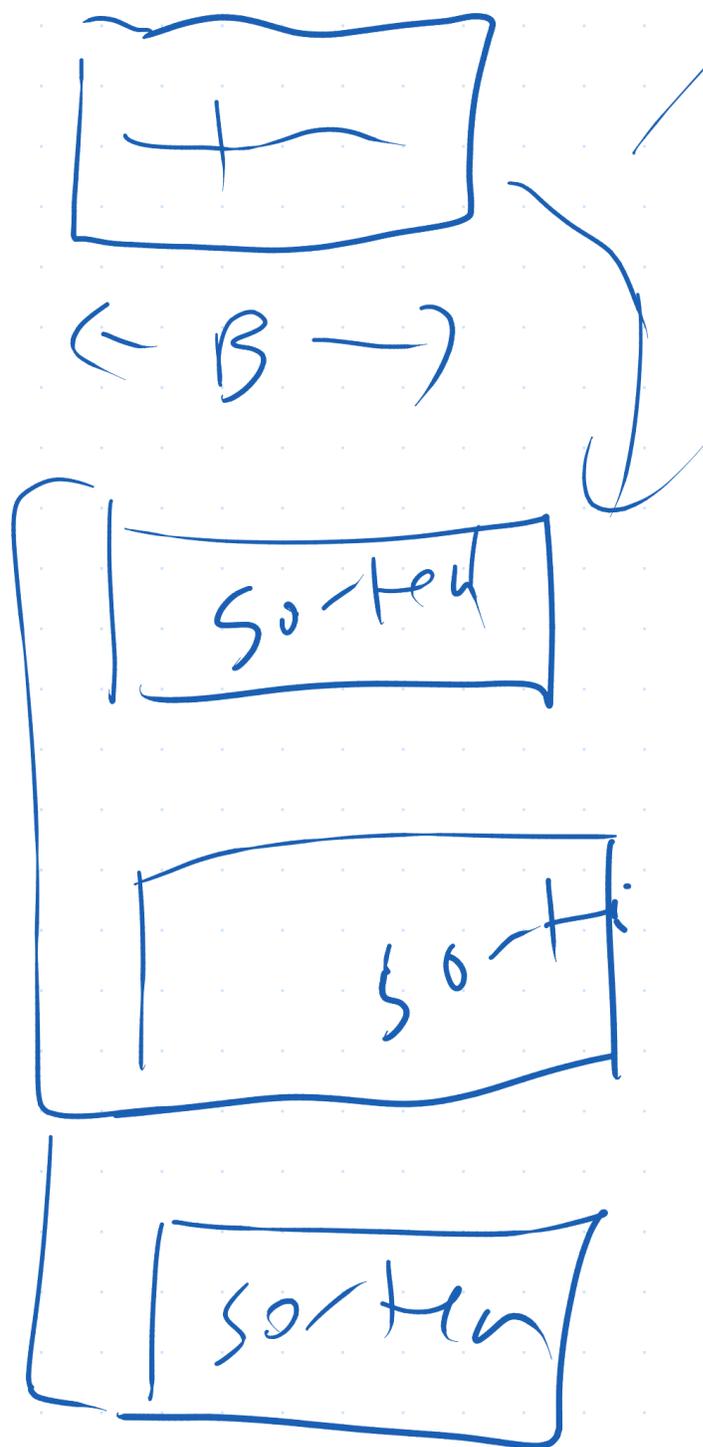


get(e)
returns
v_e

(more recent
value)

other data
on same page

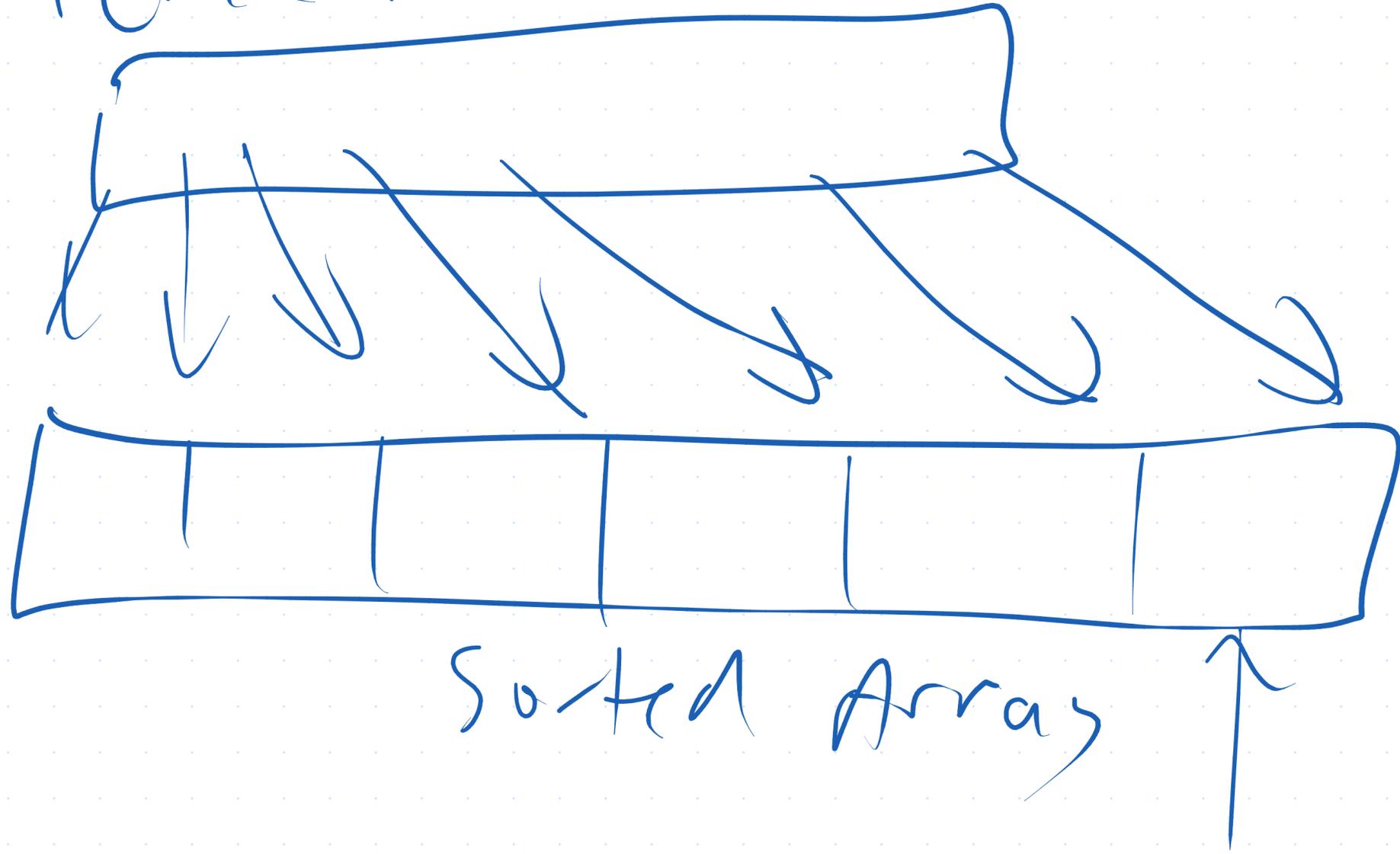
write sorted Buffer



$O(B)$ per B insertion
 $O(1)$ amortized write

$O\left(\frac{N}{B}\right)$ sorted arrays
 $O\left(\frac{N}{B}\right)$ I/Os per read

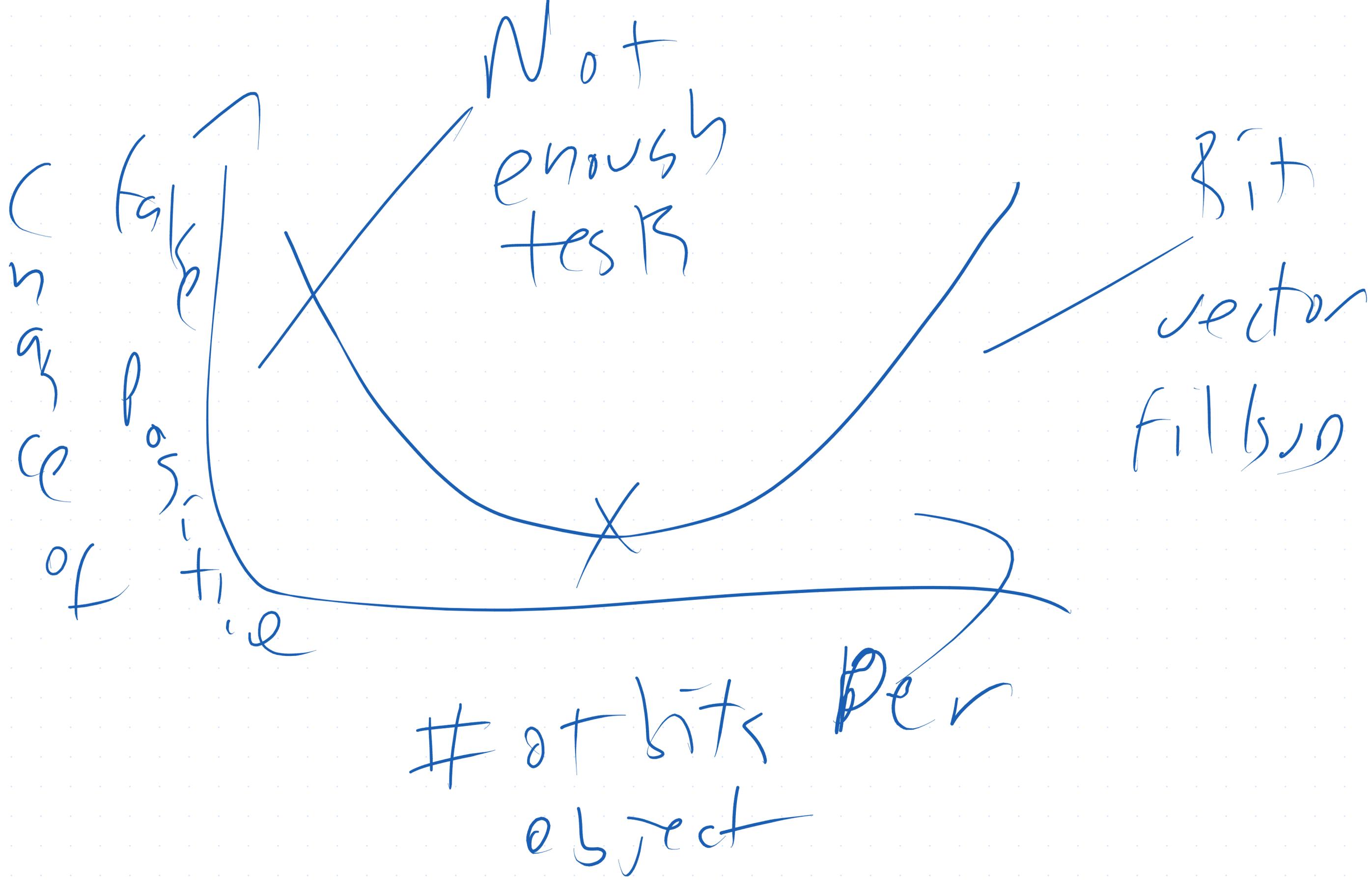
Fence pointer Table



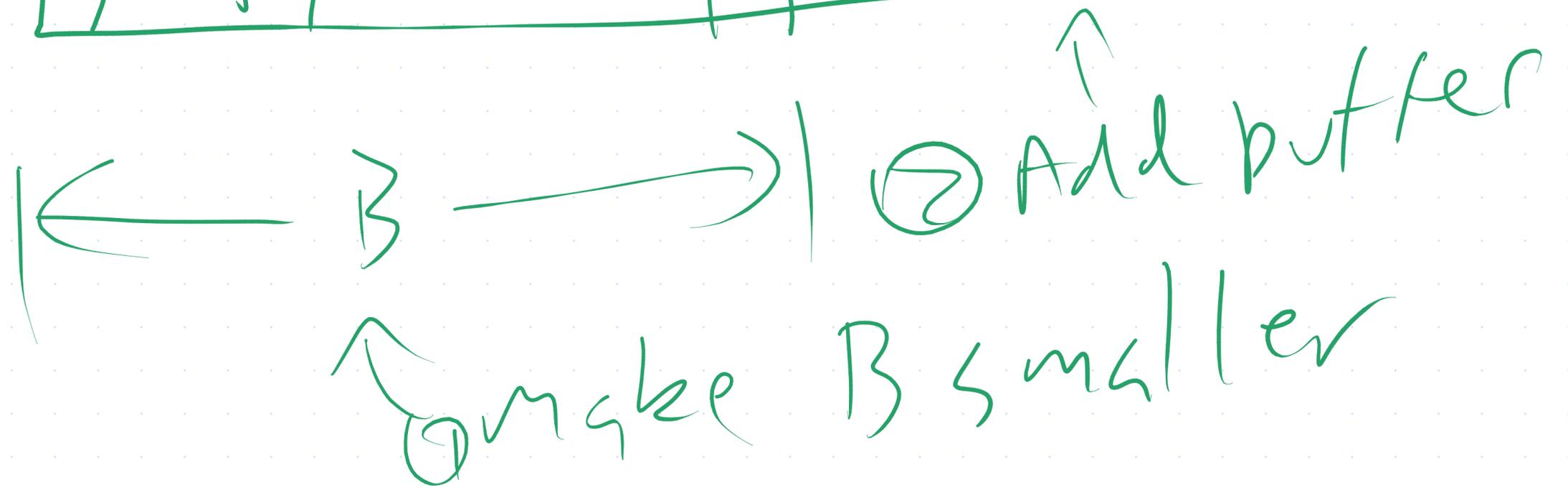
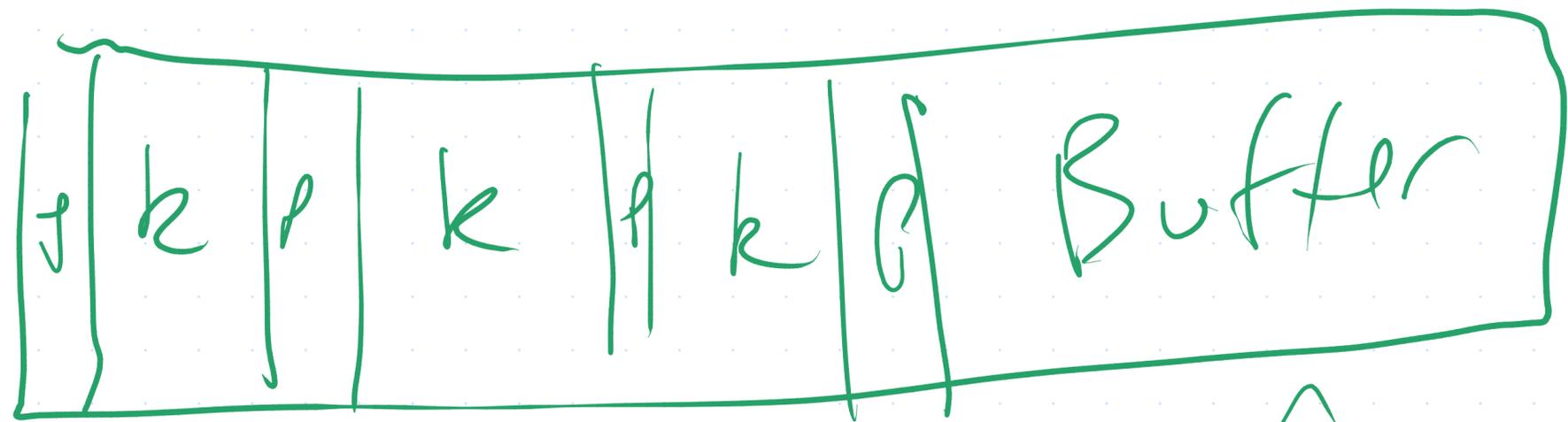
Approximate Set

$$f_0(x) = \text{true}$$

} technically
correct

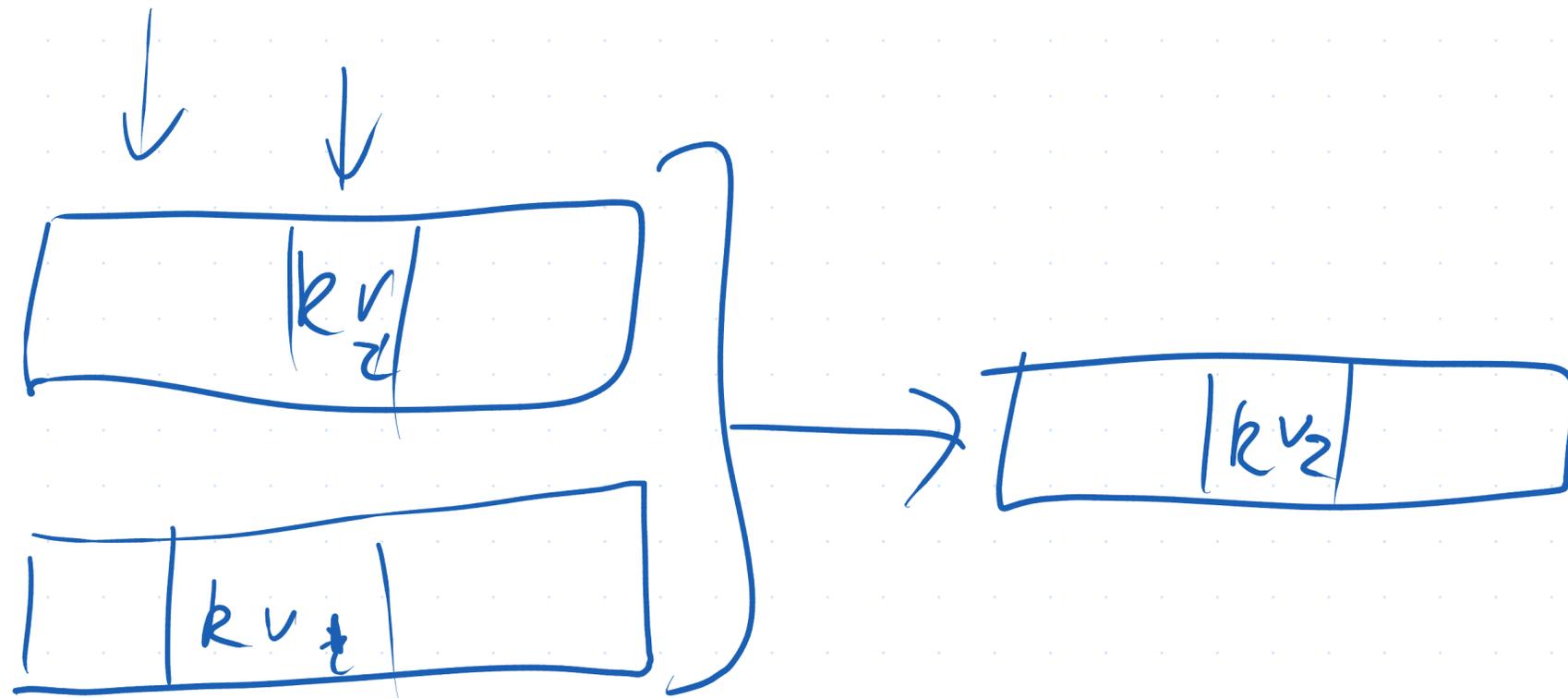


Directory Page



①





Merging "merges" values
for keys

$v_1 < v_2 \rightarrow v_2$ wins

$v_1 = \text{NULL} \rightarrow$ key skipped

Variant 1

write
fill buffer

sort buffer

merge with sorted
array on disk } $O(N)$
IOs

read

Access 1, sorted array } $O(\log N)$

Variant 2

write

fill buffer

sort buffer

write buffer to
disk } $O(B)$

read

- Access $\frac{N}{B}$ sorted array

- $O\left(\frac{N}{B} \log B\right)$

S

```
{  
  bool a;  
  bool b;  
}
```

on insert(x)

```
if  $h(x) \bmod 2 = 0$   
  a = true  
else  
  b = true
```

on lookup(x)

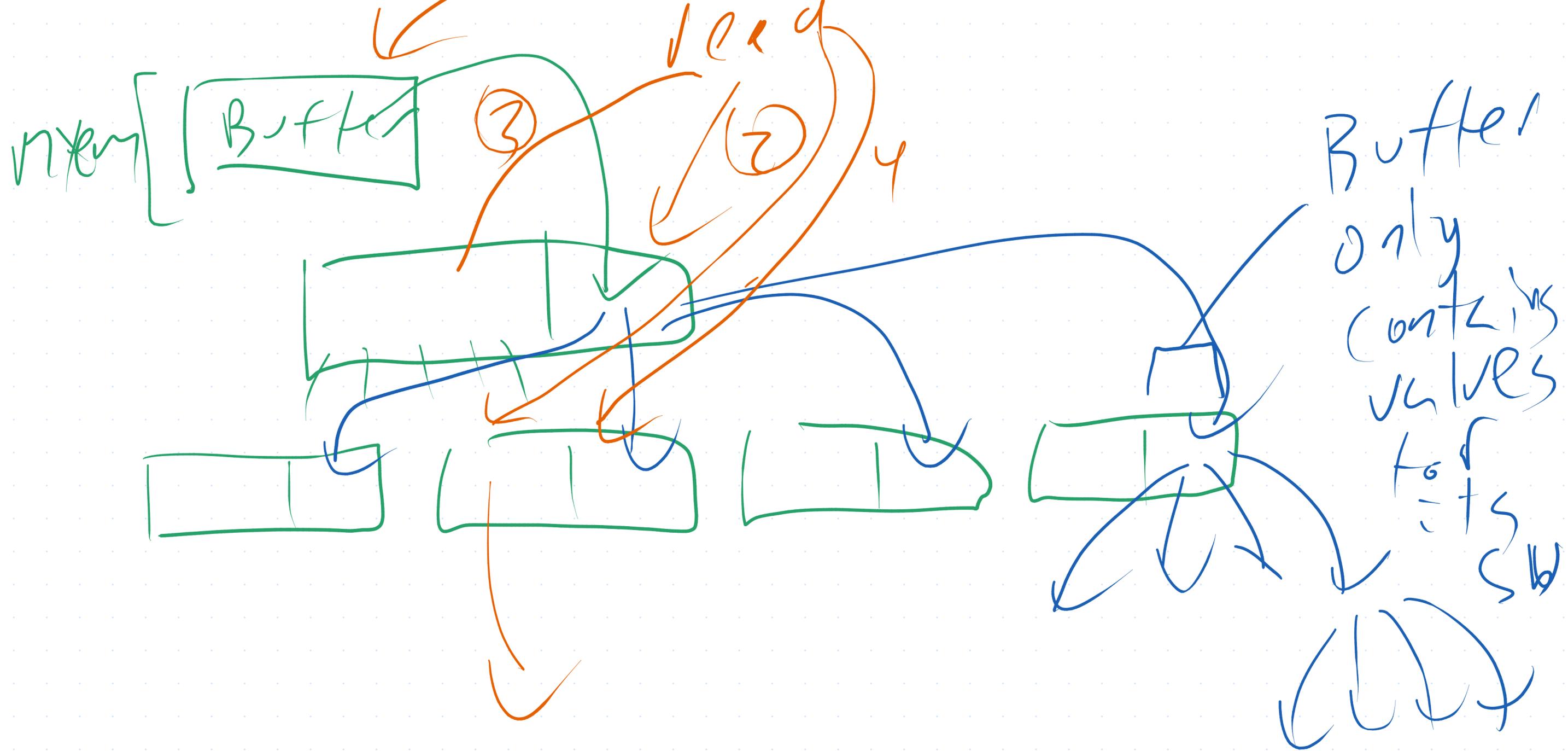
```
if  $h(x) \bmod 2 = 0$   
  return a  
else  
  return b
```

N objects

B size of best vector

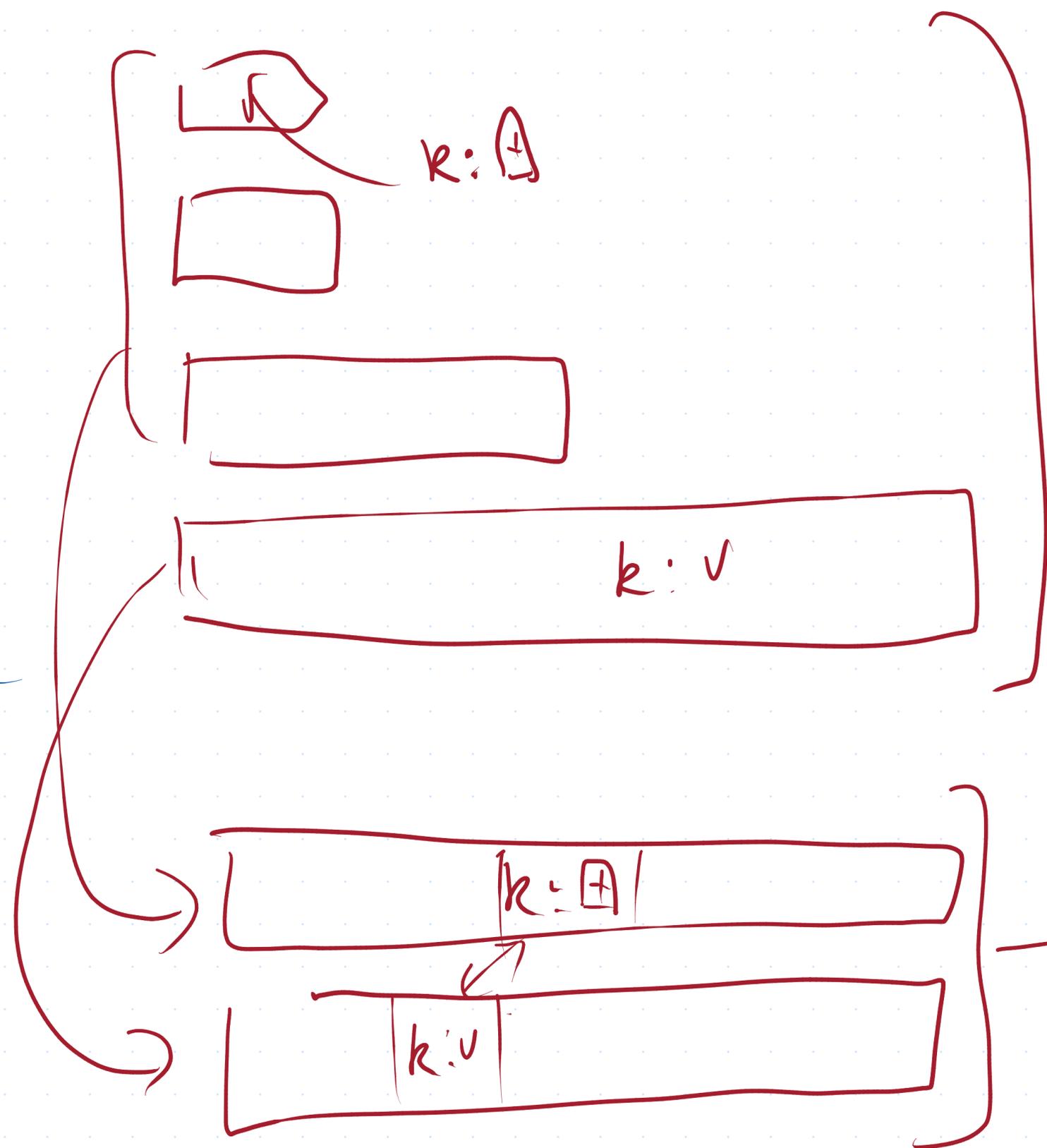
optimal # of bits per object

$$f\left(\frac{N}{B}\right)$$



$\beta - \epsilon$
 Beta-Epsilon Tree

btree



$k: A$

more recent value

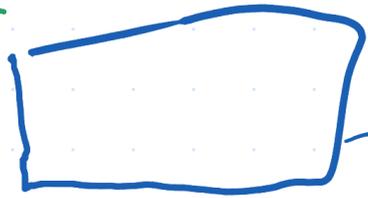


No key if last layer

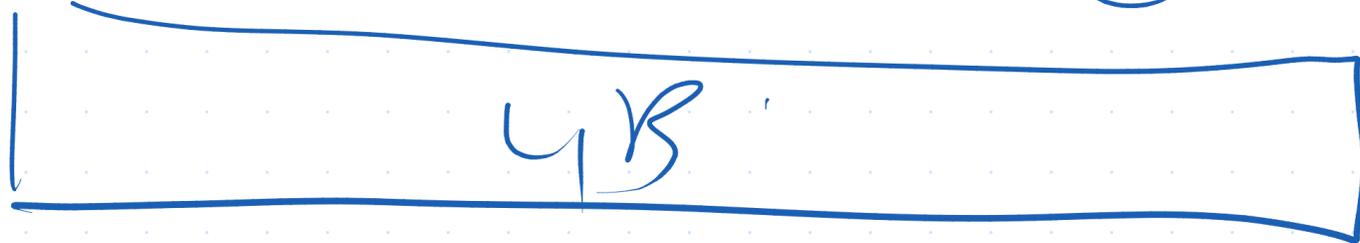
Variant 3

Log Structured Merge Tree
Index

Mem



$\leftarrow B \rightarrow$



Layer 1

Layer 2

Layer 3

disk

{
bool bits[B] = [false]
}

on insert (x)
bits[h(x) % B] = true

on lookup (x)
return bits[h(x) % B]

→ insert (42)
h(42) = 3



↑ Bit 3 must
be 1 after
insert(42)

Size of bit vector

is bits per object
→ % of rate

$O(N)$ bits?
→ bit very small
constant

Analysis

layer 1

B

layer 2

$2B$

layer 3

$4B$

layer 4

$8B$

layer i

$2^{i-1}B$

$$\sum_{i=1}^k 2^{i-1}B$$

$$= O(2^k B)$$

$$N = 2^k B$$

$$\frac{N}{B} = 2^k$$

$$k = \log\left(\frac{N}{B}\right) \curvearrowright$$

to store N
records
 $\log\left(\frac{N}{B}\right)$ layers

{ bits [R][B] = [false]

} $h_i(x)$ for $i \in [0, R)$

on insert(x)

for $i \in [0, R)$

bits[i][$h_i(x) \% B$]

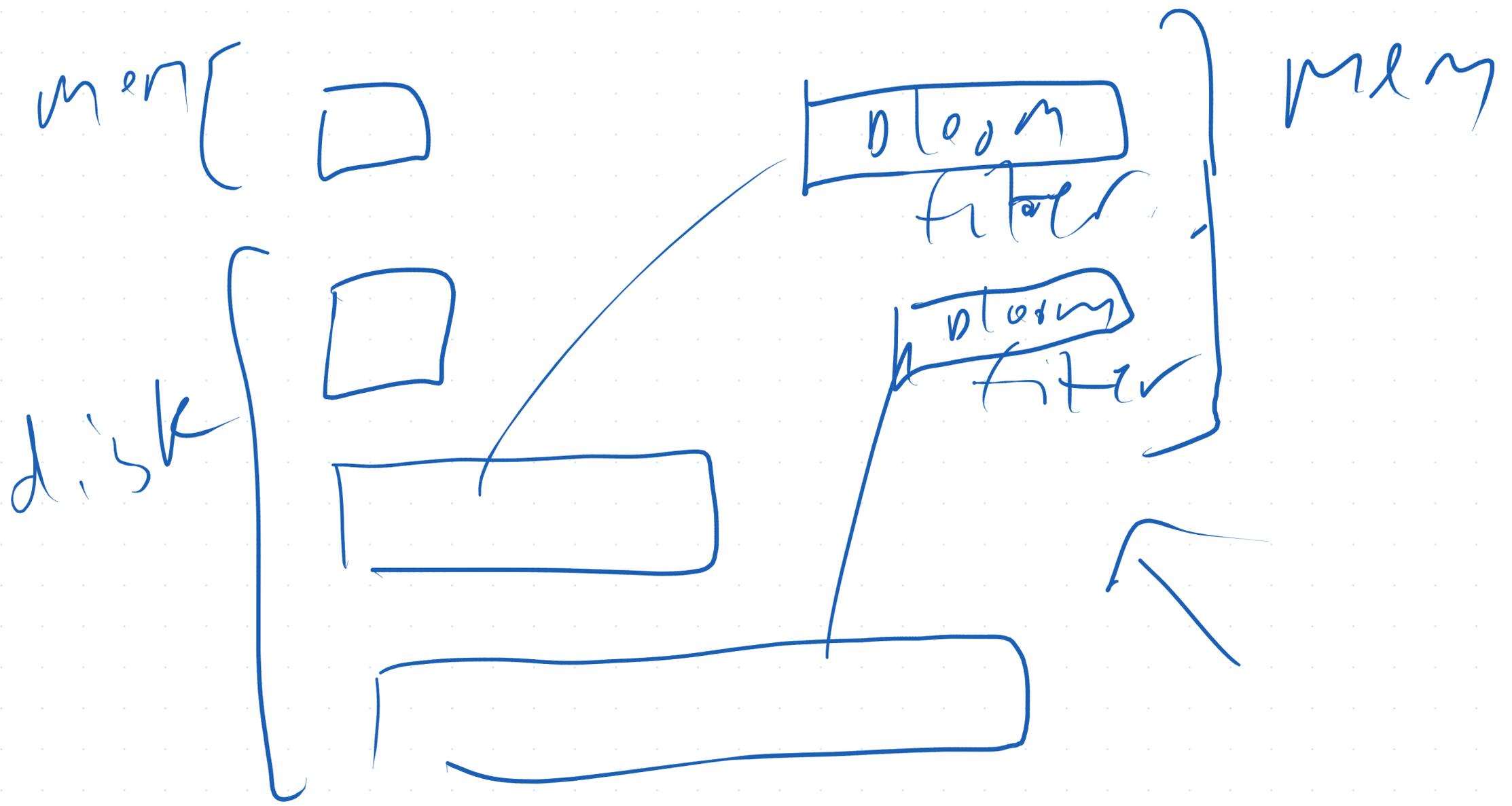
on lookup(x)

for i in $[0, R)$

if !bits[i][$h_i(x)$]

return false

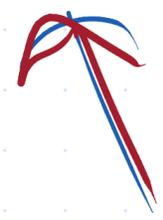
return true



}

read

$$O\left(\log\left(\frac{N}{B}\right)\right) \cdot O(\log N) = O(\log^2(N)) \quad \text{cost per read}$$
$$= (\log(N) - \log(B)) \cdot \log(N)$$
$$= \log^2(N) - \log(B)\log(N)$$



read amplification

write

Observation: Each record is copied k times (at most)

$$\uparrow \log\left(\frac{N}{B}\right)$$

← write amplification

{ bits [B]
}

Bloom Filter

on insert(x)
for $i \in \{0, \dots, k\}$

bits($h(x) \sim$) = true

Bloom Join

X

A	x
B	x
C	x
E	

~~(B₂, C₂, D₂, F₂)~~

Y

B ₂
C ₂
D ₂
F ₂

Bloom(B, C, D, F)

(A, B, C)

(B, B₂)

(C, C₂)

(B, B₂)

(C, C₂)

N records, each record copied at most $\log_2\left(\frac{N}{B}\right)$ times

$\hookrightarrow O\left(N \log\left(\frac{N}{B}\right)\right)$ - so per insertion cost is $\log\left(\frac{N}{B}\right)$ writes
(amortized)

Bloom Filter

add(x)

→ Inserts
x into
the filter

get(x)

→ If x was added
return true
always

otherwise
return false
with ϵ prob

{ buffer: (kk; B)
data; Vec(File)

)
on insert (data) add to buffer

if buffer full

write out to disk

while there are 2 arrays at a level

merge into next level down