

# Clark's Select

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

Department of Computer Science



Please sign guestbook ([www.langmead-lab.org/teaching-materials](http://www.langmead-lab.org/teaching-materials)) to tell me briefly how you are using the slides. For original Keynote files, email me ([ben.langmead@gmail.com](mailto:ben.langmead@gmail.com)).

# Clark's select

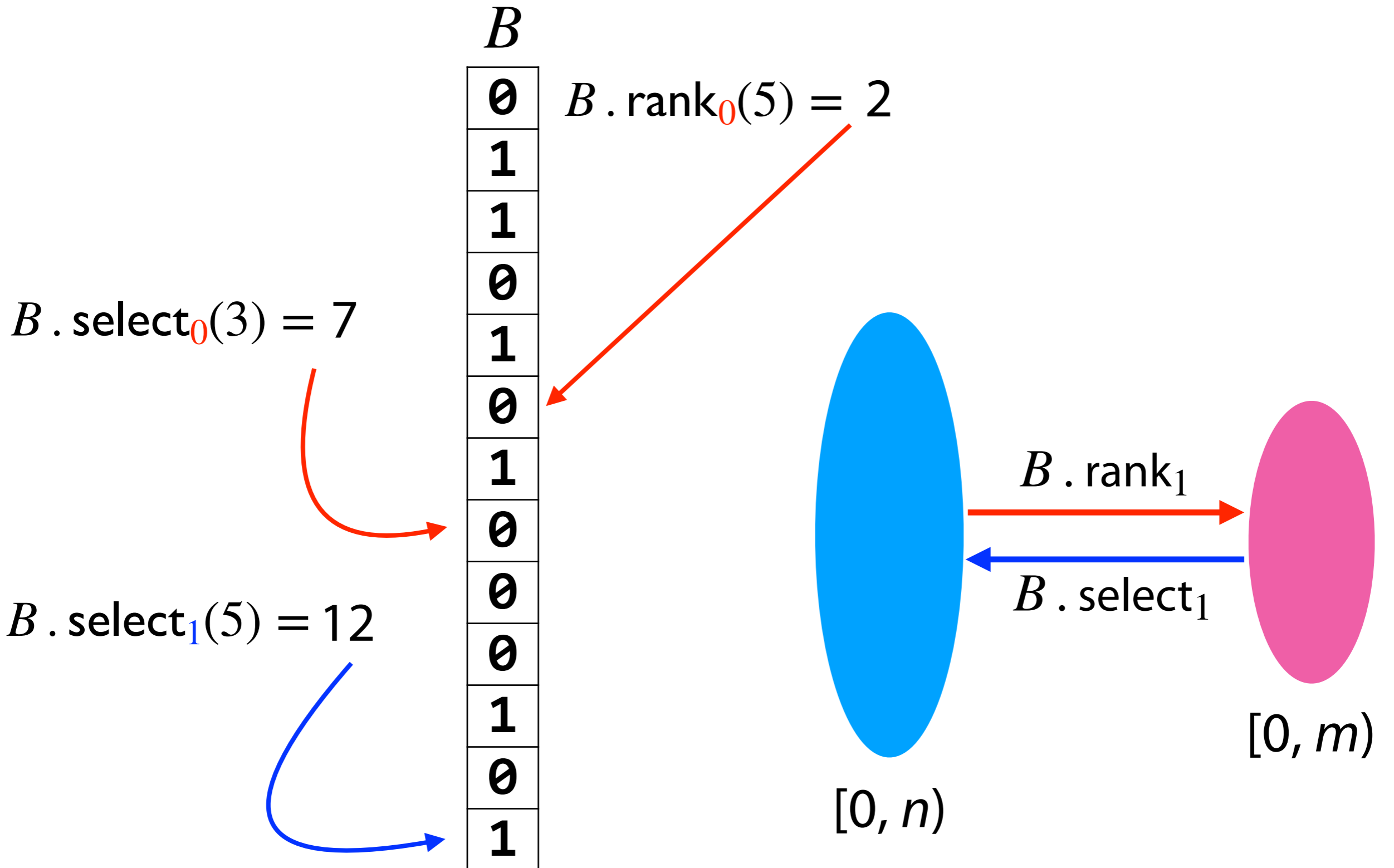
Unlike rank:

Chunks are defined by # 1s, not # bits

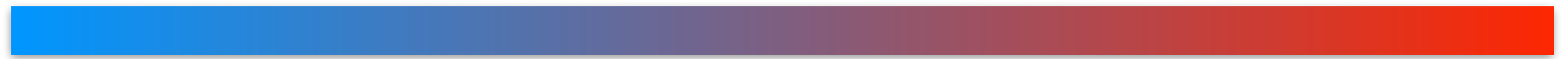
Two layers of "sparsity"

Answer is an offset into bitvector

# Bitvectors



# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)

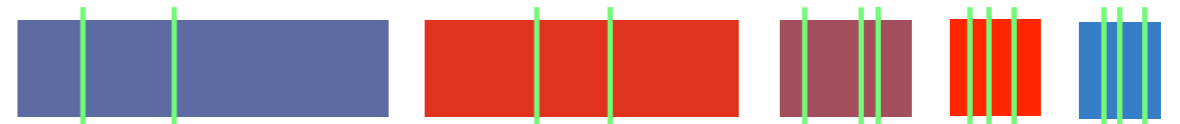


Lookup table for  
each 1-bit ✓

Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)



Lookup table for *all*  
*possible* sub-chunks ✓

# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)

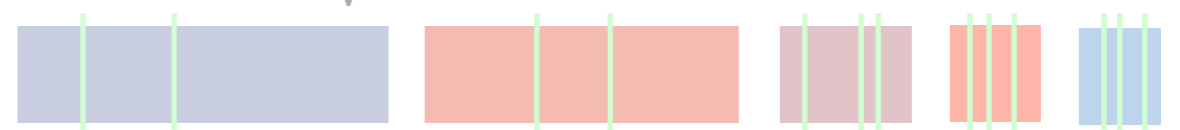


Lookup table for  
each 1-bit ✓

Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



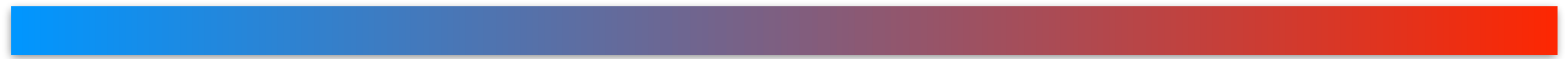
Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)



Lookup table for *all*  
*possible* sub-chunks ✓

# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)

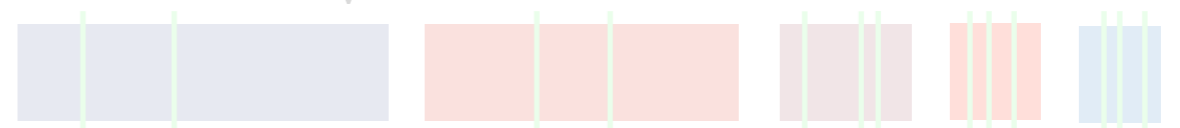


Lookup table for  
each 1-bit ✓

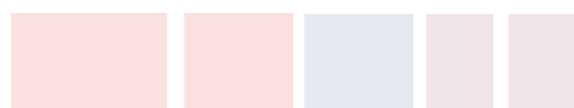
Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)



Lookup table for *all*  
*possible* sub-chunks ✓

# Clark's select



Split the string into chunks each containing  $\log^2 n$  1-bits



Larger chunks are **sparse**; 1's spread out

Shorter chunks are **dense**; 1's packed together

# Clark's select

Each chunk contains  $\log^2 n$  1-bits



We store **offset** of each chunk start

This takes:



# Clark's select

Each chunk contains  $\log^2 n$  1-bits



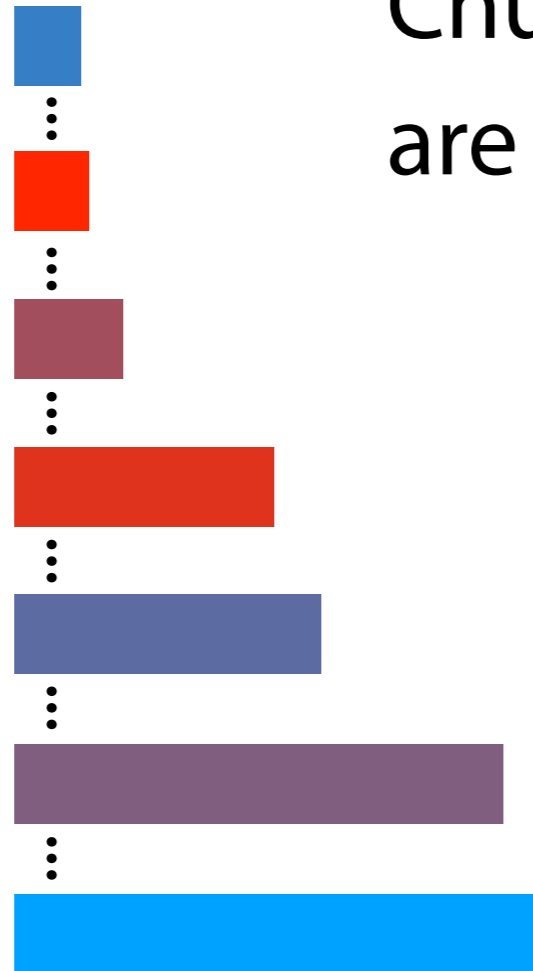
We store **offset** of each chunk start

This takes:

$$O\left(\underbrace{\frac{n}{\log^2 n} \log n}_{\text{Worst case: every bit set}}\right) = O\left(\frac{n}{\log n}\right) = \check{O}(n) \text{ bits}$$

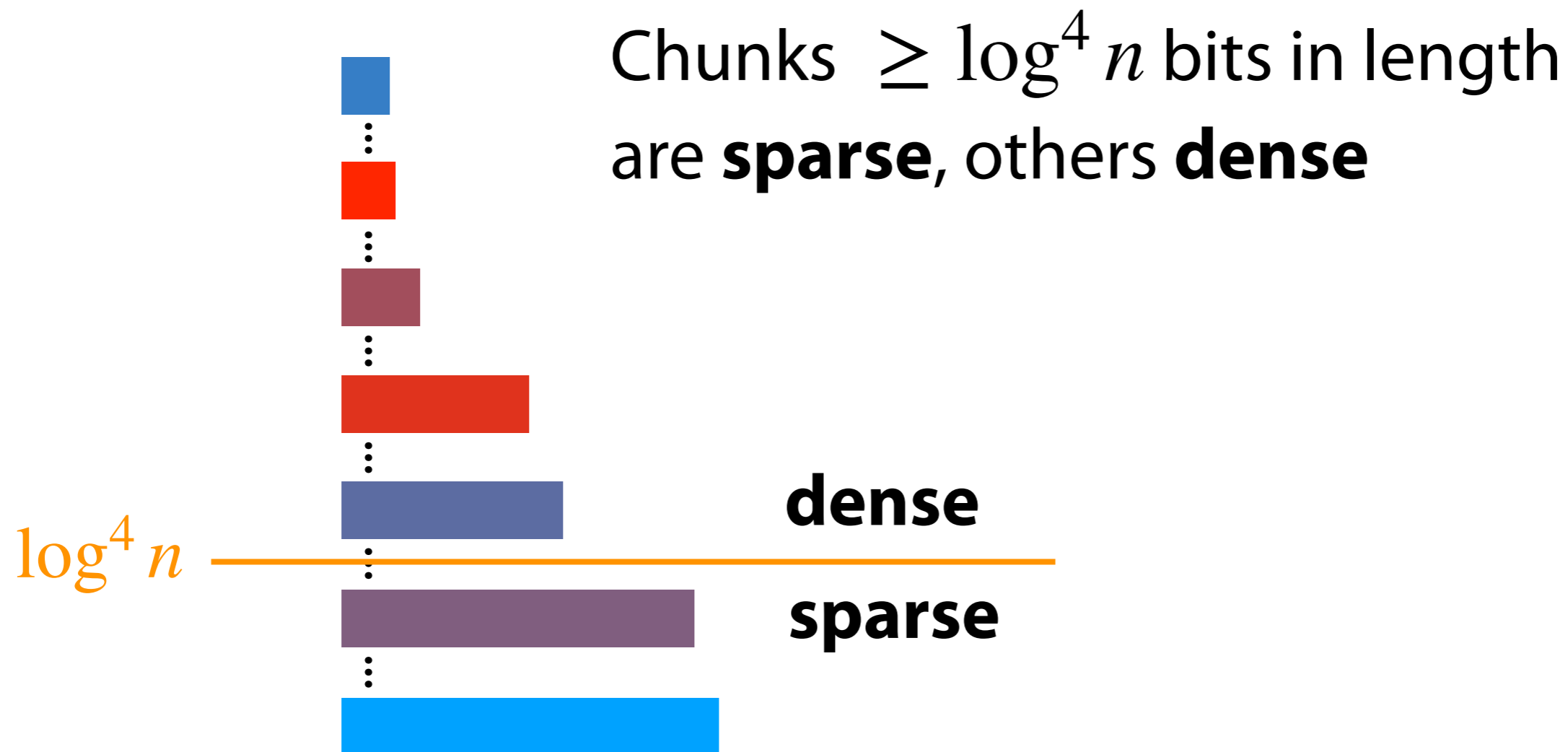
Worst case:  
every bit set

# Clark's select



Chunks  $\geq \log^4 n$  bits in length  
are **sparse**, others **dense**

# Clark's select



$\log^4 n$  is square of the # of set bits per chunk,  $\log^2 n$

("sparse" roughly means "less than  $\sqrt{\#}$  bits are set")

# Clark's select: sparse case

Pre-calculate  $B \cdot \text{select}_1$  for 1-bits in sparse chunks

# Clark's select: sparse case

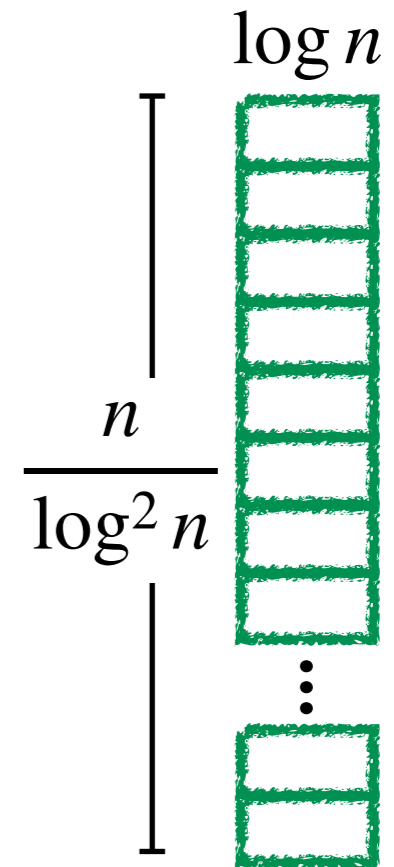
Pre-calculate  $B$ .  $\text{select}_1$  for 1-bits in sparse chunks

$$O\left(\underbrace{\frac{n}{\log^4 n}}_{\text{Max \# sparse chunks}} \cdot \underbrace{\log n}_{\text{\# bits to store 1 answer}} \cdot \underbrace{\log^2 n}_{\text{\# answers per chunk}}\right)$$

Max # sparse chunks

# bits to store 1 answer

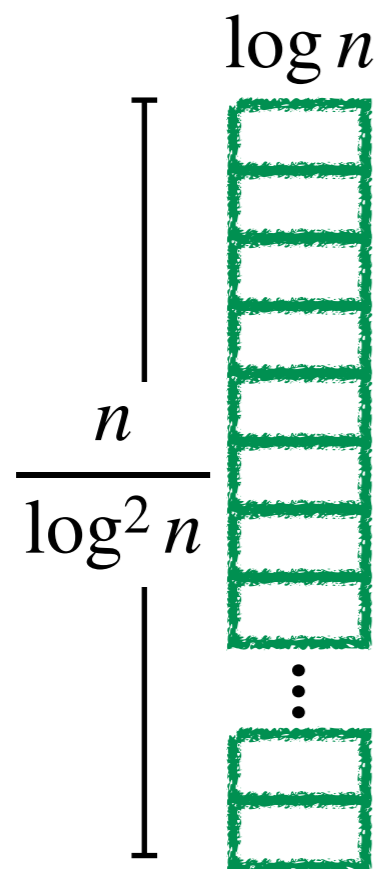
# answers per chunk



Offsets for 1-bits in sparse chunks

$$= O\left(\frac{n}{\log n}\right) = \check{O}(n)$$

# Clark's select



Offsets for 1-bits  
in sparse chunks

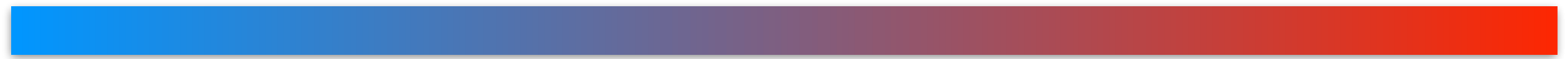
So far, strategy for select is:

- (a) find what chunk it's in (division)
- (b) if chunk is **sparse** ( $\geq \log^4 n$  bits)
  - (b.i) look up in **sparse offset table**
- (c) if chunk is **dense** ( $< \log^4 n$  bits)

**TODO**

Space is  $\tilde{O}(n)$  so far

# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)

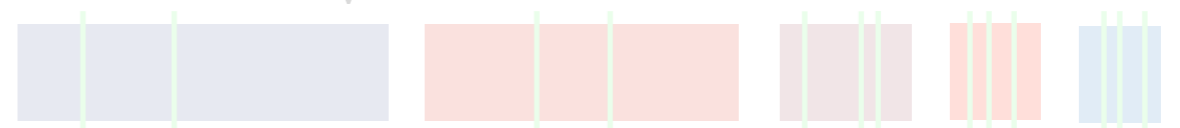


Lookup table for  
each 1-bit ✓

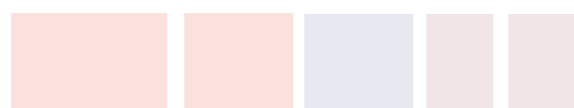
Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



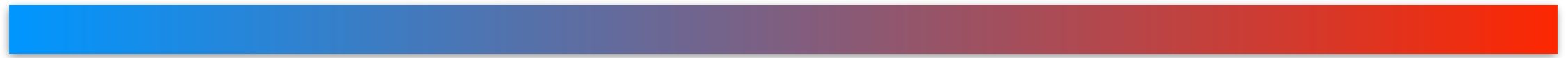
Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)



Lookup table for *all*  
*possible* sub-chunks ✓

# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)

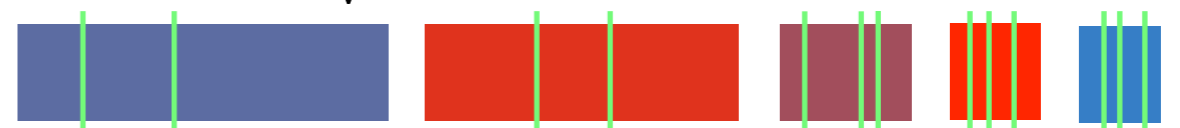


Lookup table for  
each 1-bit ✓

Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)

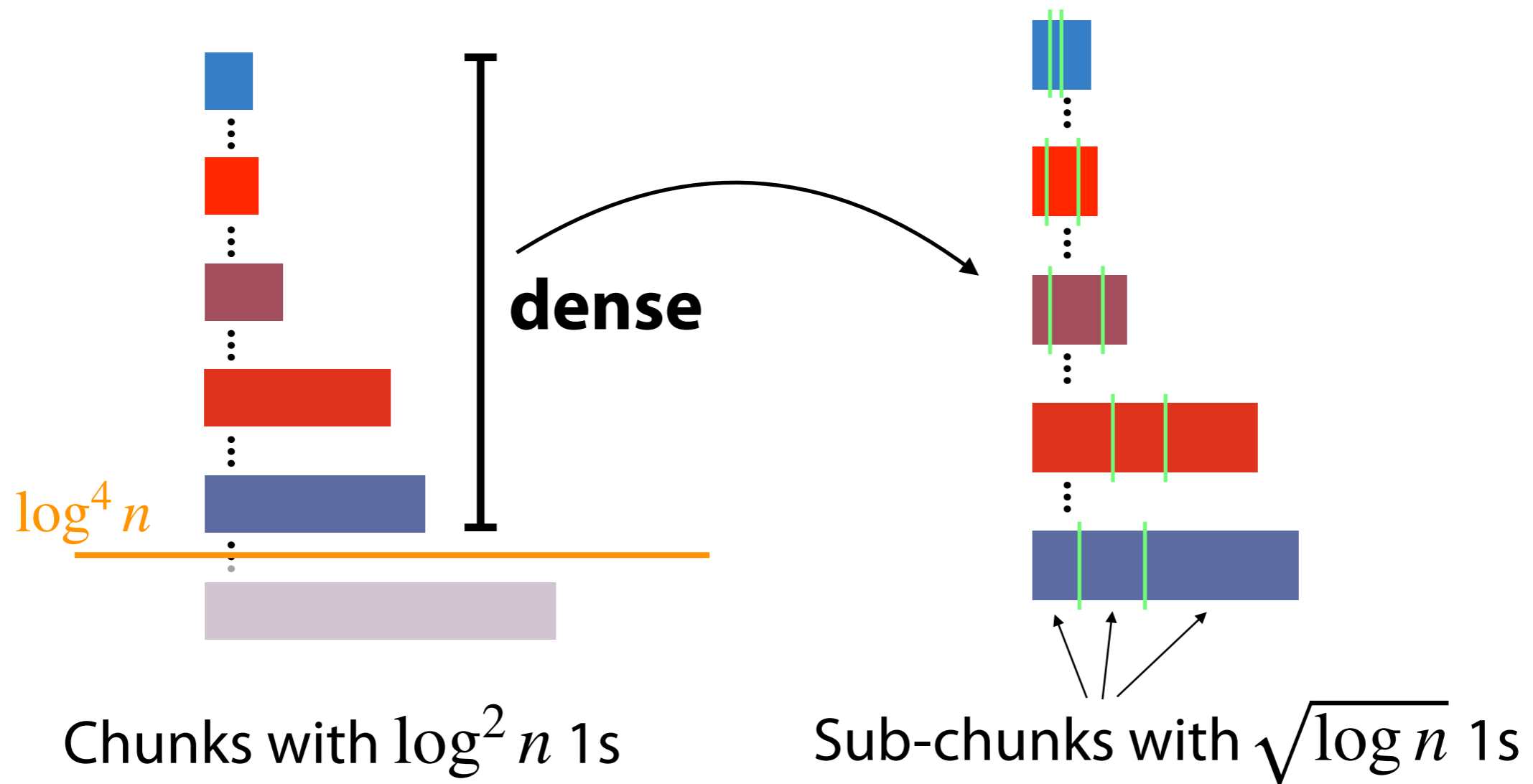


Lookup table for *all possible* sub-chunks ✓



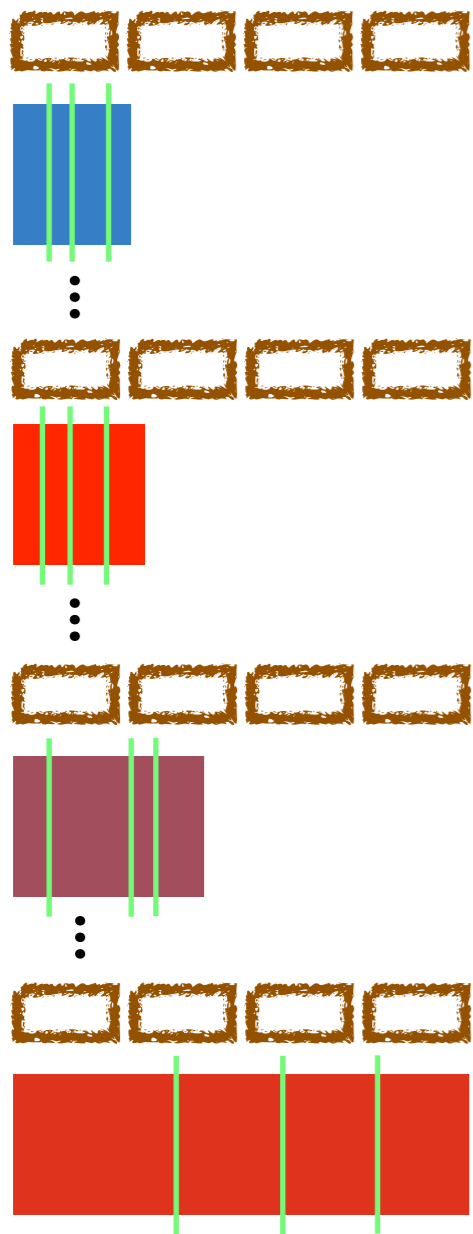
# Clark's select: dense case

**Dense** chunks are shorter than  $\log^4 n$  bits; further subdivide to sub-chunks of  $\sqrt{\log n}$  1-bits each

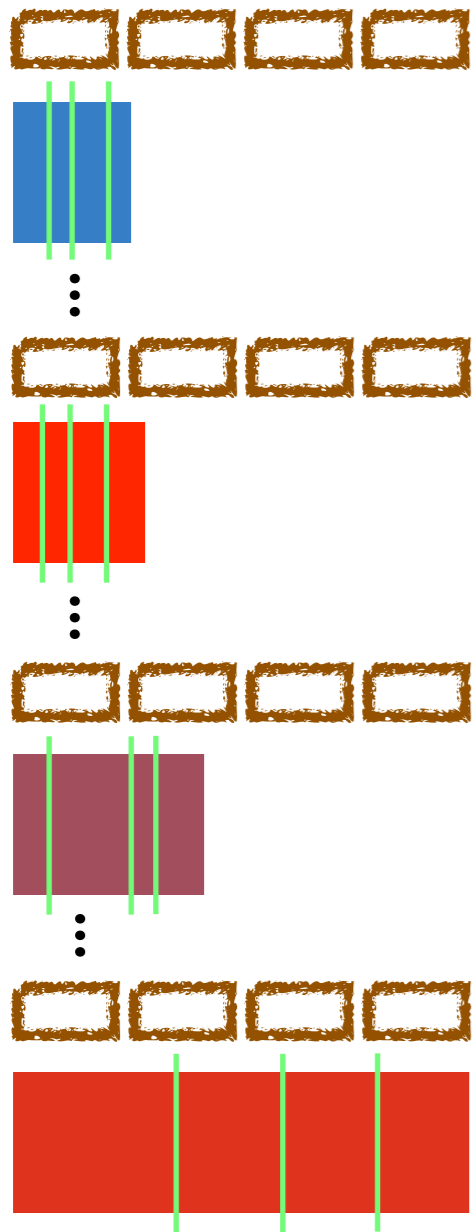


# Clark's select: dense case

Store **relative offset** per sub-chunk



# Clark's select: dense case

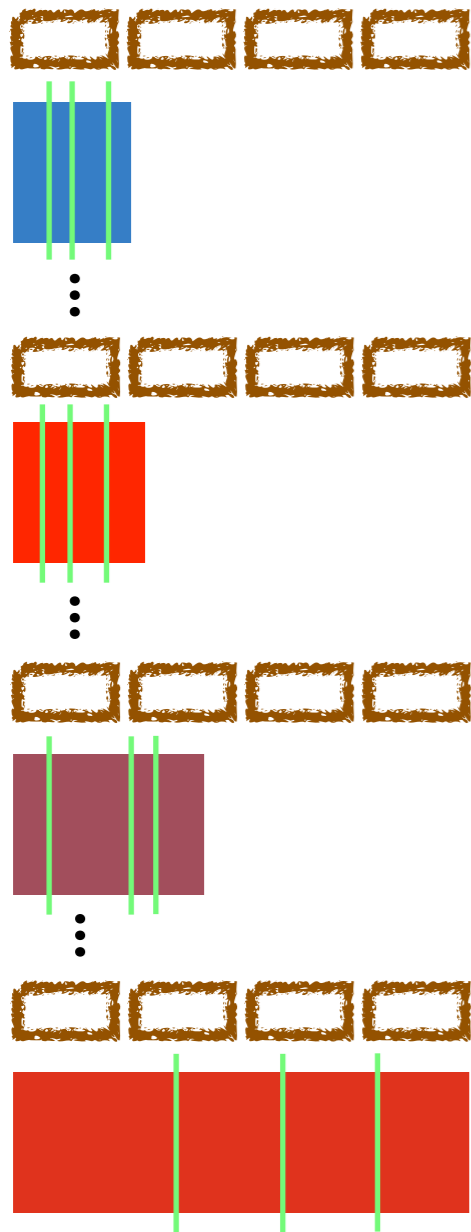


Store **relative offset** per sub-chunk

At most  $n/\sqrt{\log n}$  sub-chunks

Containing chunk has  $< \log^4 n$  bits, so **relative offset** fits in

# Clark's select: dense case



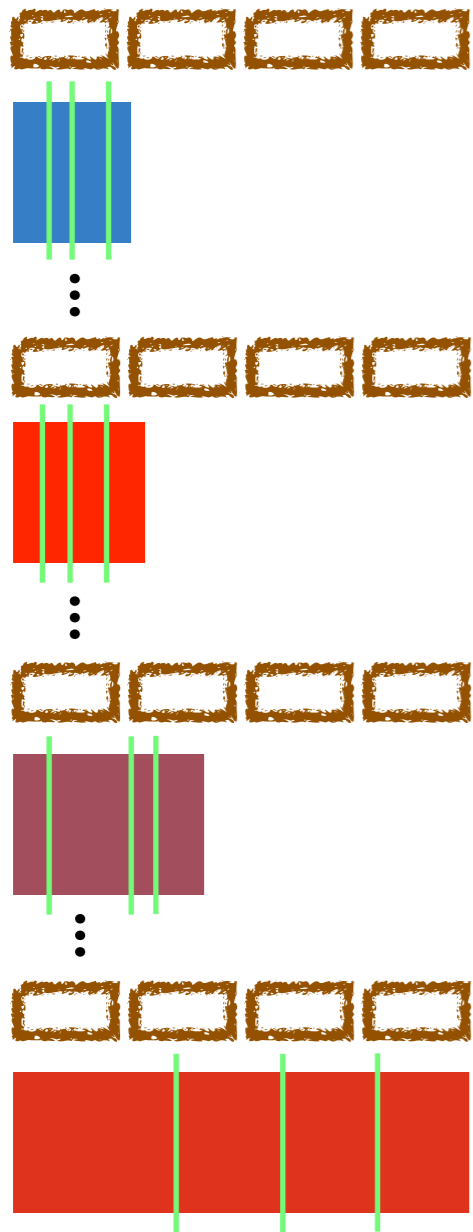
Store **relative offset** per sub-chunk

At most  $n/\sqrt{\log n}$  sub-chunks

Containing chunk has  $< \log^4 n$  bits, so **relative offset** fits in

$O(\log \log^4 n) = O(\log \log n)$  bits

# Clark's select: dense case



Store **relative offset** per sub-chunk

At most  $n/\sqrt{\log n}$  sub-chunks

Containing chunk has  $< \log^4 n$  bits, so **relative offset** fits in

$O(\log \log^4 n) = O(\log \log n)$  bits

Overall:  $O\left(\frac{n \log \log n}{\sqrt{\log n}}\right) = \check{o}(n)$

# Clark's select

So far:

**(a)** find what chunk it's in (division)

**(b)** if chunk is **sparse**

**(b.i)** look up in **sparse offset table**

**(c)** if chunk is **dense**

**(c.i)** look up **chunk's offset**

**(c.ii)** find what sub-chunk it's in (division by  $\sqrt{\log n}$ )

**(c.iii)** look up sub-chunk's **relative offset**

**TODO: need to look *within* sub-chunks**

# Clark's select

So far:

**(a)** find what chunk it's in (division)

**(b)** if chunk is **sparse**

**(b.i)** look up in **sparse offset table**

**(c)** if chunk is **dense**

**(c.i)** look up **chunk's offset**

**(c.ii)** find what sub-chunk it's in (division by  $\sqrt{\log n}$ )

**(c.iii)** look up sub-chunk's **relative offset**

**TODO:** need to look *within* sub-chunks

# Clark's select

So far:

(a) find what chunk it's in (division)

(b) if chunk is **sparse**

(b.i) look up in **sparse offset table**

(c) if chunk is **dense**

(c.i) look up **chunk's offset**

(c.ii) find what sub-chunk it's in (division by  $\sqrt{\log n}$ )

(c.iii) look up sub-chunk's **relative offset**

**TODO: need to look *within* sub-chunks**



# Clark's select: dense/sparse case

Sub-chunks with  $\geq 1/2 \log n$  bits are **sparse**;  
we simply store relative offsets for every 1-bit

# Clark's select: dense/sparse case

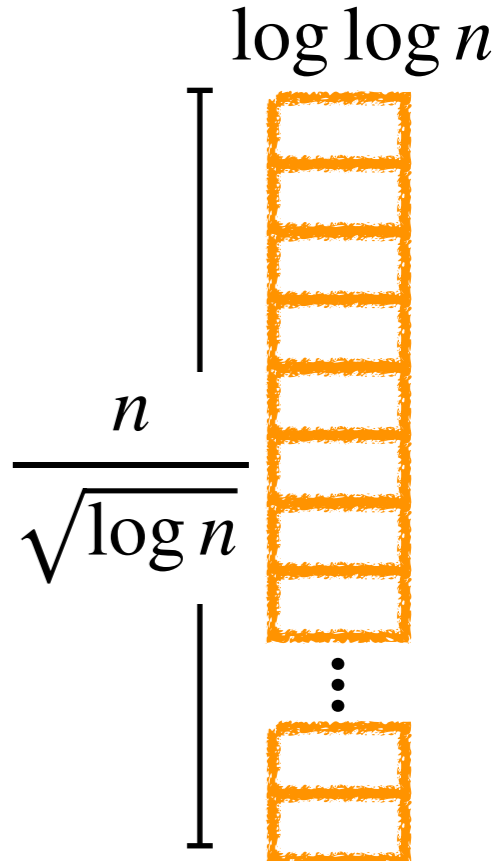
Sub-chunks with  $\geq 1/2 \log n$  bits are **sparse**;  
we simply store relative offsets for every 1-bit

Overall:  $O\left(\underbrace{\frac{n}{1/2 \log n}}_{\text{Max \# sparse sub-chunks}} \underbrace{\log \log n}_{\substack{\text{\# bits to store} \\ \text{1 answer (rel.} \\ \text{to chunk)}}} \underbrace{\sqrt{\log n}}_{\substack{\text{\# 1-bits} \\ \text{per chunk}}}\right)$

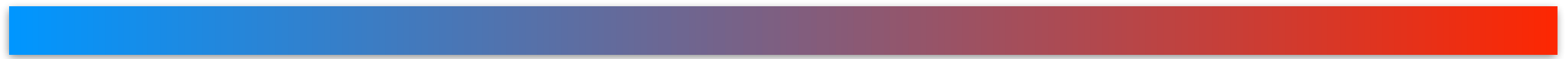
# Clark's select: dense/sparse case

Sub-chunks with  $\geq 1/2 \log n$  bits are **sparse**;  
we simply store relative offsets for every 1-bit

Overall:  $O\left(\underbrace{\frac{n}{1/2 \log n}}_{\text{Max \# sparse sub-chunks}} \underbrace{\log \log n}_{\text{\# bits to store 1 answer (rel. to chunk)}} \underbrace{\sqrt{\log n}}_{\text{\# 1-bits per chunk}}\right)$

$$= O\left(\frac{n \sqrt{\log n} \log \log n}{\log n}\right) = O\left(\frac{n \log \log n}{\sqrt{\log n}}\right) = \check{o}(n)$$


# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)



Lookup table for  
each 1-bit ✓

Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



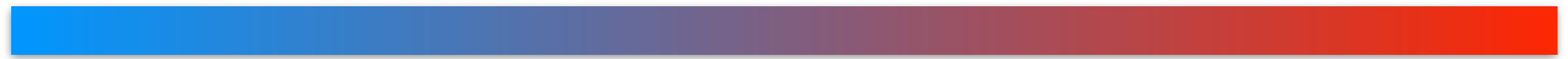
Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)



Lookup table for *all possible* sub-chunks ✓

# Clark's select



Split into  $\log^2 n$ -weight chunks



Sparse ( $\geq \log^4 n$ -length)

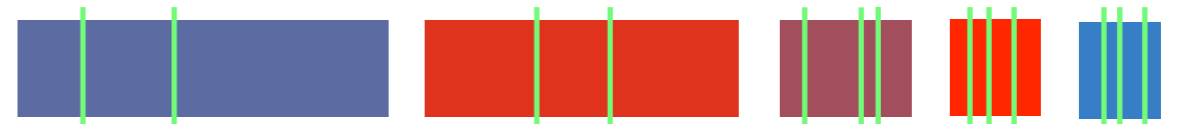


Lookup table for  
each 1-bit ✓

Dense ( $< \log^4 n$ -length)



Split into  $\sqrt{\log n}$ -weight sub-chunks



Sparse ( $\geq 1/2 \log n$ -length)



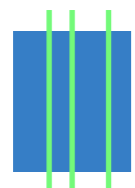
Lookup table for  
each 1-bit ✓

Dense ( $< 1/2 \log n$ -length)

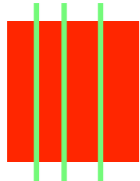


Lookup table for *all*  
*possible* sub-chunks ✓

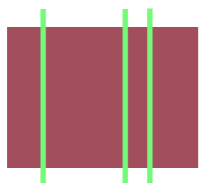
# Clark's select: dense/dense case



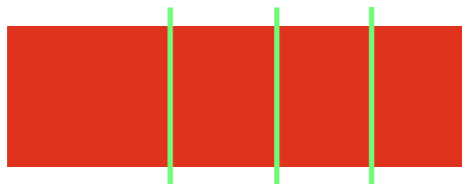
⋮



⋮



⋮



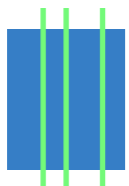
Sub-chunks  $< 1/2 \log n$  bits are **dense**;  
pre-calculate answers for all such chunks, like rank:

$$2^{1/2 \log n}$$

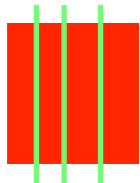
possible  
bitvectors

# Clark's select: dense/dense case

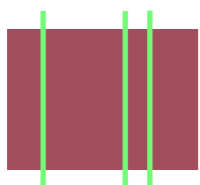
Sub-chunks  $< 1/2 \log n$  bits are **dense**;  
pre-calculate answers for all such chunks, like rank:



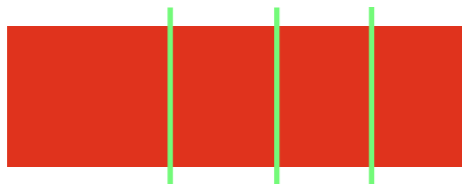
⋮



⋮



⋮

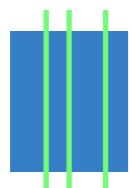


$$2^{1/2 \log n} \cdot \sqrt{\log n}$$

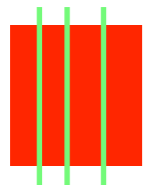
possible bitvectors      possible 1-bits

# Clark's select: dense/dense case

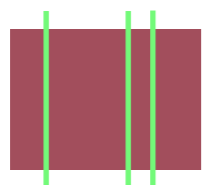
Sub-chunks  $< 1/2 \log n$  bits are **dense**;  
pre-calculate answers for all such chunks, like rank:



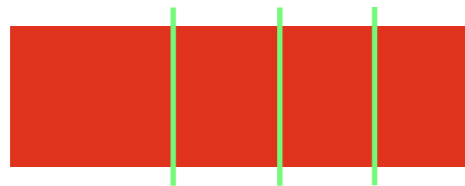
⋮



⋮



⋮



$$2^{1/2 \log n} \cdot \sqrt{\log n} \cdot \log \log n$$

possible  
bitvectors

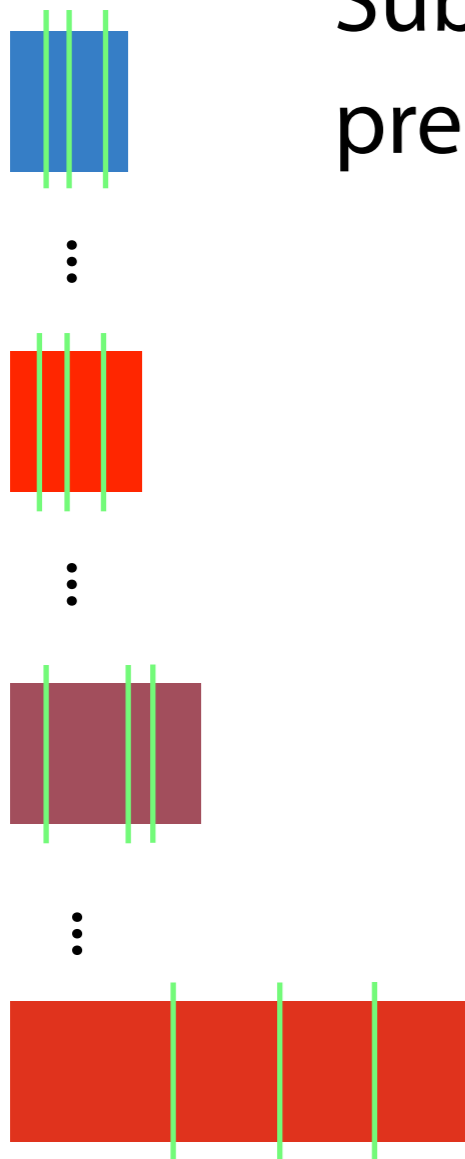
possible  
1-bits

answer



# Clark's select: dense/dense case

Sub-chunks  $< 1/2 \log n$  bits are **dense**;  
pre-calculate answers for all such chunks, like rank:

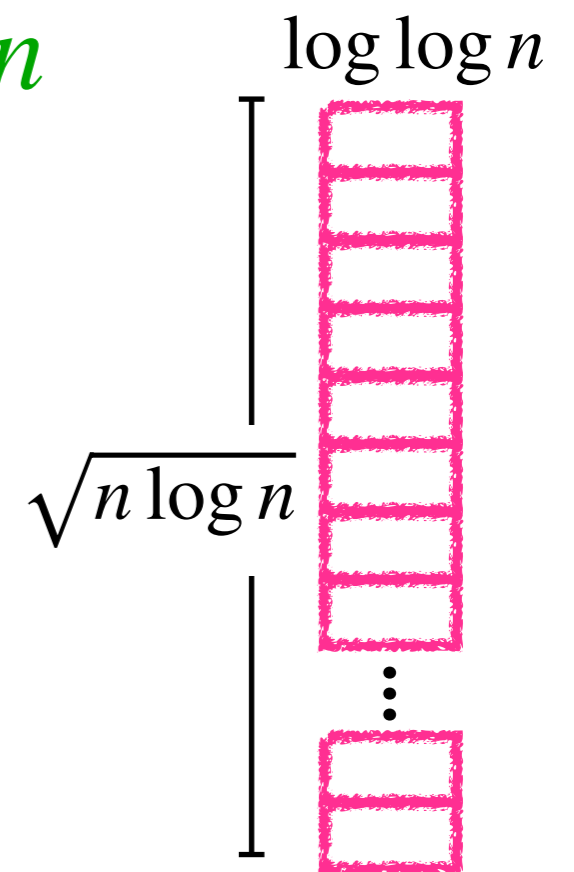


$$2^{1/2 \log n} \cdot \sqrt{\log n} \cdot \log \log n$$

possible bitvectors      possible 1-bits      answer

$$= O\left(\sqrt{n \log n} \log \log n\right)$$

$$= \check{o}(n)$$



# Clark's select

**(a)** find what chunk it's in (division by  $\log^2 n$ )

**(b)** if chunk is **sparse** ( $\geq \log^4 n$  bits)

**(b.i)** look up answer in **sparse offset table**

**(c)** if chunk is **dense** ( $< \log^4 n$  bits)

**(c.i)** look up **chunk's offset**

**(c.ii)** find what sub-chunk it's in (divide by  $\sqrt{\log n}$ )

**(c.iii)** look up sub-chunk's **relative offset**

**(c.iv)** if sub-chunk is **sparse** ( $\geq 1/2 \log n$  bits)

**(c.iv.A)** look up answer in **sparse 1-bit table**

**(c.iv.B)** return **(c.i) + (c.iii) + (c.iv.A)**

**(c.v)** if sub-chunk is **dense**

**(c.v.A)** look up answer in **all possible dense/dense table**

**(c.v.B)** return **(c.i) + (c.iii) + (c.v.A)**

# Clark's select

- (a) find what chunk it's in (division by  $\log^2 n$ )
- (b) if chunk is **sparse** ( $\geq \log^4 n$  bits)
  - (b.i) look up answer in **sparse offset table**
- (c) if chunk is **dense** ( $< \log^4 n$  bits)
  - (c.i) look up **chunk's offset**
  - (c.ii) find what sub-chunk it's in (divide by  $\sqrt{\log n}$ )
  - (c.iii) look up sub-chunk's **relative offset**
  - (c.iv) if sub-chunk is **sparse** ( $\geq 1/2 \log n$  bits)
    - (c.iv.A) look up answer in **sparse 1-bit table**
    - (c.iv.B) return (c.i) + (c.iii) + (c.iv.A)
  - (c.v) if sub-chunk is **dense**
    - (c.v.A) look up answer in **all possible dense/dense table**
    - (c.v.B) return (c.i) + (c.iii) + (c.v.A)

# Clark's select

- (a) find what chunk it's in (division by  $\log^2 n$ )
- (b) if chunk is **sparse** ( $\geq \log^4 n$  bits)
  - (b.i) look up answer in **sparse offset table**
- (c) if chunk is **dense** ( $< \log^4 n$  bits)
  - (c.i) look up **chunk's offset**
  - (c.ii) find what sub-chunk it's in (divide by  $\sqrt{\log n}$ )
  - (c.iii) look up sub-chunk's **relative offset**
  - (c.iv) if sub-chunk is **sparse** ( $\geq 1/2 \log n$  bits)
    - (c.iv.A) look up answer in **sparse 1-bit table**
    - (c.iv.B) return (c.i) + (c.iii) + (c.iv.A)
  - (c.v) if sub-chunk is **dense**
    - (c.v.A) look up answer in **all possible dense/dense table**
    - (c.v.B) return (c.i) + (c.iii) + (c.v.A)

# Clark's select

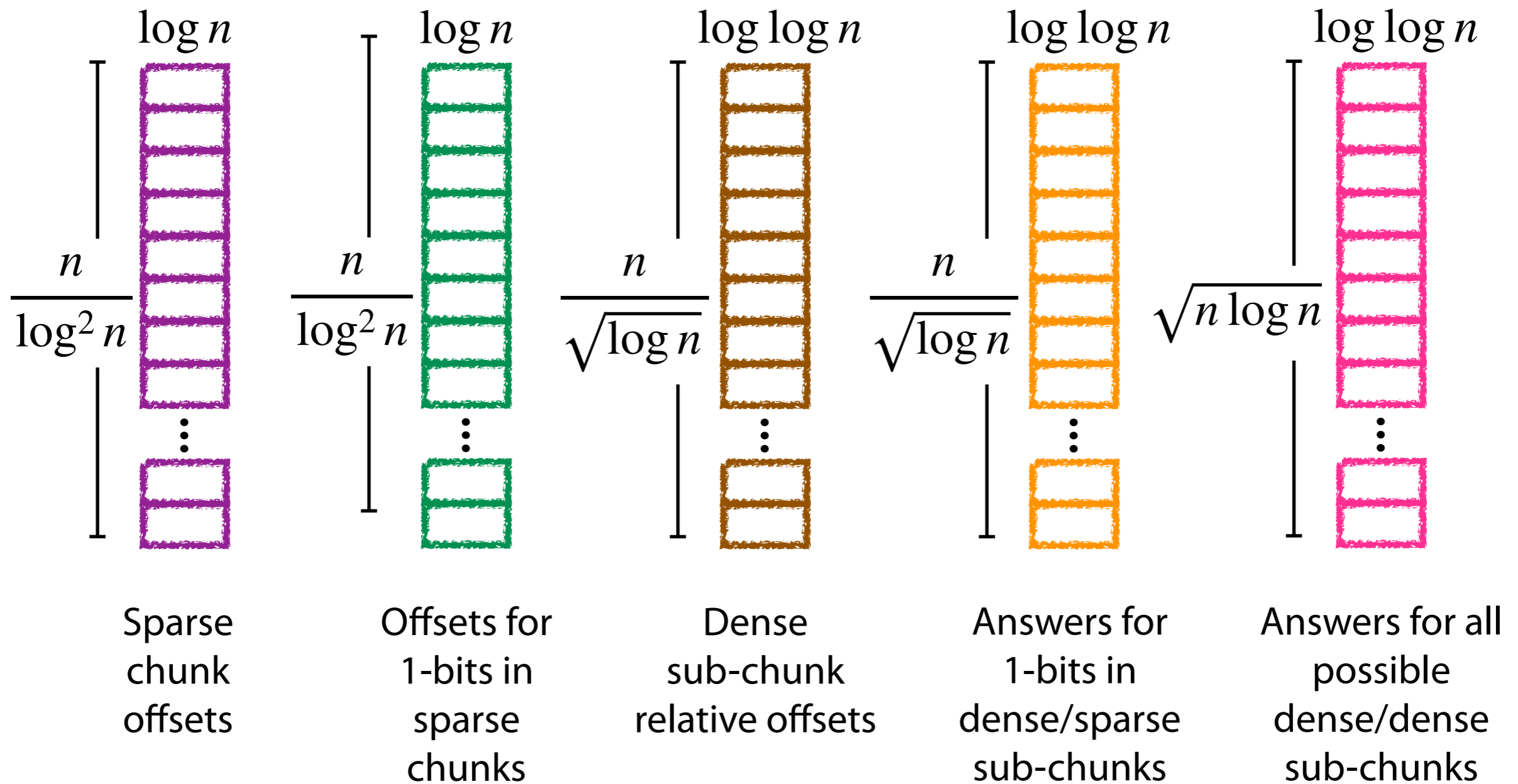
- (a) find what chunk it's in (division by  $\log^2 n$ )
- (b) if chunk is **sparse** ( $\geq \log^4 n$  bits)
  - (b.i) look up answer in **sparse offset table**
- (c) if chunk is **dense** ( $< \log^4 n$  bits)
  - (c.i) look up **chunk's offset**
  - (c.ii) find what sub-chunk it's in (divide by  $\sqrt{\log n}$ )
  - (c.iii) look up sub-chunk's **relative offset**
  - (c.iv) if sub-chunk is **sparse** ( $\geq 1/2 \log n$  bits)
    - (c.iv.A) look up answer in **sparse 1-bit table**
    - (c.iv.B) return (c.i) + (c.iii) + (c.iv.A)
  - (c.v) if sub-chunk is **dense**
    - (c.v.A) look up answer in **all possible dense/dense table**
    - (c.v.B) return (c.i) + (c.iii) + (c.v.A)

# Clark's select

- (a) find what chunk it's in (division by  $\log^2 n$ )
- (b) if chunk is **sparse** ( $\geq \log^4 n$  bits)
  - (b.i) look up answer in **sparse offset table**
- (c) if chunk is **dense** ( $< \log^4 n$  bits)
  - (c.i) look up **chunk's offset**
  - (c.ii) find what sub-chunk it's in (divide by  $\sqrt{\log n}$ )
  - (c.iii) look up sub-chunk's **relative offset**
  - (c.iv) if sub-chunk is **sparse** ( $\geq 1/2 \log n$  bits)
    - (c.iv.A) look up answer in **sparse 1-bit table**
    - (c.iv.B) return (c.i) + (c.iii) + (c.iv.A)
  - (c.v) if sub-chunk is **dense**
    - (c.v.A) look up answer in **all possible dense/dense table**
    - (c.v.B) return (c.i) + (c.iii) + (c.v.A)

# Clark's select

Overall, space is  $\tilde{O}(n)$



# Bitvectors

The dream! 🛏️ 🤔

The reality! 🎆 🥰 🎉

	Time	Space (bits)	Note
$B$ . access	$O(1)$	$n$	Lookup
$B$ . select <sub>1</sub>	$O(1)$	$\check{O}(n)$	✓ Clark
$B$ . rank <sub>1</sub>	$O(1)$	$\check{O}(n)$	✓ Jacobson