

# Elias-Fano Encoding

A powerful tool for data structure design

Giulio Ermanno Pibiri  
giulio.pibiri@di.unipi.it  
*University of Pisa, and ISTI-CNR*



Tokyo, 10/04/2018

# Problem

Consider a sequence  $S[0,n)$  of  $n$  *positive* and *monotonically increasing integers*, i.e.,  $S[i-1] \leq S[i]$  for  $1 \leq i \leq n-1$ , possibly repeated.

How to represent it as a *bit vector* in which each original integer is *self-delimited*, using as few as possible bits?

# Problem

Consider a sequence  $S[0,n)$  of  $n$  *positive* and *monotonically increasing integers*, i.e.,  $S[i-1] \leq S[i]$  for  $1 \leq i \leq n-1$ , possibly repeated.

How to represent it as a *bit vector* in which each original integer is *self-delimited*, using as few as possible bits?

Huge research corpora describing different space/time trade-offs.

- Elias gamma/delta [Elias-1974]
- Variable Byte [Salomon-2007]
- Varint-G8IU [Stepanov et al.-2011]
- Simple-9/16 [Anh and Moffat 2005-2010]
- PForDelta (PFD) [Zukowski et al.-2006]
- OptPFD [Yan et al.-2009]
- Binary Interpolative Coding [Moffat and Stuiver-2000]

Given a *textual collection*  $D$ , each document can be seen as a (multi-)set of terms. The set of terms occurring in  $D$  is the *lexicon*  $T$ .

For each term  $t$  in  $T$  we store in a list  $L_t$  the identifiers of the documents in which  $t$  appears.

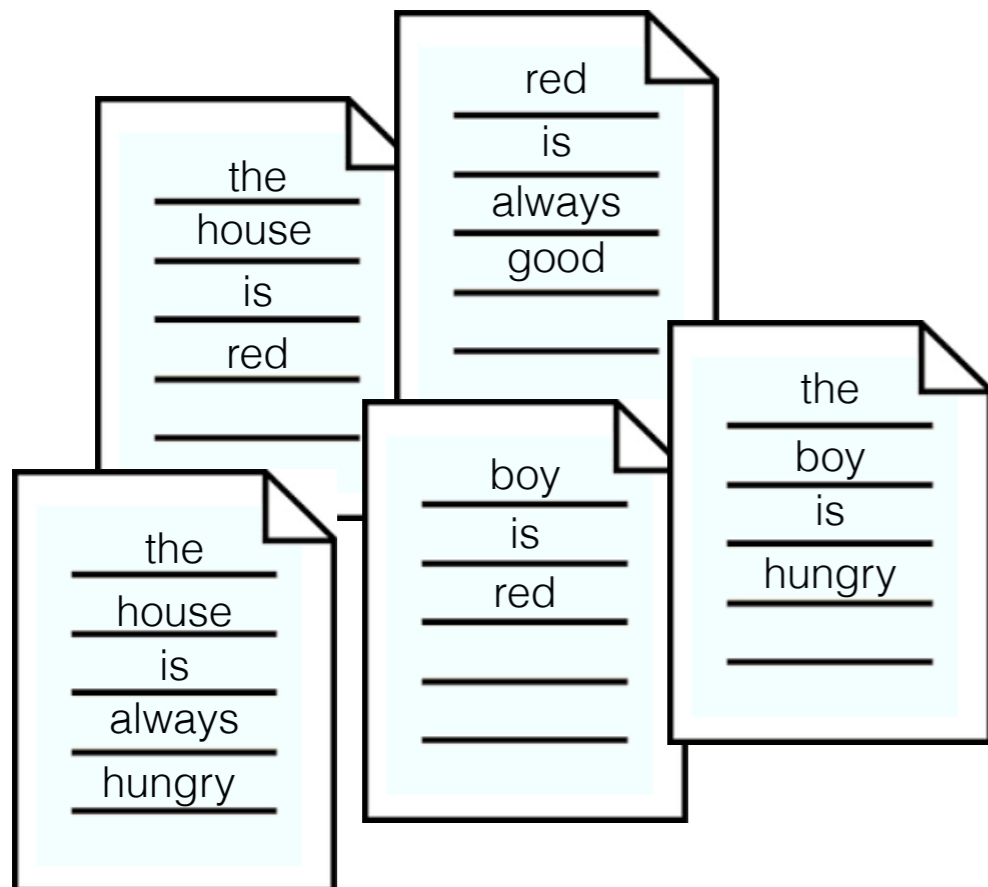
The collection of all inverted lists  $\{L_{t_1}, \dots, L_{t_T}\}$  is the inverted index.

# Inverted Indexes

Given a *textual collection*  $D$ , each document can be seen as a (multi-)set of terms. The set of terms occurring in  $D$  is the *lexicon*  $T$ .

For each term  $t$  in  $T$  we store in a list  $L_t$  the identifiers of the documents in which  $t$  appears.

The collection of all inverted lists  $\{L_{t_1}, \dots, L_{t_T}\}$  is the inverted index.

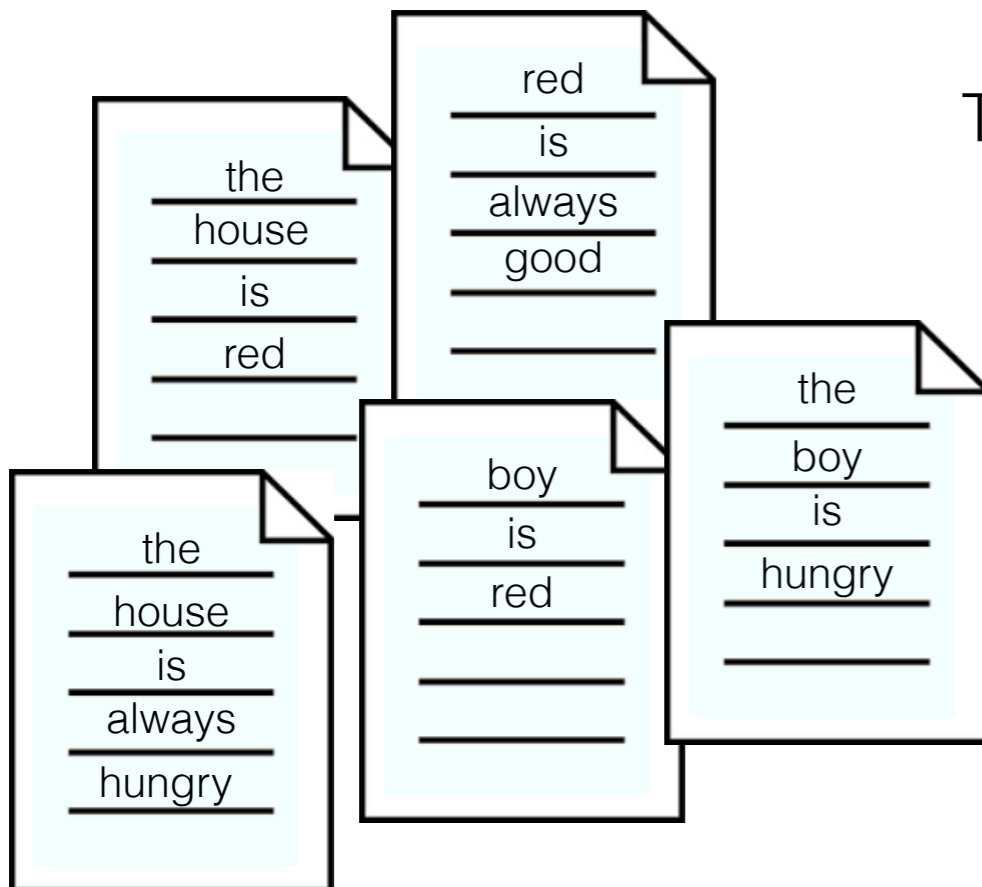


# Inverted Indexes

Given a *textual collection*  $D$ , each document can be seen as a (multi-)set of terms. The set of terms occurring in  $D$  is the *lexicon*  $T$ .

For each term  $t$  in  $T$  we store in a list  $L_t$  the identifiers of the documents in which  $t$  appears.

The collection of all inverted lists  $\{L_{t_1}, \dots, L_{t_T}\}$  is the inverted index.



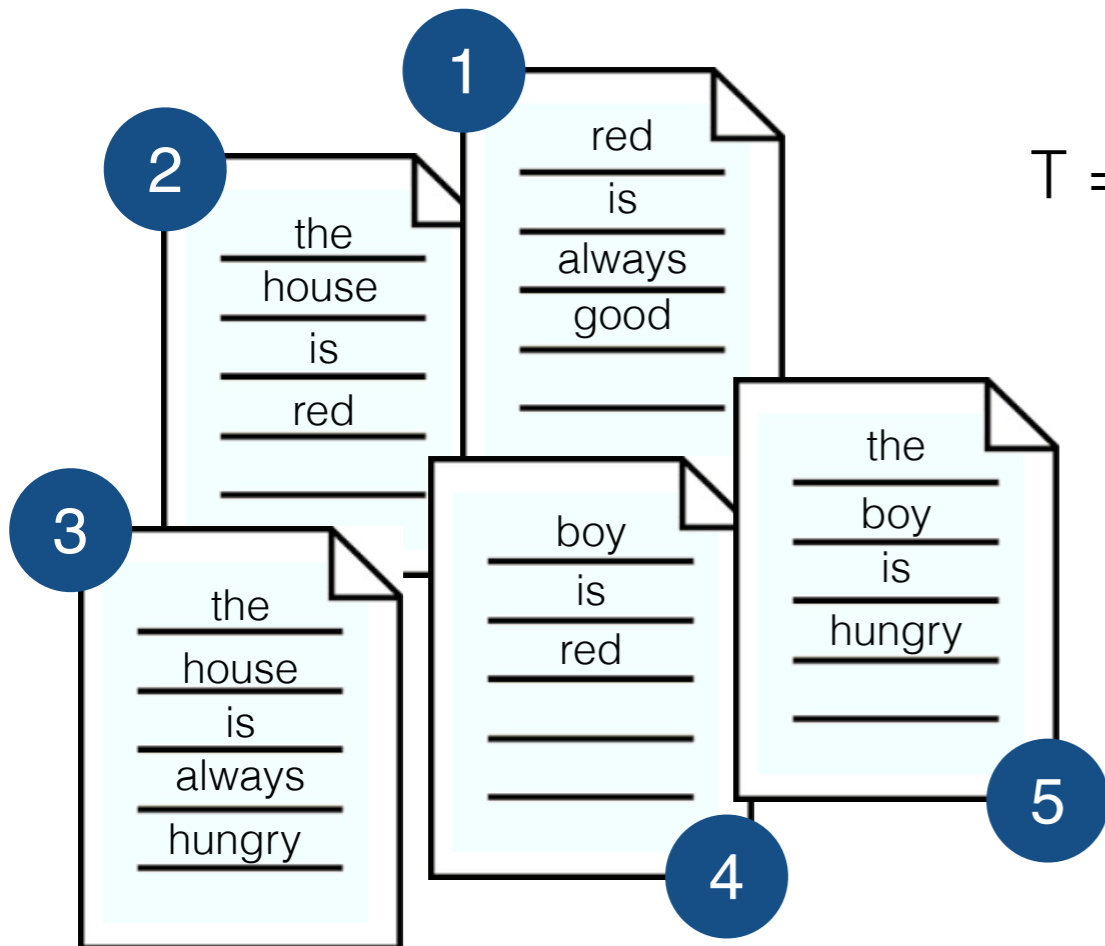
$t_1$        $t_2$        $t_3$        $t_4$        $t_5$        $t_6$        $t_7$        $t_8$   
 $T = \{\text{always, boy, good, house, hungry, is, red, the}\}$

# Inverted Indexes

Given a *textual collection*  $D$ , each document can be seen as a (multi-)set of terms. The set of terms occurring in  $D$  is the *lexicon*  $T$ .

For each term  $t$  in  $T$  we store in a list  $L_t$  the identifiers of the documents in which  $t$  appears.

The collection of all inverted lists  $\{L_{t_1}, \dots, L_{t_T}\}$  is the inverted index.



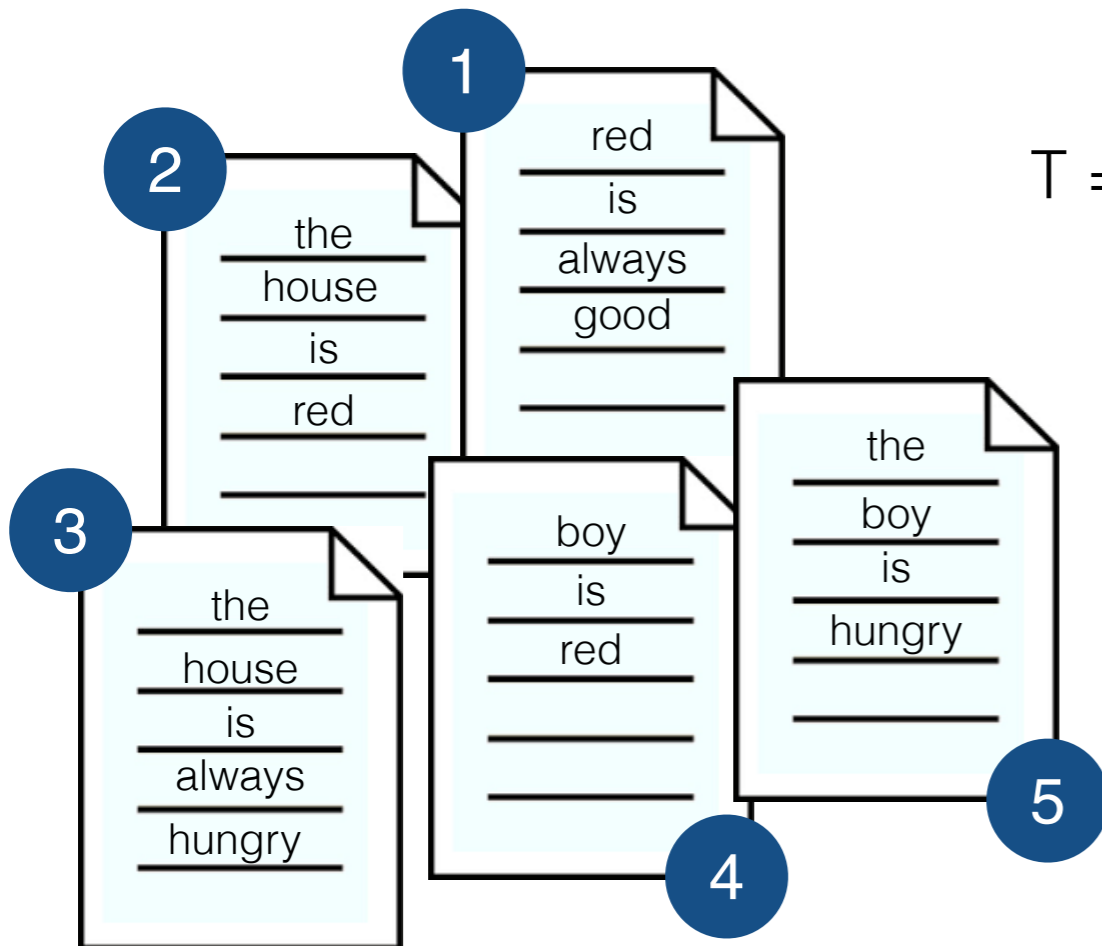
$t_1$        $t_2$        $t_3$        $t_4$        $t_5$        $t_6$        $t_7$        $t_8$   
 $T = \{\text{always, boy, good, house, hungry, is, red, the}\}$

# Inverted Indexes

Given a *textual collection*  $D$ , each document can be seen as a (multi-)set of terms. The set of terms occurring in  $D$  is the *lexicon*  $T$ .

For each term  $t$  in  $T$  we store in a list  $L_t$  the identifiers of the documents in which  $t$  appears.

The collection of all inverted lists  $\{L_{t_1}, \dots, L_{t_T}\}$  is the inverted index.



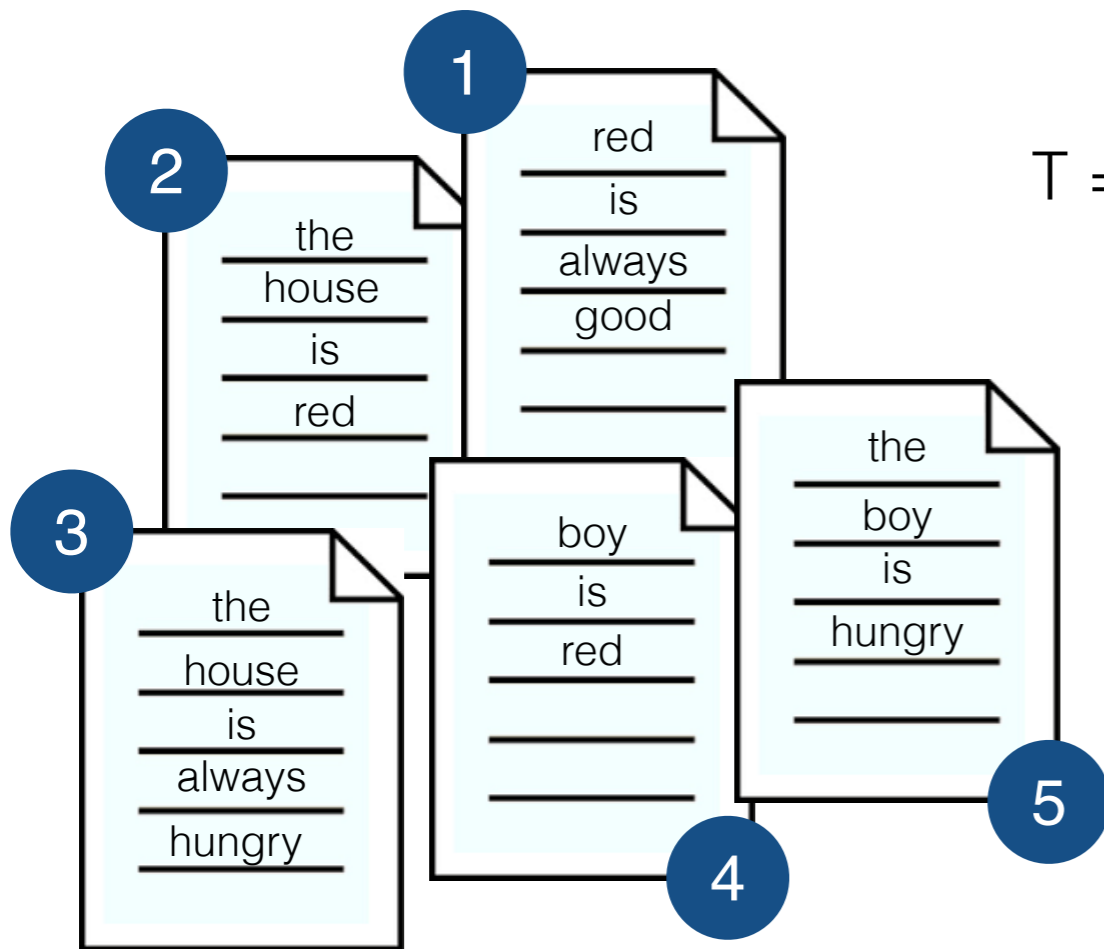
$T = \{t_1, t_2, t_3, t_4, t_5, t_6, t_7, t_8\}$   
 $T = \{\text{always, boy, good, house, hungry, is, red, the}\}$

$L_{t_1} = [1, 3]$   
 $L_{t_2} = [4, 5]$   
 $L_{t_3} = [1]$   
 $L_{t_4} = [2, 3]$   
 $L_{t_5} = [3, 5]$   
 $L_{t_6} = [1, 2, 3, 4, 5]$   
 $L_{t_7} = [1, 2, 4]$   
 $L_{t_8} = [2, 3, 5]$



Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$T = \{t_1, t_2, t_3, t_4, t_5, t_6, t_7, t_8\}$   
 $T = \{\text{always, boy, good, house, hungry, is, red, the}\}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

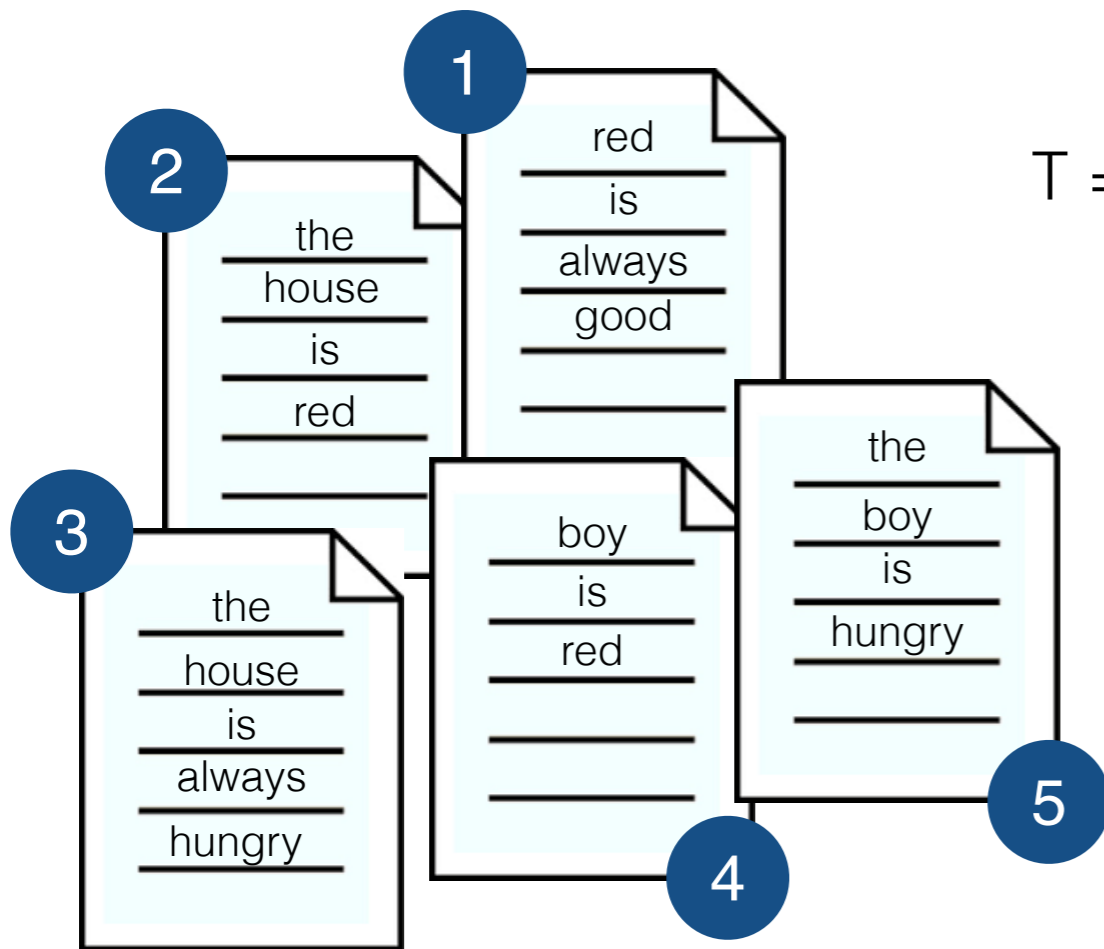
$L_{t_5} = [3, 5]$

$L_{t_6} = [1, 2, 3, 4, 5]$

$L_{t_7} = [1, 2, 4]$

$L_{t_8} = [2, 3, 5]$

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$T = \{ \begin{matrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 & t_7 & t_8 \\ \text{always, boy, good, house, hungry, is, red, the} \end{matrix} \}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

$L_{t_5} = [3, 5]$

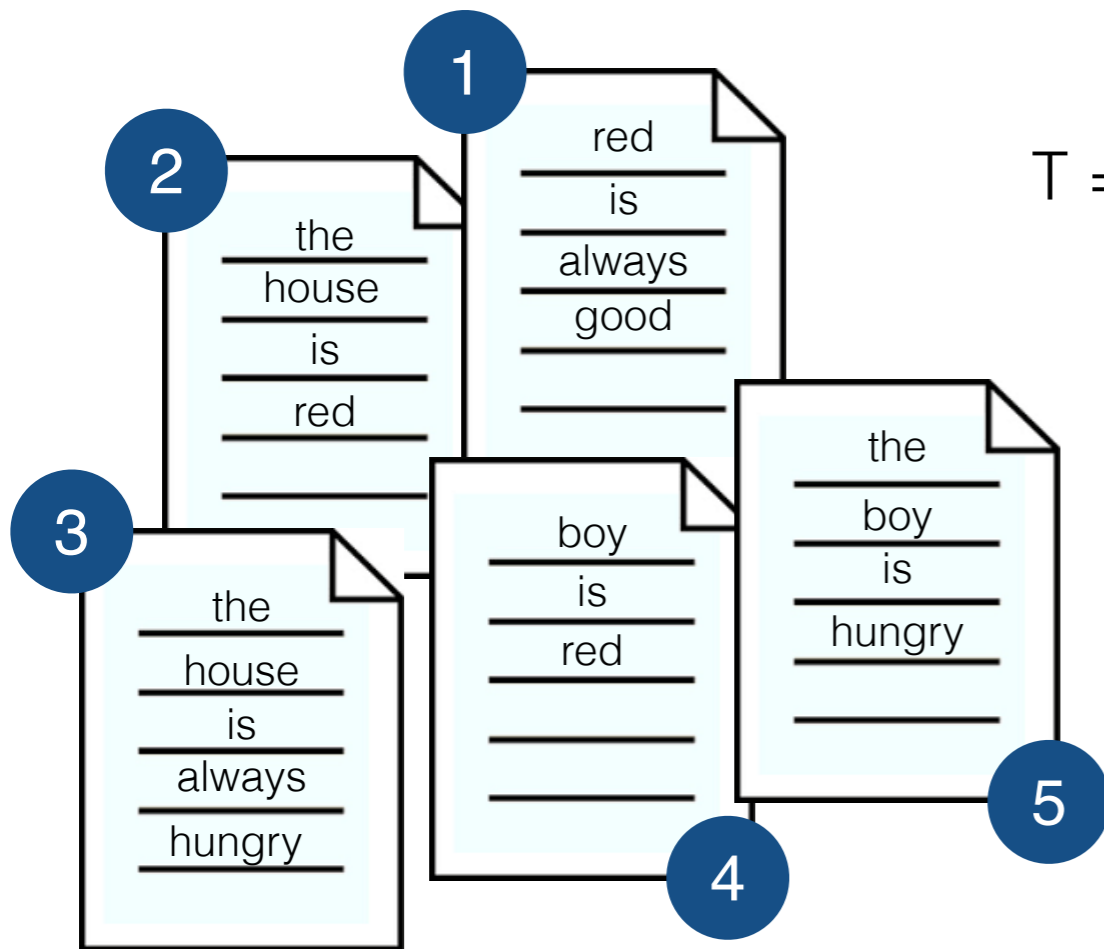
$L_{t_6} = [1, 2, 3, 4, 5]$

$L_{t_7} = [1, 2, 4]$

$L_{t_8} = [2, 3, 5]$

$q = \{ \text{boy, is, the} \}$

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.

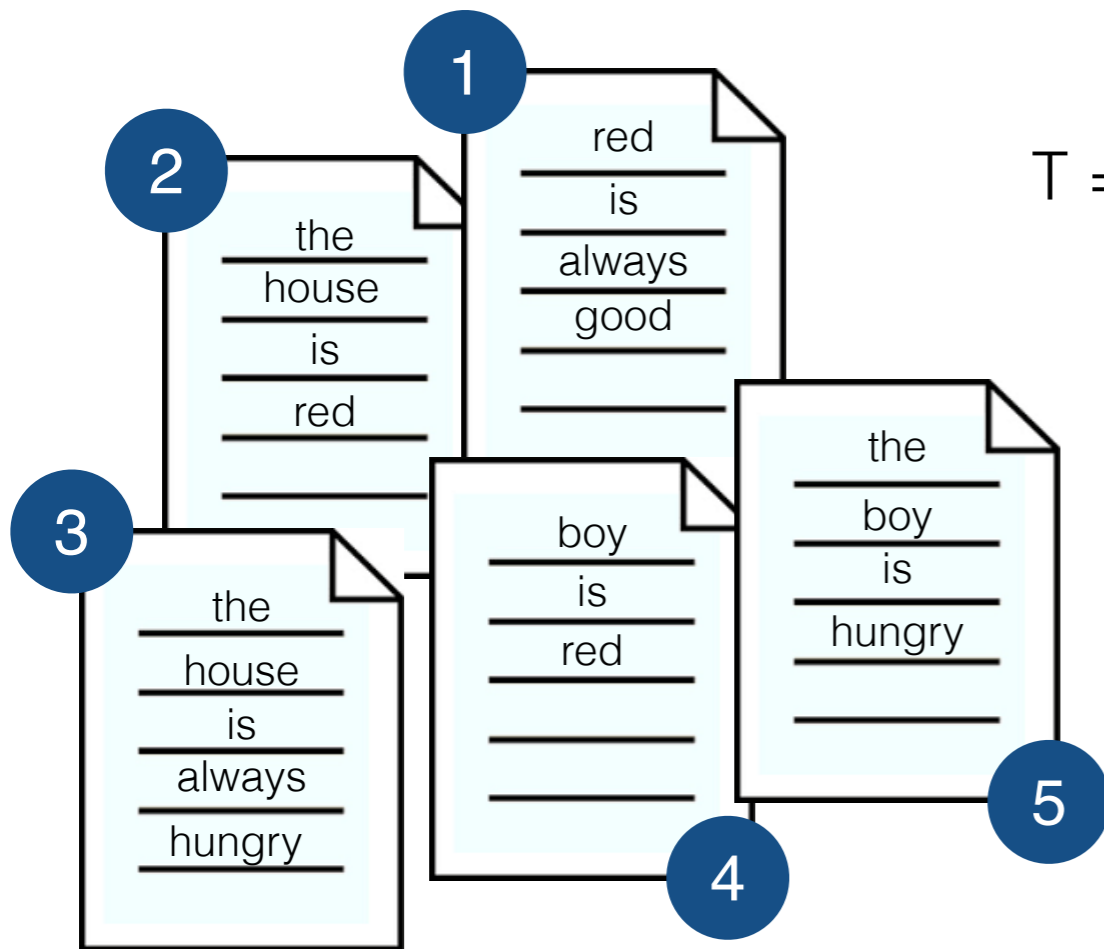


$t_1$      $t_2$      $t_3$      $t_4$      $t_5$      $t_6$      $t_7$      $t_8$   
 $T = \{\text{always, boy, good, house, hungry, is, red, the}\}$

- $L_{t_1} = [1, 3]$
- $L_{t_2} = [4, 5]$
- $L_{t_3} = [1]$
- $L_{t_4} = [2, 3]$
- $L_{t_5} = [3, 5]$
- $L_{t_6} = [1, 2, 3, 4, 5]$
- $L_{t_7} = [1, 2, 4]$
- $L_{t_8} = [2, 3, 5]$

$q = \{\text{boy, is, the}\}$

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$T = \{ \begin{matrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 & t_7 & t_8 \\ \text{always, boy, good, house, hungry, is, red, the} \end{matrix} \}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

$L_{t_5} = [3, 5]$

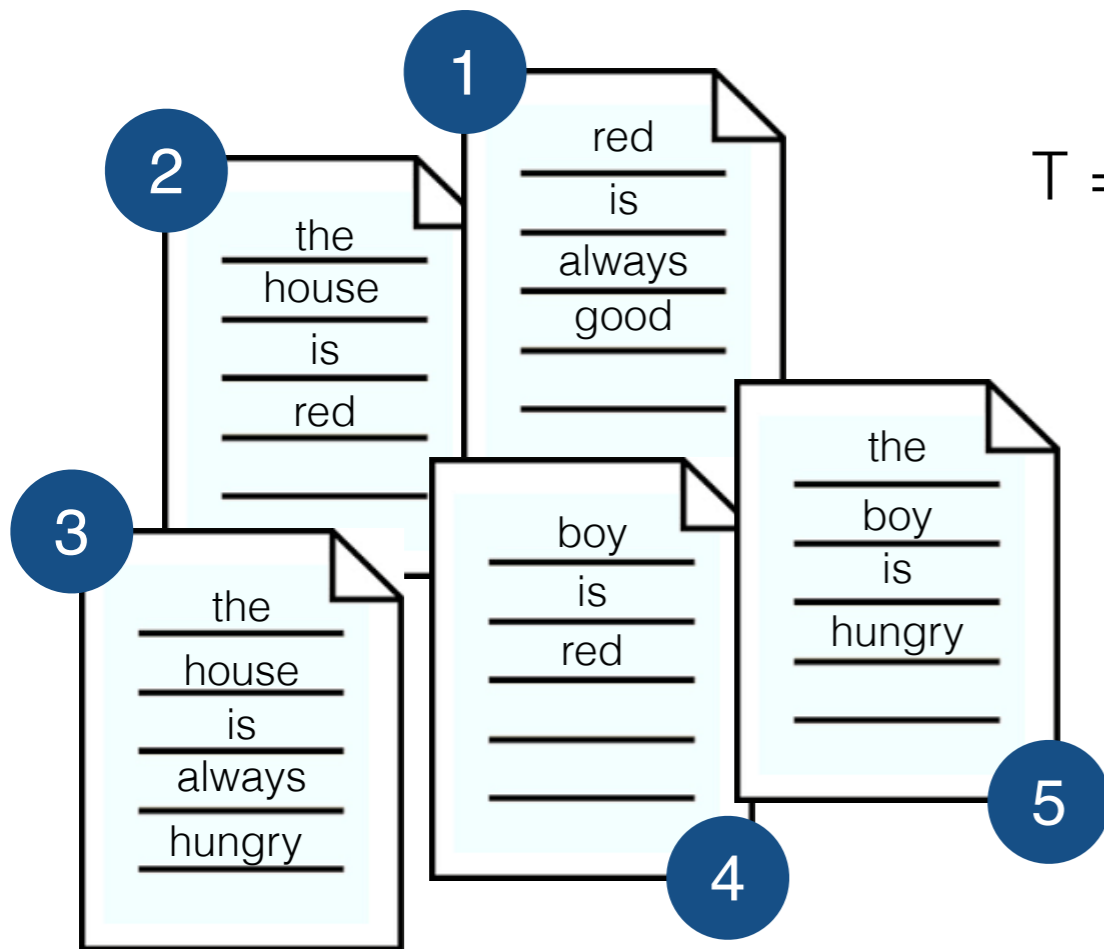
$L_{t_6} = [1, 2, 3, 4, 5]$

$L_{t_7} = [1, 2, 4]$

$L_{t_8} = [2, 3, 5]$

$q = \{\text{boy, is, the}\}$

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$T = \{ \begin{matrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 & t_7 & t_8 \\ \text{always, boy, good, house, hungry, is, red, the} \end{matrix} \}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

$L_{t_5} = [3, 5]$

$L_{t_6} = [1, 2, 3, 4, 5]$

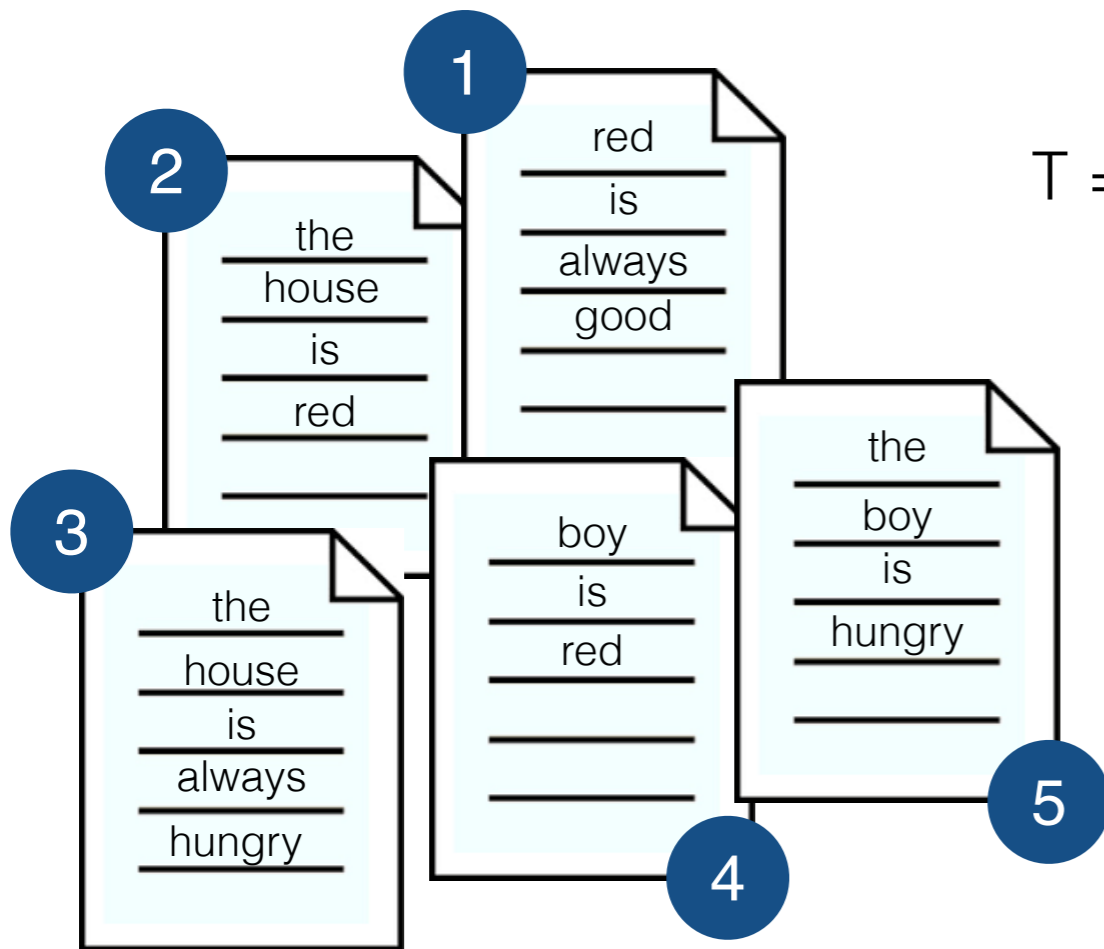
$L_{t_7} = [1, 2, 4]$

$L_{t_8} = [2, 3, 5]$

$q = \{\text{boy, is, the}\}$

$q = \{\text{good, hungry}\}$

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$t_1$     $t_2$     $t_3$     $t_4$     $t_5$     $t_6$     $t_7$     $t_8$   
 $T = \{\text{always, boy, good, house, hungry, is, red, the}\}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

$L_{t_5} = [3, 5]$

$L_{t_6} = [1, 2, 3, 4, 5]$

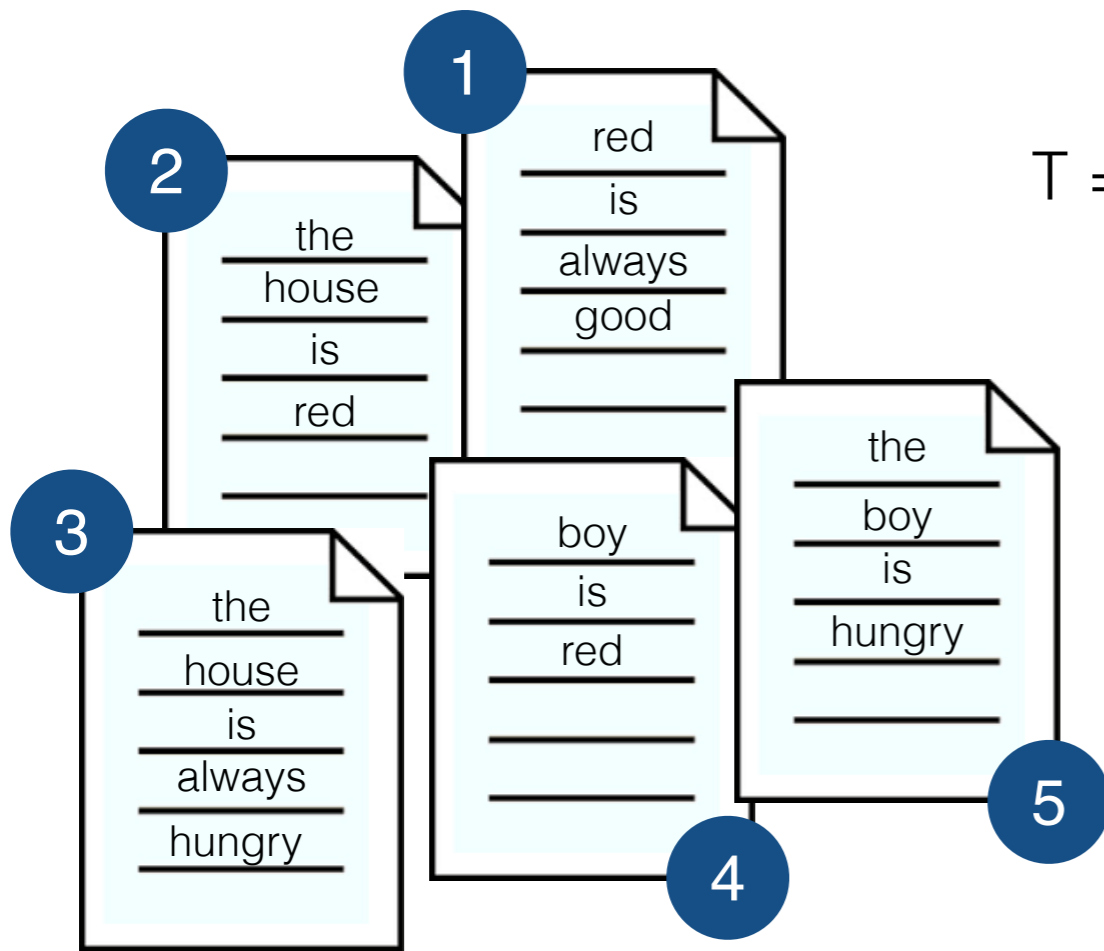
$L_{t_7} = [1, 2, 4]$

$L_{t_8} = [2, 3, 5]$

$q = \{\text{boy, is, the}\}$

$q = \{\text{good, hungry}\}$

Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$T = \{ \begin{matrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 & t_7 & t_8 \\ \text{always, boy, good, house, hungry, is, red, the} \end{matrix} \}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

$L_{t_5} = [3, 5]$

$L_{t_6} = [1, 2, 3, 4, 5]$

$L_{t_7} = [1, 2, 4]$

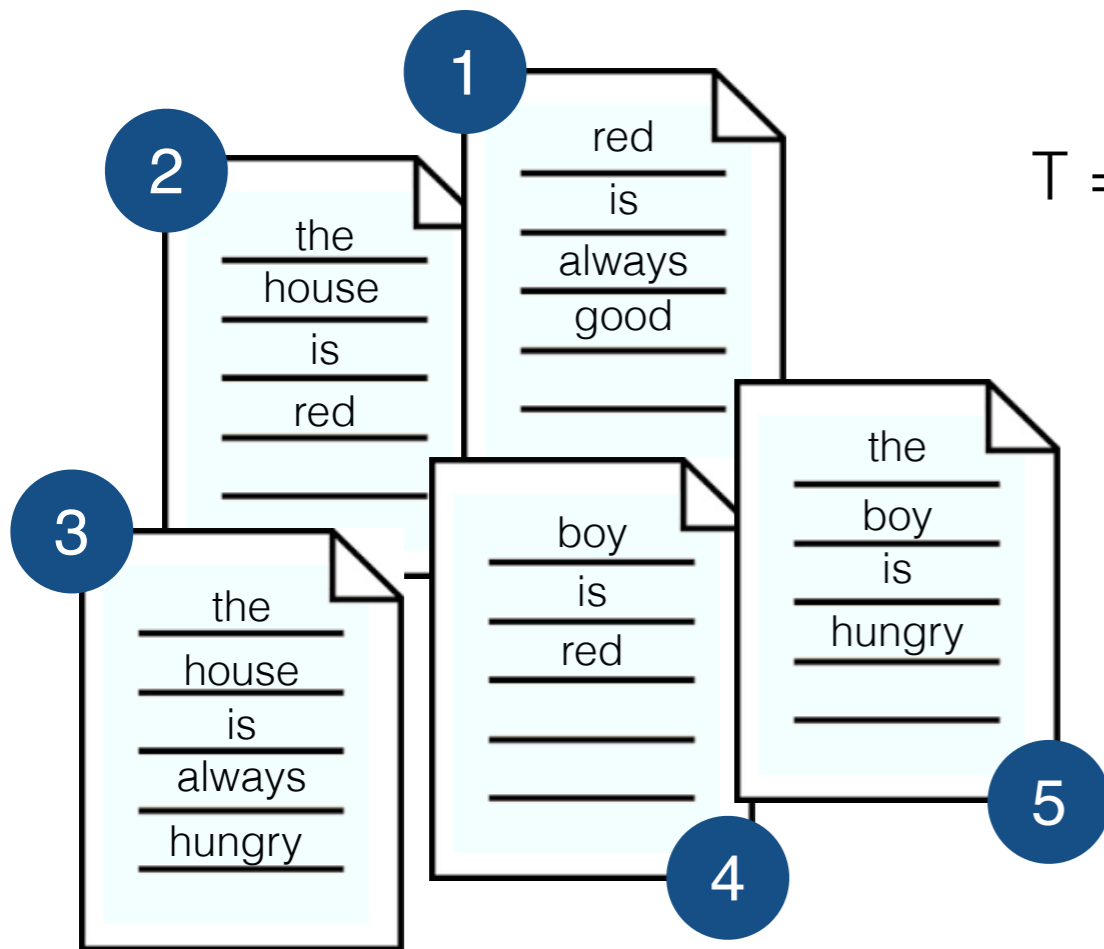
$L_{t_8} = [2, 3, 5]$

$q = \{\text{boy, is, the}\}$

$q = \{\text{good, hungry}\}$



Inverted Indexes owe their popularity to the *efficient resolution of queries*, such as: “return me all documents in which terms  $\{t_1, \dots, t_k\}$  occur”.



$T = \{ \begin{matrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 & t_7 & t_8 \\ \text{always, boy, good, house, hungry, is, red, the} \end{matrix} \}$

$L_{t_1} = [1, 3]$

$L_{t_2} = [4, 5]$

$L_{t_3} = [1]$

$L_{t_4} = [2, 3]$

$L_{t_5} = [3, 5]$

$L_{t_6} = [1, 2, 3, 4, 5]$

$L_{t_7} = [1, 2, 4]$

$L_{t_8} = [2, 3, 5]$

$q = \{\text{boy, is, the}\}$

$q = \{\text{good, hungry}\}$

**intersection**

# Genesis - 1970s



Peter Elias  
[1923 - 2001]



Robert Fano  
[1917 - 2016]

Robert Fano. *On the number of bits required to implement an associative memory*. Memorandum 61, Computer Structures Group, MIT (1971).

Peter Elias. *Efficient Storage and Retrieval by Content and Address of Static Files*. Journal of the ACM (JACM) 21, 2, 246–260 (1974).

# Genesis - 1970s



Peter Elias  
[1923 - 2001]



Robert Fano  
[1917 - 2016]

Robert Fano. *On the number of bits required to implement an associative memory*. Memorandum 61, Computer Structures Group, MIT (1971).

Peter Elias. *Efficient Storage and Retrieval by Content and Address of Static Files*. Journal of the ACM (JACM) 21, 2, 246–260 (1974).



Sebastiano Vigna. *Quasi-succinct indices*.

In Proceedings of the 6-th ACM International Conference on Web Search and Data Mining (WSDM), 83-92 (2013).

40 years later!

# Elias-Fano Encoding

3	1
4	2
7	3
13	4
14	5
15	6
21	7
43	8

# Elias-Fano Encoding

	3	1
	4	2
	7	3
	13	4
	14	5
	15	6
	21	7
u =	43	8

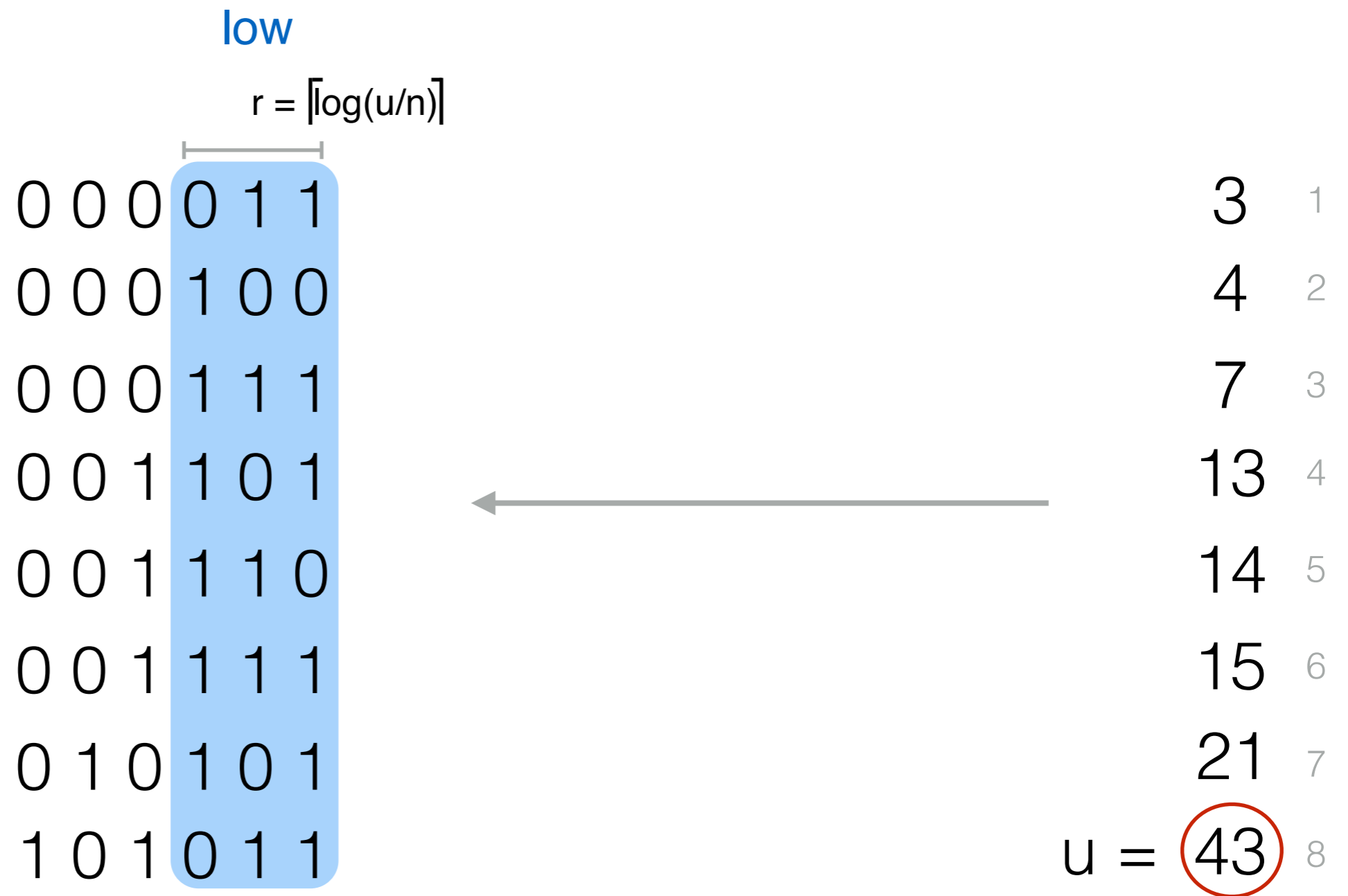
# Elias-Fano Encoding

0 0 0 0 1 1  
0 0 0 1 0 0  
0 0 0 1 1 1  
0 0 1 1 0 1  
0 0 1 1 1 0  
0 0 1 1 1 1  
0 1 0 1 0 1  
1 0 1 0 1 1

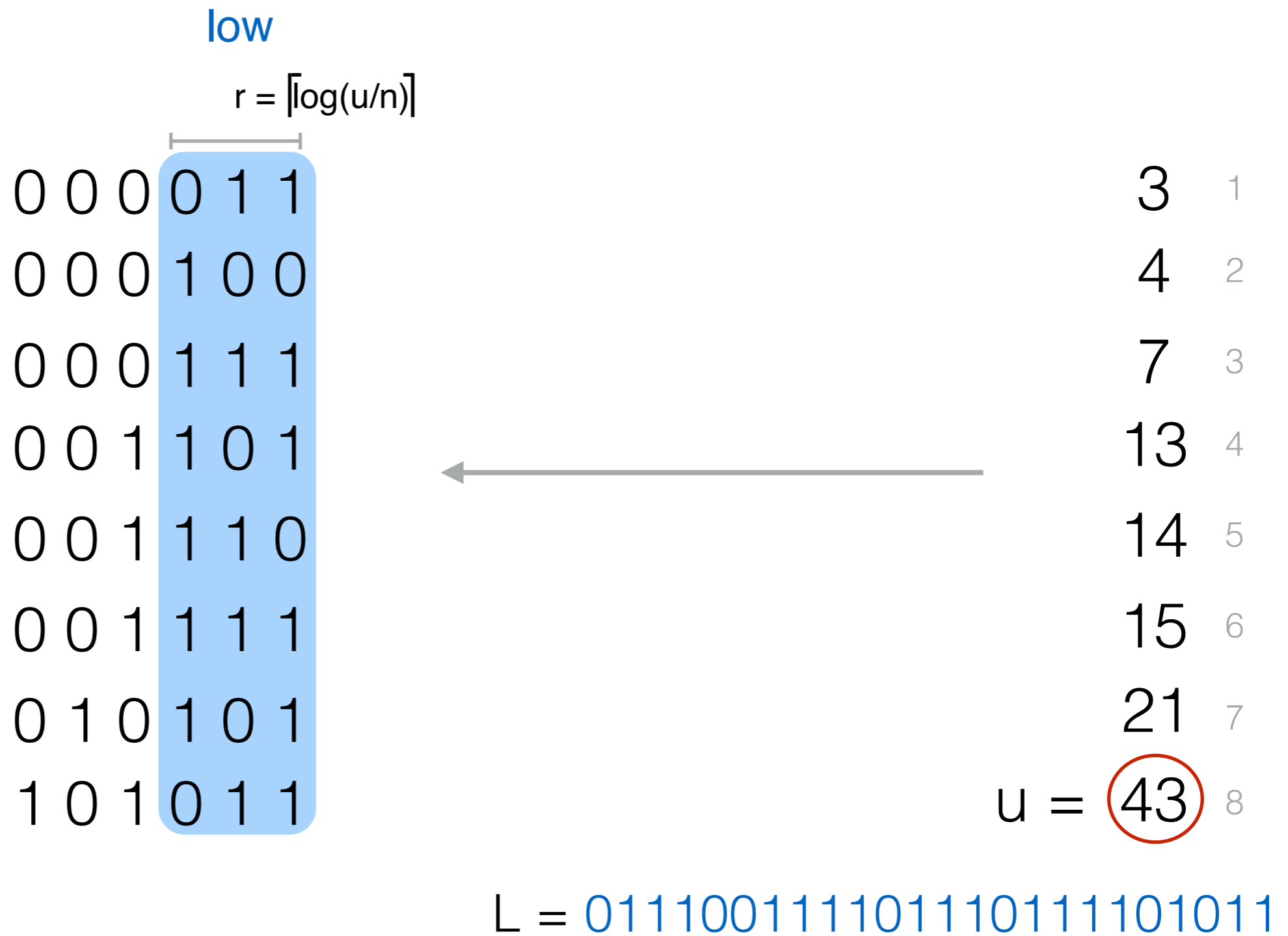


3 1  
4 2  
7 3  
13 4  
14 5  
15 6  
21 7  
u = 43 8

# Elias-Fano Encoding

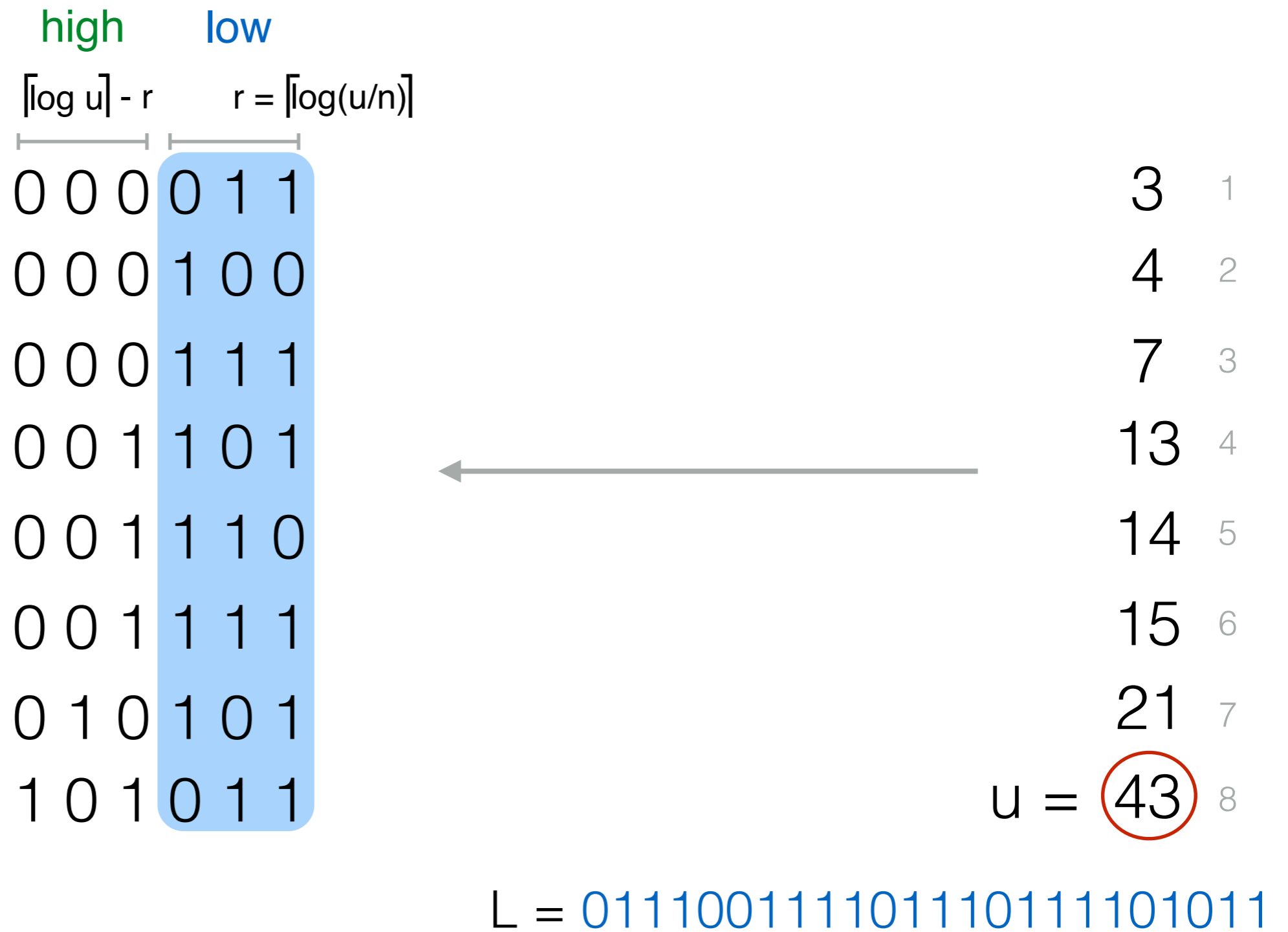


# Elias-Fano Encoding

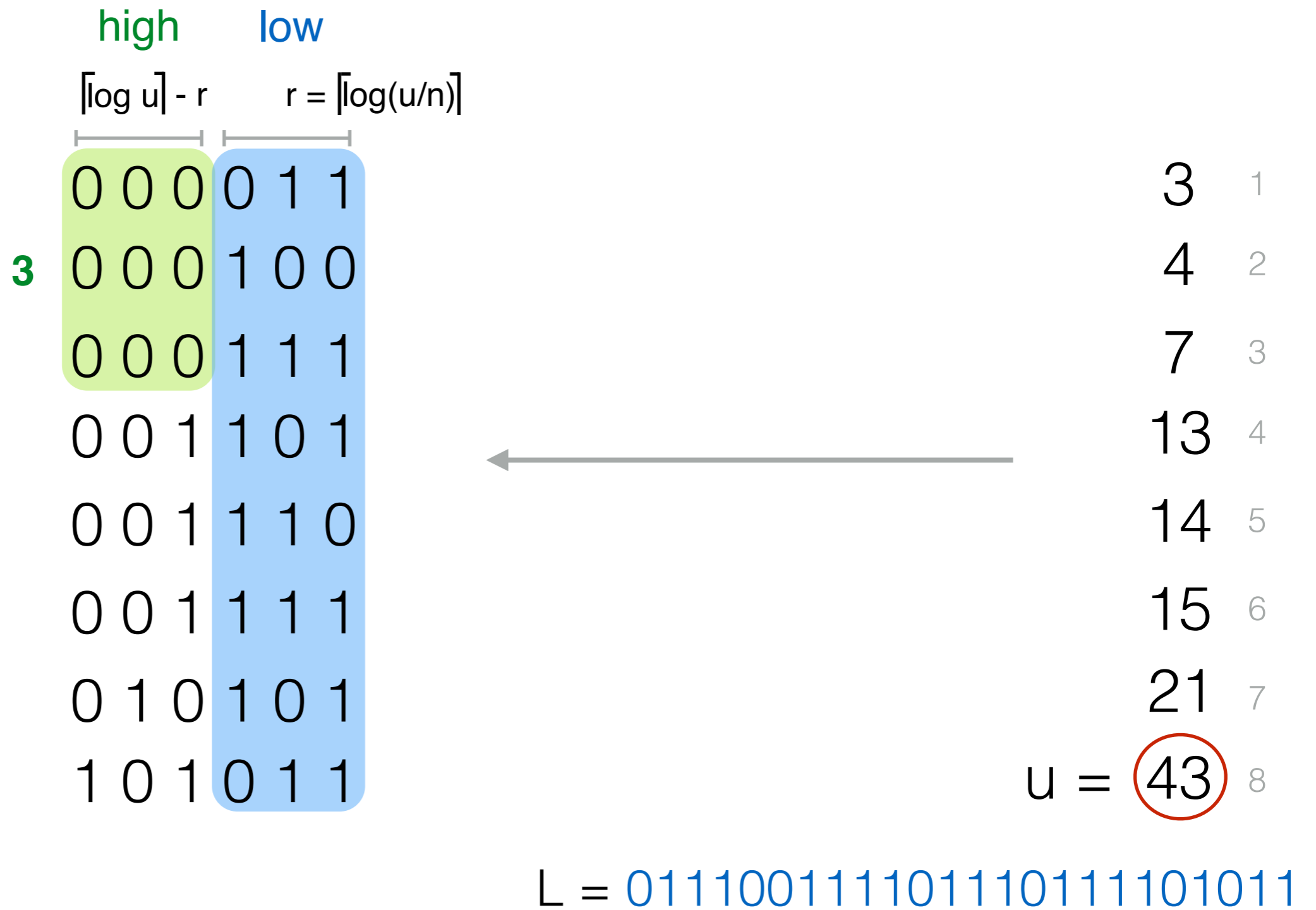




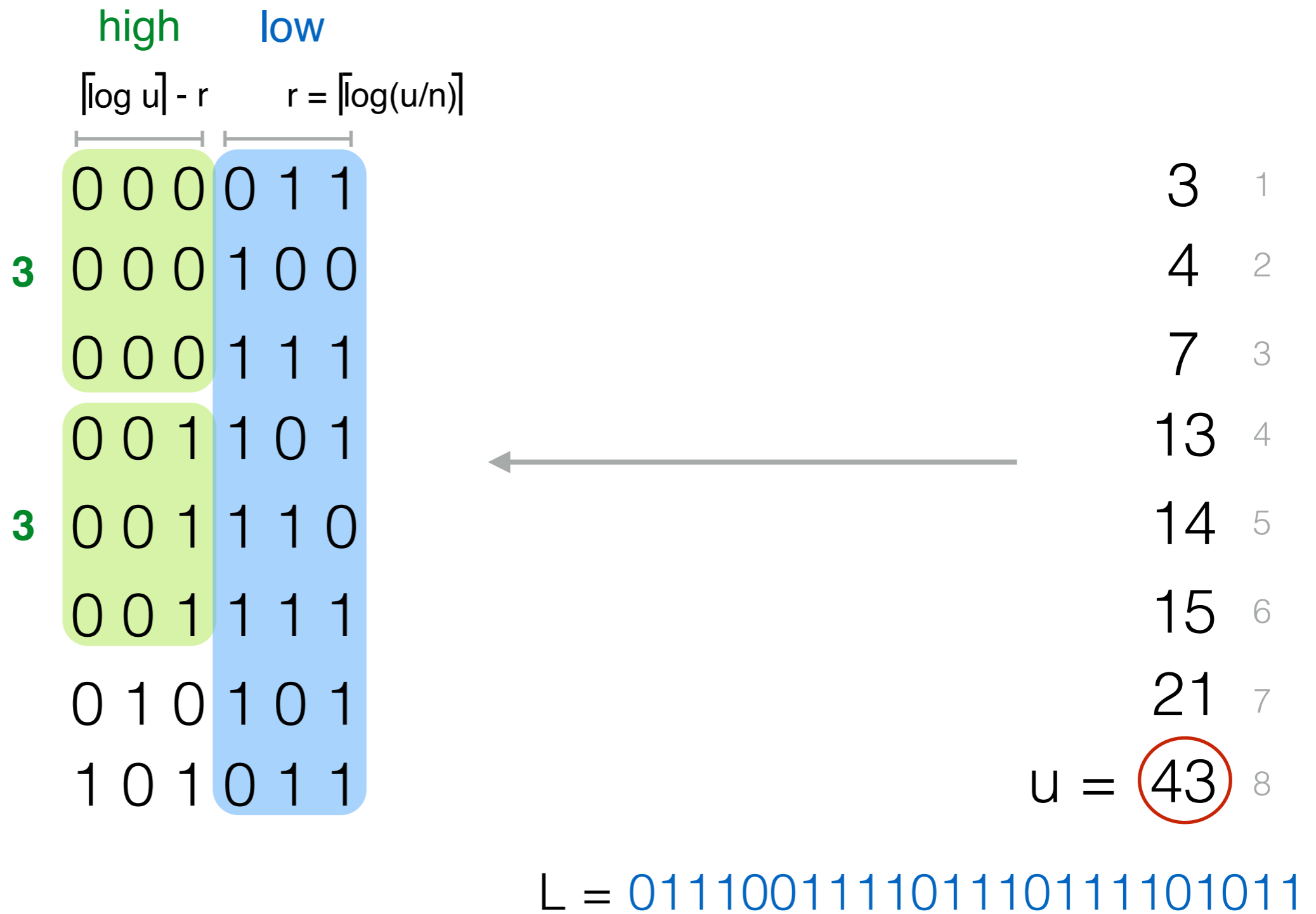
# Elias-Fano Encoding



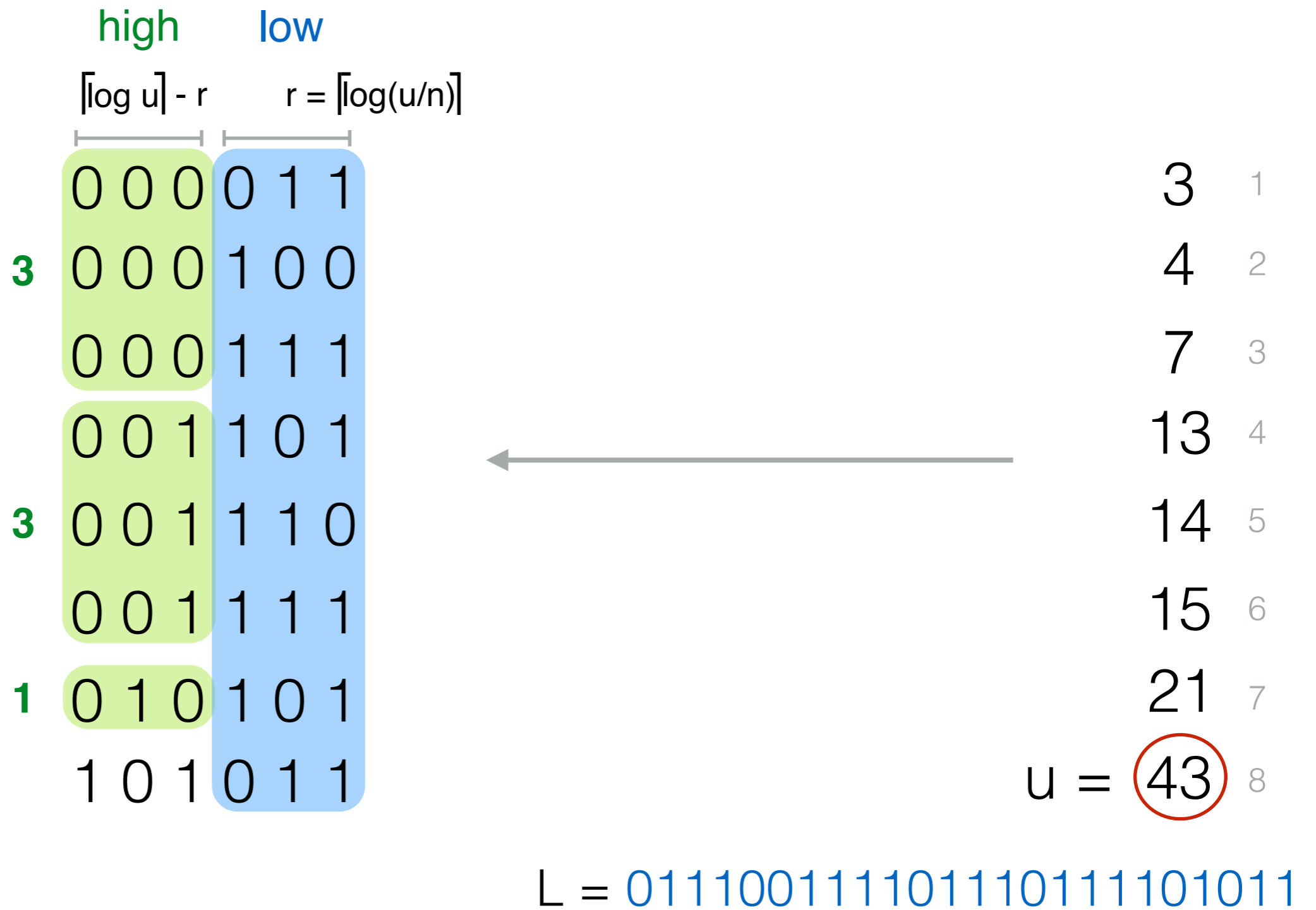
# Elias-Fano Encoding



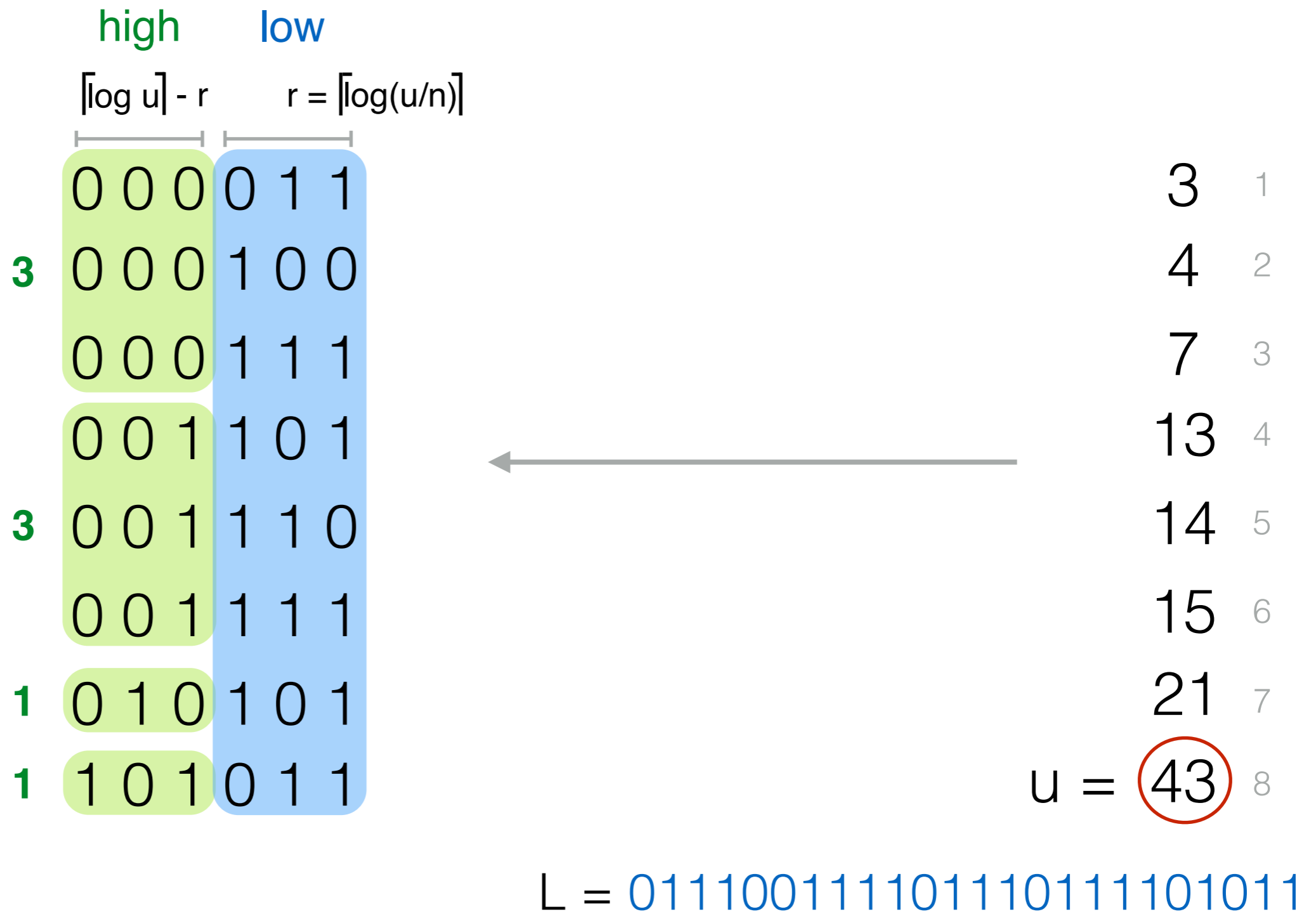
# Elias-Fano Encoding



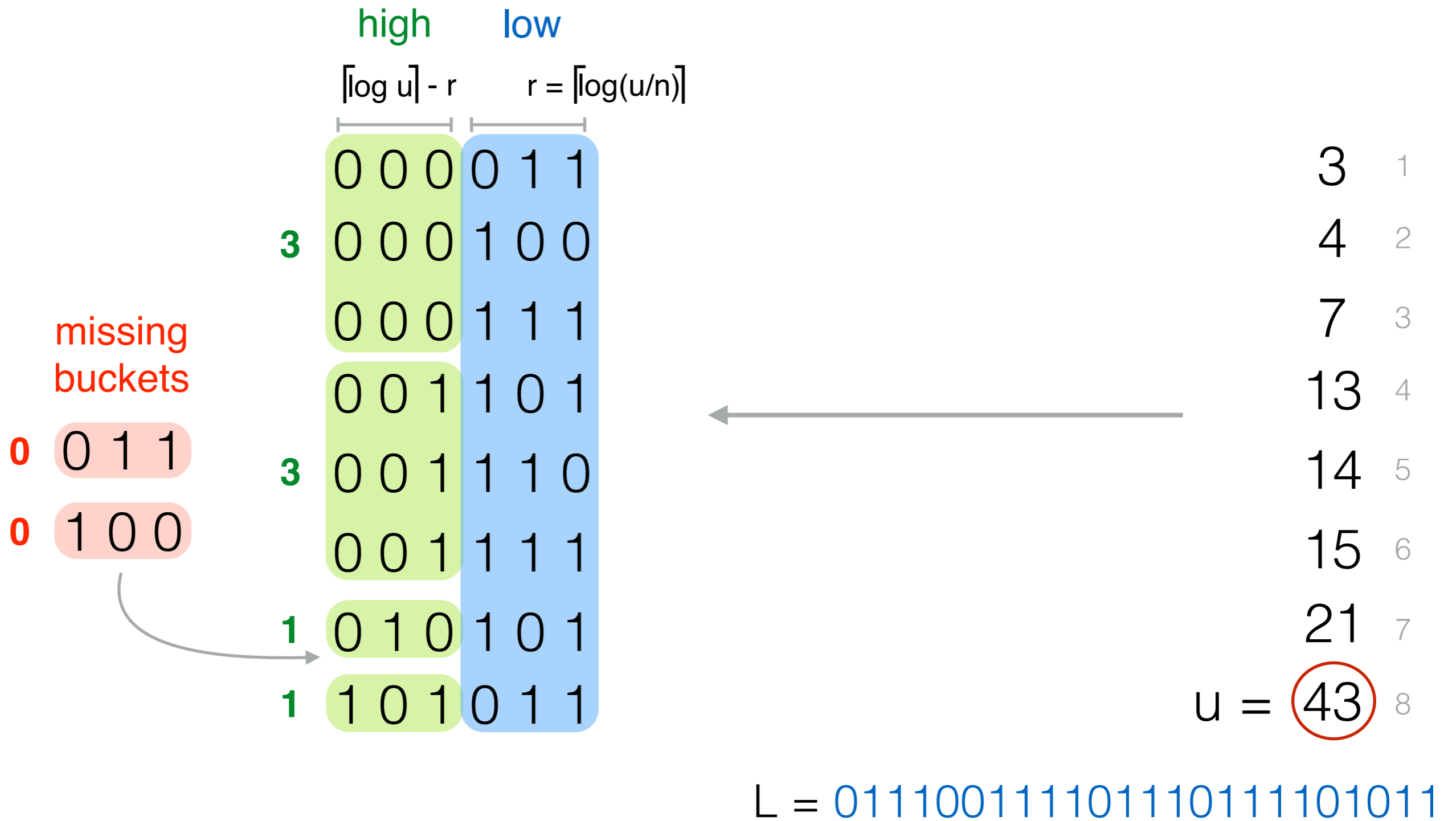
# Elias-Fano Encoding



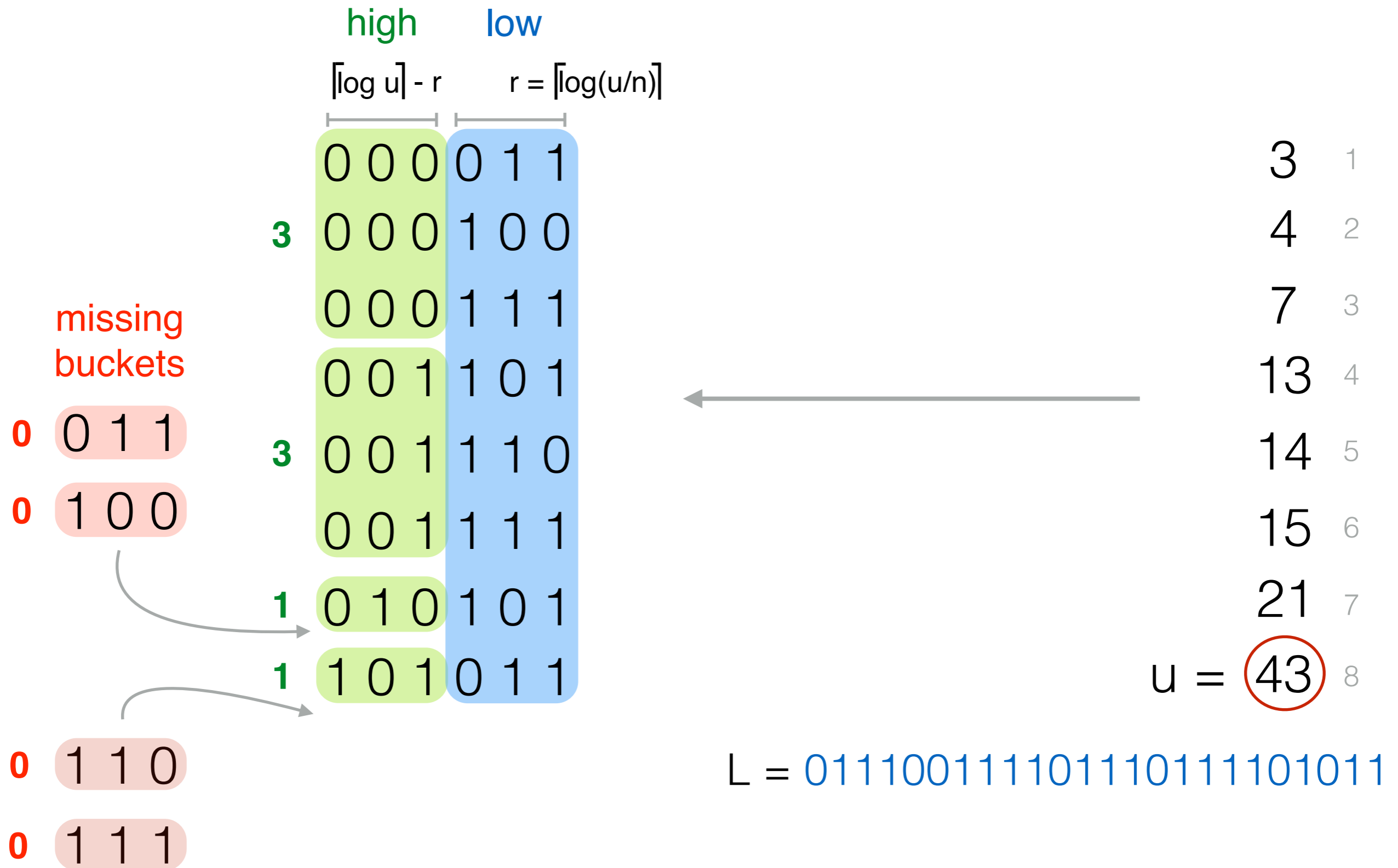
# Elias-Fano Encoding



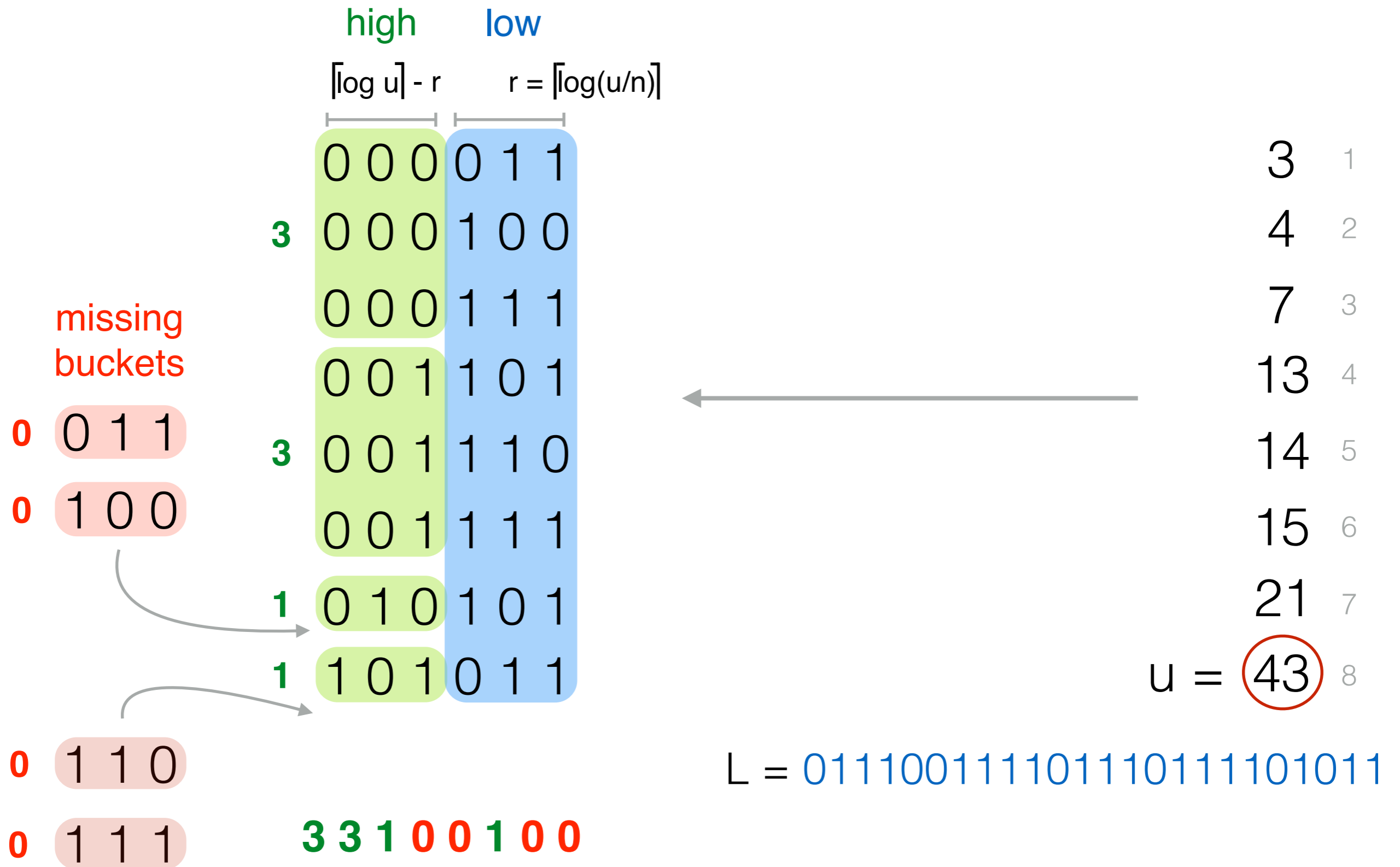
# Elias-Fano Encoding



# Elias-Fano Encoding

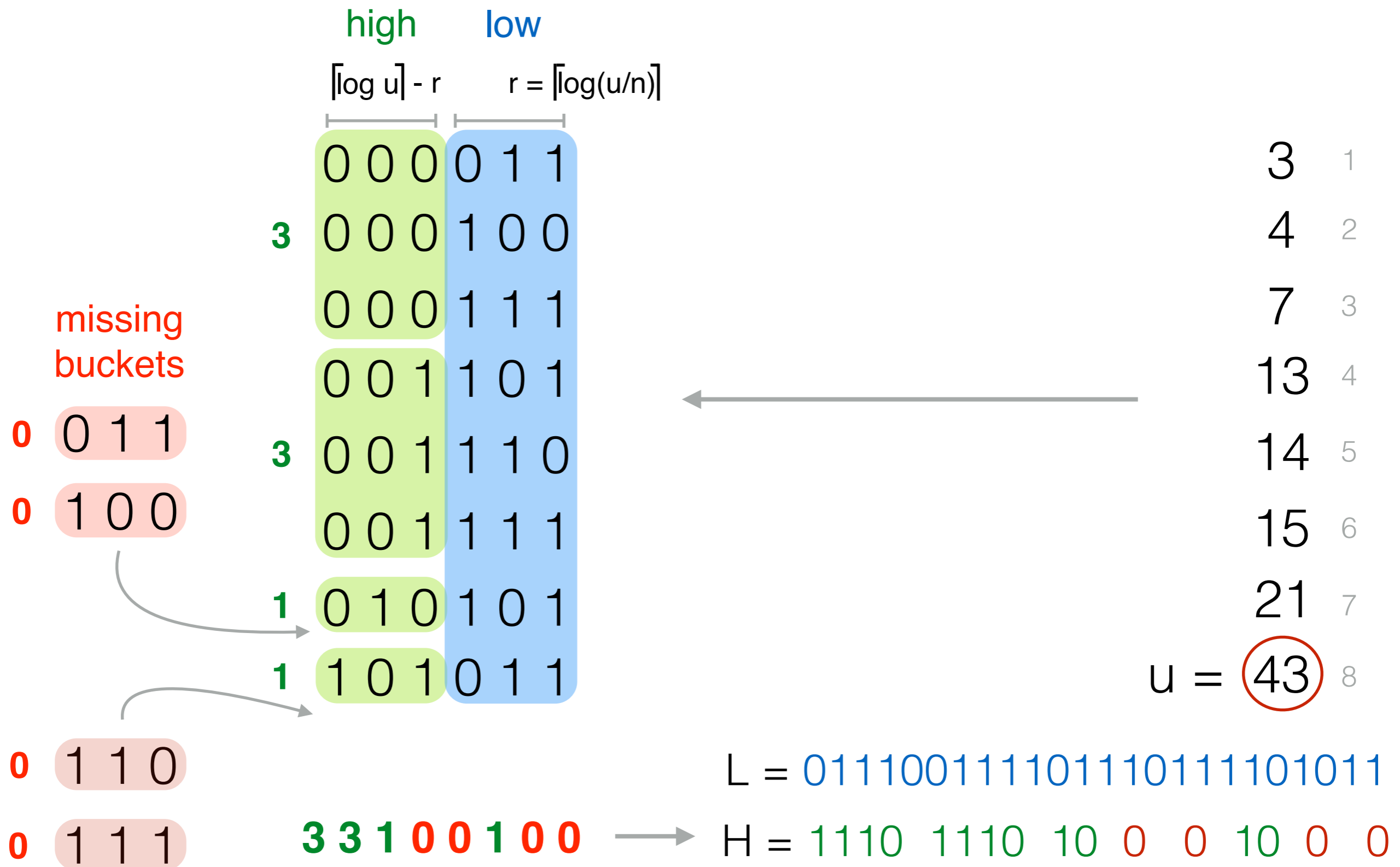


# Elias-Fano Encoding





# Elias-Fano Encoding



$$EF(S[0,n)) = ?$$

$$EF(S[0,n)) = ?$$

 $\lceil \log(u/n) \rceil$ 

L = 011100111101110111101011

H = 1110 1110 10 0 0 10 0 0

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil$$

$\lceil \log(u/n) \rceil$

L = 011100111101110111101011

H = 1110 1110 10 0 0 10 0 0

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil$$

$\lceil \log(u/n) \rceil$

L = 011100111101110111101011

H = 1110 1110 10 0 0 10 0 0

n ones

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil$$

$\lceil \log(u/n) \rceil$   
L = 011100111101110111101011  
H = 1110 1110 10 0 0 10 0 0  
n ones

We store a 0 whenever we change bucket.

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil$$

$\overbrace{\quad}^{\lceil \log(u/n) \rceil}$   
L = 011100111101110111101011  
H = 1110 1110 10 0 0 10 0 0

n ones  
 $2^{\lceil \log n \rceil}$  zeros

We store a 0 whenever we change bucket.

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

$\overbrace{\quad}^{\lceil \log(u/n) \rceil}$   
L = 011100111101110111101011  
H = 1110 1110 10 0 0 10 0 0

We store a 0 whenever we change bucket.

n ones  
 $2^{\lceil \log n \rceil}$  zeros



$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

$\overbrace{\quad}^{\lceil \log(u/n) \rceil}$   
L = 011100111101110111101011  
H = 1110 1110 10 0 0 10 0 0

n ones

$2^{\lceil \log n \rceil}$  zeros

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

$\overbrace{\quad}^{\lceil \log(u/n) \rceil}$   
L = 011100111101110111101011  
H = 1110 1110 10 0 0 10 0 0

n ones

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

$\lceil \log(u/n) \rceil$   
┌───┐  
L = 011100111101110111101011  
H = 1110 1110 10 0 0 10 0 0

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

00000000000000000000



$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

000**1**0000000000000000

3

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

000**1**00**1**000000000000

3

6

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

000**1**00**1**000**1**00000000

3

6

10

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

000**1**00**1**000**11**0000000

3

6

1011

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

000**1**00**1**000**11**00000**1**  
3      6      1011      17

$$EF(S[0,n]) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| ?$$

000**1**00**1**000**11**00000**1**  
3      6      1011      17

With possible repetitions!  
(*weak* monotonicity)

$$EF(S[0,n]) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| = \binom{u+n}{n}$$

000**1**00**1**000**11**00000**1**  
3      6      1011      17

With possible repetitions!  
(*weak* monotonicity)

$$EF(S[0,n)) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| = \binom{u+n}{n}$$

$$\left\lceil \log \binom{u+n}{n} \right\rceil \approx n \log \frac{u+n}{n}$$

000**1**00**1**000**11**000000**1**  
3          6          1011          17

With possible repetitions!  
(*weak* monotonicity)



$$EF(S[0,n]) = n \left\lceil \log \frac{u}{n} \right\rceil + 2n \text{ bits}$$

Is it good or not?

(less than half a bit away [Elias-1974])

## Information Theoretic Lower Bound

The minimum number of bits needed to describe a set  $\mathcal{X}$  is

$$\left\lceil \log |\mathcal{X}| \right\rceil \text{ bits.}$$

$\mathcal{X}$  is the set of all monotone sequence of length  $n$  drawn from a universe  $u$ .

$$|\mathcal{X}| = \binom{u+n}{n}$$

000**1**00**1**000**11**00000**1**  
3      6      1011      17

With possible repetitions!  
(*weak* monotonicity)

$$\left\lceil \log \binom{u+n}{n} \right\rceil \approx n \log \frac{u+n}{n}$$

**optimal**



**Access** to each  $S[i]$  in  $O(1)$  worst-case

**Access** to each  $S[i]$  in  $O(1)$  worst-case

**Predecessor**( $x$ ) =  $\max\{S[i] \mid S[i] < x\}$

**Successor**( $x$ ) =  $\min\{S[i] \mid S[i] \geq x\}$

queries in  $O\left(\log \frac{u}{n}\right)$  worst-case

Access to each  $S[i]$  in  $O(1)$  worst-case

$\text{Predecessor}(x) = \max\{S[i] \mid S[i] < x\}$

$\text{Successor}(x) = \min\{S[i] \mid S[i] \geq x\}$

queries in  $O\left(\log \frac{u}{n}\right)$  worst-case

# Rank/Select on bitmaps

## Definition

Given a bitvector  $B$  of  $n$  bits:

$\text{Rank}_{0/1}(i) = \#$  of 0/1 in  $B[0, i)$

$\text{Select}_{0/1}(i) =$  position of  $i$ -th 0/1

# Rank/Select on bitmaps

## Definition

Given a bitvector  $B$  of  $n$  bits:

$\text{Rank}_{0/1}(i) = \#$  of 0/1 in  $B[0, i)$

$\text{Select}_{0/1}(i) =$  position of  $i$ -th 0/1

## Examples

$B = 101011010101111010110101$

# Rank/Select on bitmaps

## Definition

Given a bitvector  $B$  of  $n$  bits:

$\text{Rank}_{0/1}(i) = \# \text{ of } 0/1 \text{ in } B[0, i)$

$\text{Select}_{0/1}(i) = \text{position of } i\text{-th } 0/1$

## Examples

$B = 101011010101111010110101$

$\text{Rank}_0(5) = 2$



# Rank/Select on bitmaps

## Definition

Given a bitvector  $B$  of  $n$  bits:

$\text{Rank}_{0/1}(i) = \#$  of 0/1 in  $B[0, i)$

$\text{Select}_{0/1}(i) =$  position of  $i$ -th 0/1

## Examples

$B = 101011010101111010110101$

$\text{Rank}_0(5) = 2$

$\text{Rank}_1(7) = 4$

# Rank/Select on bitmaps

## Definition

Given a bitvector  $B$  of  $n$  bits:

$\text{Rank}_{0/1}(i) = \#$  of 0/1 in  $B[0, i)$

$\text{Select}_{0/1}(i) =$  position of  $i$ -th 0/1

## Examples

$B = 101011010101111010110101$

$\text{Rank}_0(5) = 2$      $\text{Select}_0(5) = 10$

$\text{Rank}_1(7) = 4$

# Rank/Select on bitmaps

## Definition

Given a bitvector  $B$  of  $n$  bits:

$\text{Rank}_{0/1}(i) = \#$  of 0/1 in  $B[0, i)$

$\text{Select}_{0/1}(i) =$  position of  $i$ -th 0/1

## Examples

$B = 101011010101111010110101$

$\text{Rank}_0(5) = 2$        $\text{Select}_0(5) = 10$

$\text{Rank}_1(7) = 4$        $\text{Select}_1(7) = 11$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$

1 2 3 4 5 6 7 8

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = ?$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = ?$

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = ?$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = ?$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$



# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = ?$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i)$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = ?$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Rank}_0(\text{Select}_1(i))$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001000$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Rank}_0(\text{Select}_1(i))$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001000$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$$\begin{aligned} \text{Access}(i) &= \text{Rank}_0(\text{Select}_1(i)) \\ &= \text{Select}_1(i) - i \end{aligned}$$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001000$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

$\text{Access}(7) = S[7] = ?$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$



# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

$\text{Access}(7) = S[7] = ?$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

$\text{Access}(7) = S[7] = 010000$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

$\text{Access}(7) = S[7] = 010\ 101$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$

# Random Access

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{Access}(4) = S[4] = 001\ 101$

$\text{Access}(7) = S[7] = 010\ 101$

Recall: we store a 0 whenever we change bucket.

$H = 1110111010001000$

$L = 011100111101110111101011$

$r = \lceil \log(u/n) \rceil$

Complexity:  $O(1)$

$\text{Access}(i) = \text{Select}_1(i) - i \ll r \mid L[(i-1)r, ir)$

# Available Implementations

Library	Author(s)	Link	Language
folly	Facebook, Inc.	<a href="https://github.com/facebook/folly">https://github.com/facebook/folly</a>	C++
sdsl	Simon Gog	<a href="https://github.com/simongog/sdsl-lite">https://github.com/simongog/sdsl-lite</a>	C++
ds2i	Giuseppe Ottaviano Rossano Venturini Nicola Tonello	<a href="https://github.com/ot/ds2i">https://github.com/ot/ds2i</a>	C++
Sux	Sebastiano Vigna	<a href="http://sux.di.unimi.it">http://sux.di.unimi.it</a>	Java/C++

# Killer applications

## 1. Inverted Indexes

# Killer applications

1. Inverted Indexes

2. Social Networks

# Killer applications

## 1. Inverted Indexes

## 2. Social Networks

### **Unicorn: A System for Searching the Social Graph**

Michael Curtiss, Iain Becker, Tudor Bosman, Sergey Doroshenko,  
Lucian Grijincu, Tom Jackson, Sandhya Kunnatur, Soren Lassen, Philip Pronin,  
Sriram Sankar, Guanghao Shen, Gintaras Woss, Chao Yang, Ning Zhang

Facebook, Inc.

#### **ABSTRACT**

Unicorn is an online, in-memory social graph-aware indexing system designed to search trillions of edges between tens of billions of users and entities on thousands of commodity servers. Unicorn is based on standard concepts in information retrieval. Unicorn is based on standard concepts in information retrieval. Unicorn is based on standard concepts in information retrieval.

rative of the evolution of Unicorn's architecture, as well as documentation for the major features and components of the system.

To the best of our knowledge, no other online graph retrieval system has ever been built with the scale of Unicorn. To the best of our knowledge, no other online graph retrieval system has ever been built with the scale of Unicorn. To the best of our knowledge, no other online graph retrieval system has ever been built with the scale of Unicorn.



# Killer applications

## 1. Inverted Indexes

## 2. Social Networks

### Unicorn: A System for Searching the Social Graph

Michael Curtiss, Iain Becker, Tudor Bosman, Sergey Doroshenko,  
Lucian Grijincu, Tom Jackson, Sandhya Kunnatur, Soren Lassen, Philip Pronin,  
Sriram Sankar, Guanghao Shen, Gintaras Woss, Chao Yang, Ning Zhang

Facebook, Inc.

#### ABSTRACT

Unicorn is an online, in-memory social graph-aware indexing system designed to search trillions of edges between tens of billions of users and entities on thousands of commodity servers. Unicorn is based on standard concepts in information retrieval. Unicorn is based on standard concepts in information retrieval. Unicorn is based on standard concepts in information retrieval.

rative of the evolution of Unicorn's architecture, as well as documentation for the major features and components of the system.

To the best of our knowledge, no other online graph retrieval system has ever been built with the scale of Unicorn. To the best of our knowledge, no other online graph retrieval system has ever been built with the scale of Unicorn. To the best of our knowledge, no other online graph retrieval system has ever been built with the scale of Unicorn.

# Killer applications

## 1. Inverted Indexes

## 2. Social Networks

### Unicorn: A System for Searching the Social Graph

Michael Curtiss, Iain Becker, Tudor Bosman, Sergey Doroshenko,  
Lucian Grijincu, Tom Jackson, Sandhya Kunnatur, Soren Lassen, Philip Pronin,  
Sriram Sankar, Guanghao Shen, Gintaras Woss, Chao Yang, Ning Zhang

Facebook, Inc.

#### ABSTRACT

Unicorn is an online, in-memory social graph-aware indexing system designed to search trillions of edges between tens of billions of users and entities on thousands of commodity servers. Unicorn is based on standard concepts in information retrieval. Unicorn is based on standard concepts in information retrieval. Unicorn is based on standard concepts in information retrieval.

#### Open Source

All Unicorn index server and aggregator code is written in C++. Unicorn relies extensively on modules in Facebook's "Folly" Open Source Library [5]. As part of the effort of releasing Graph Search, we have open-sourced a C++ implementation of the Elias-Fano index representation [31] as part of Folly.

# Killer applications

1. Inverted Indexes
2. Social Networks
3. Compressed Tries for N-Grams

# ***N*-grams - Introduction**

Strings of  $N$  words.

$N$  typically ranges from 1 to 5.

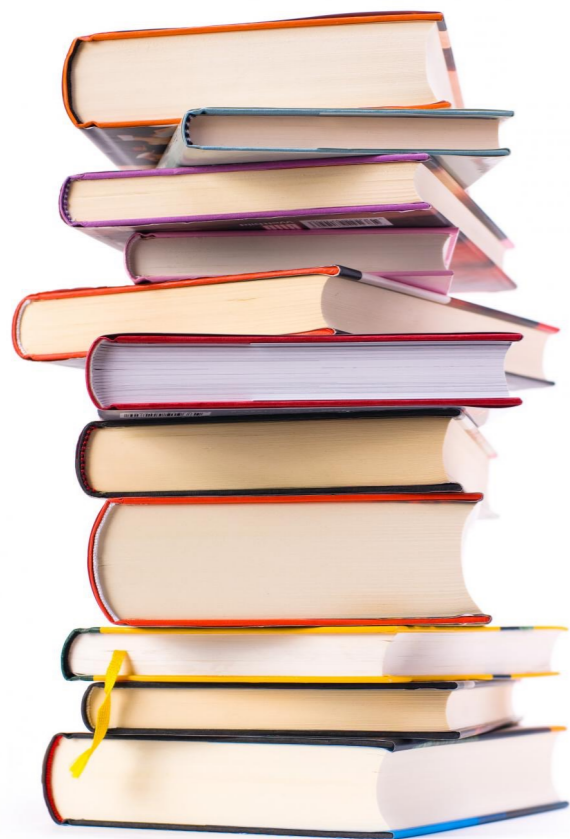
Extracted from text using a *sliding window* approach.

# *N*-grams - Introduction

Strings of  $N$  words.

$N$  typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.

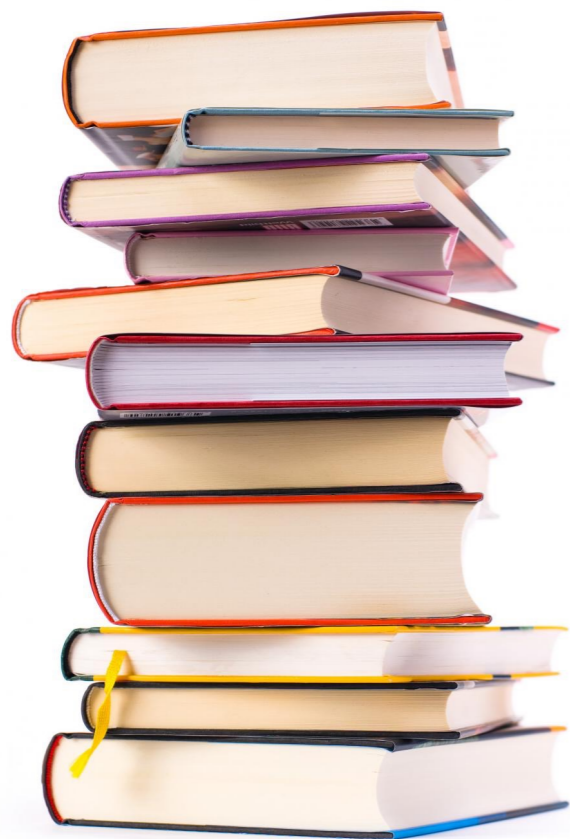


# N-grams - Introduction

Strings of  $N$  words.

$N$  typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.



Google Books

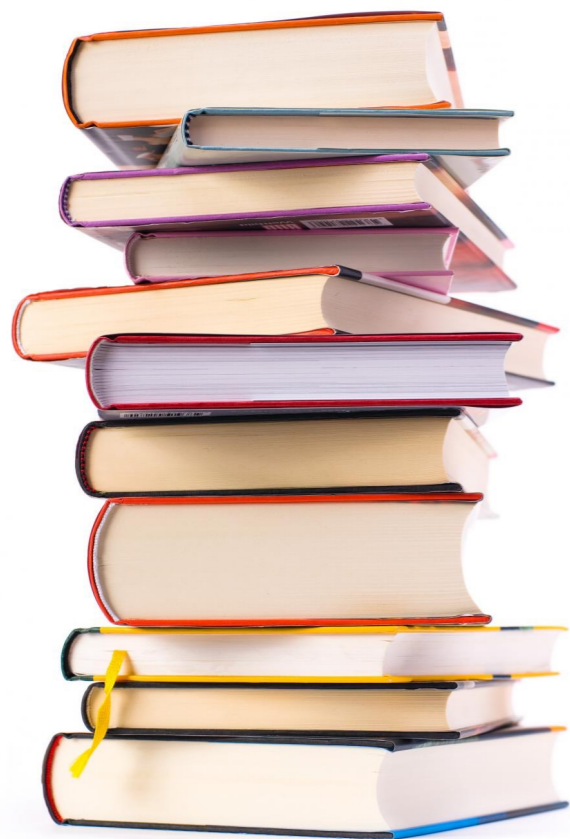
≈ 6% of the books ever published

# N-grams - Introduction

Strings of  $N$  words.

$N$  typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.



## Google Books

≈ 6% of the books ever published

$N$	number of grams
1	24,359,473
2	667,284,771
3	7,397,041,901
4	1,644,807,896
5	1,415,355,596

More than 11 billion grams.

# *N*-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



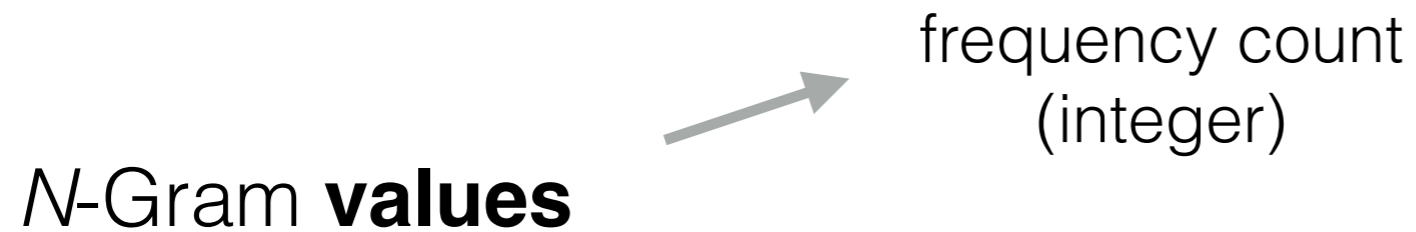
# *N*-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.

*N*-Gram **values**

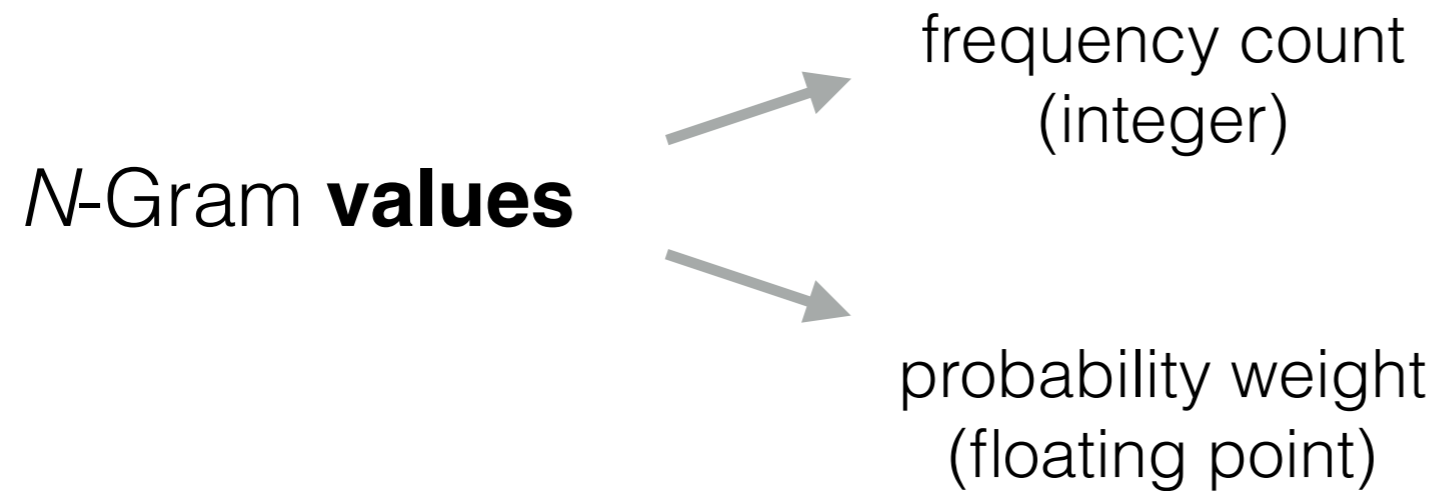
# *N*-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



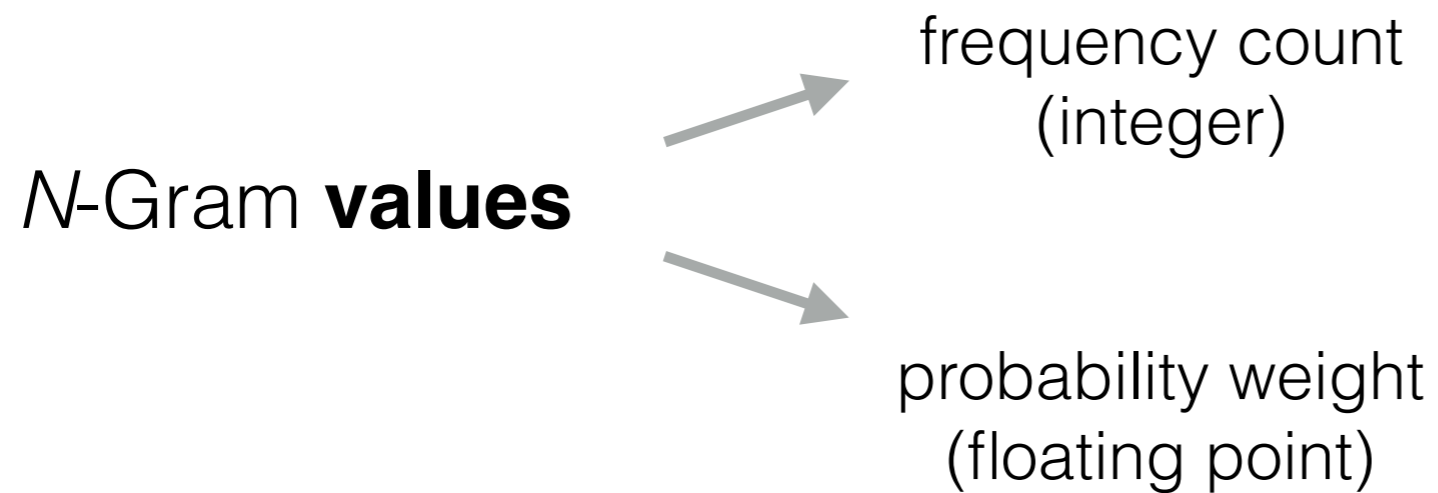
# *N*-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



# *N*-grams - Challenge

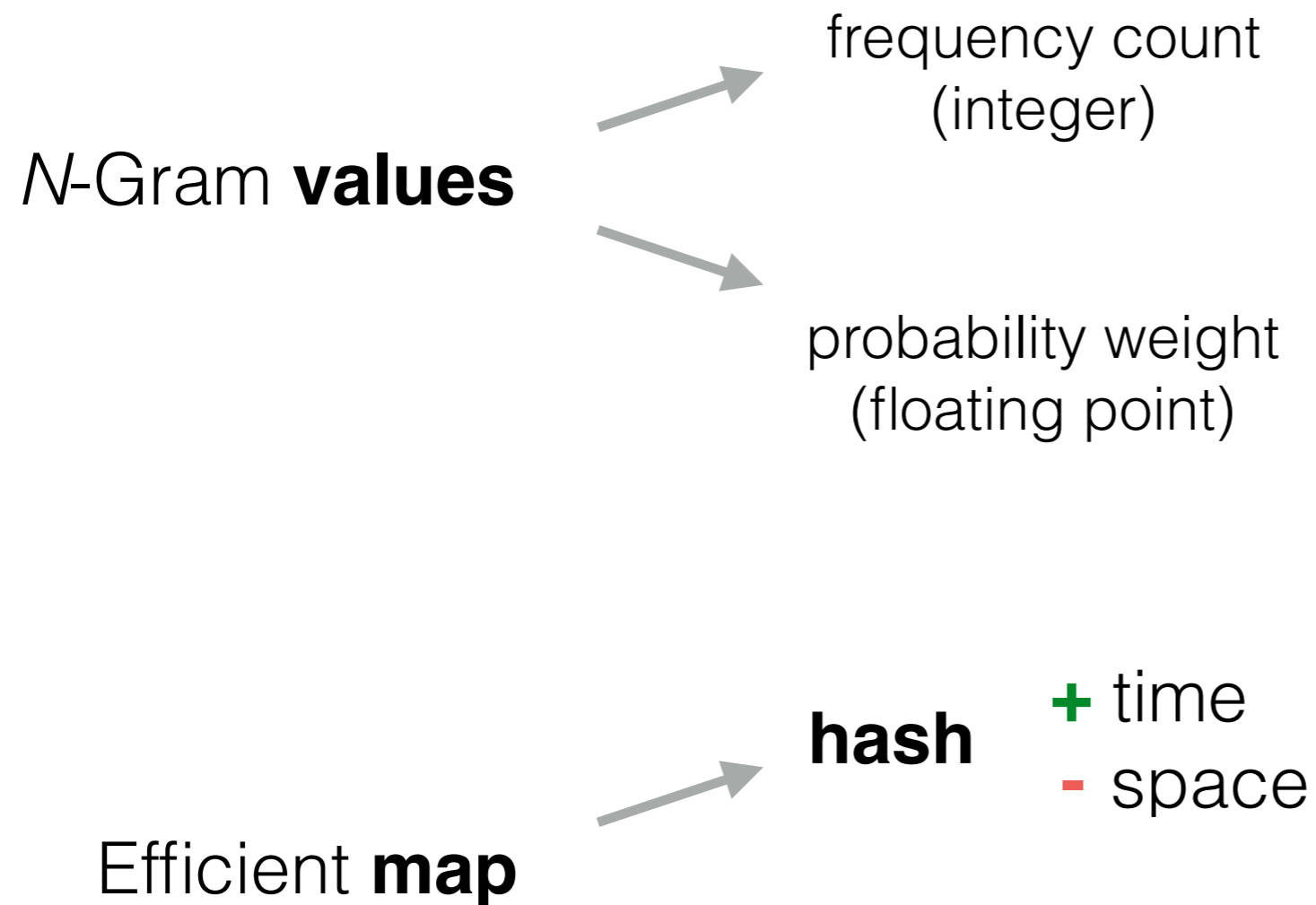
Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



Efficient **map**

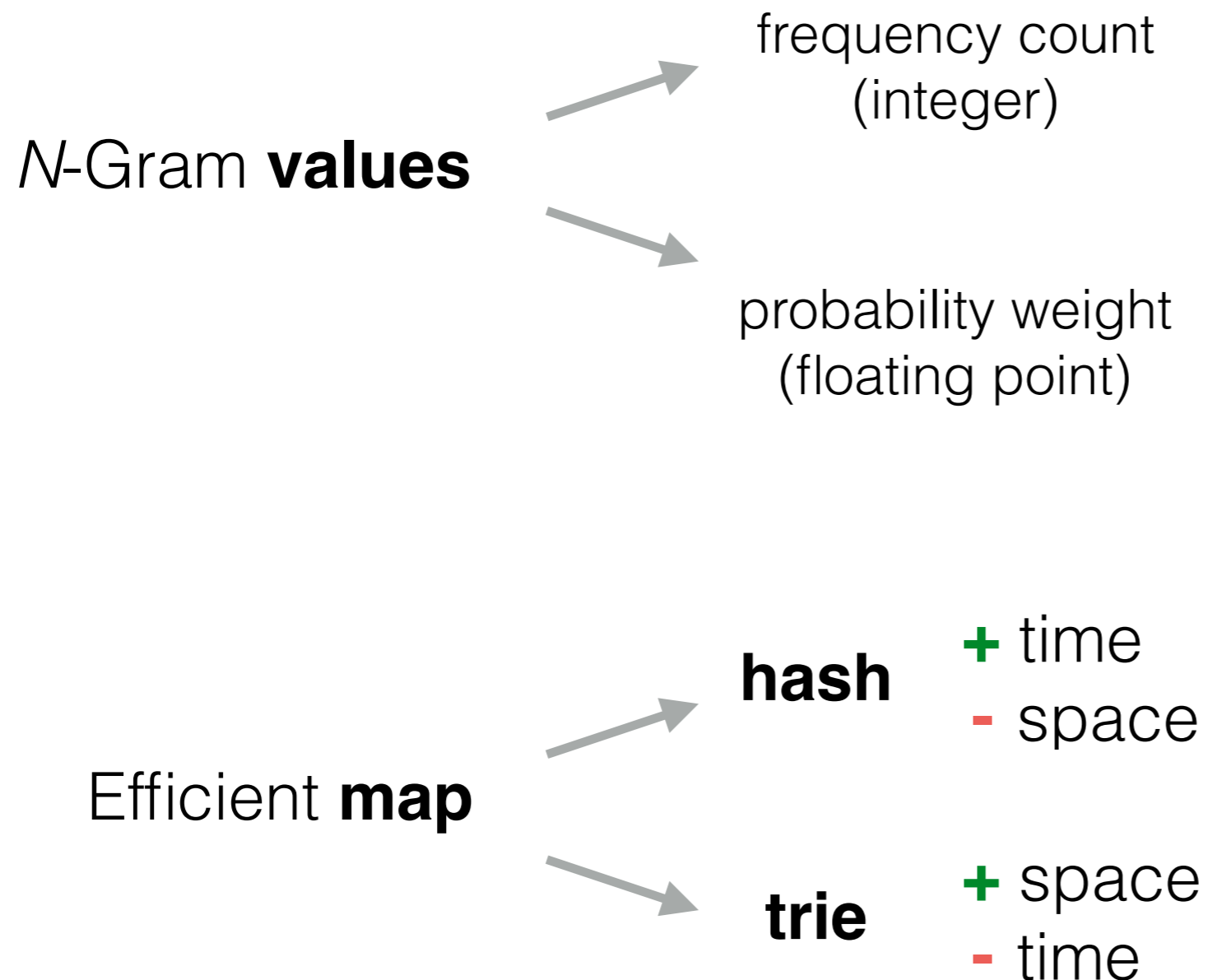
# N-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



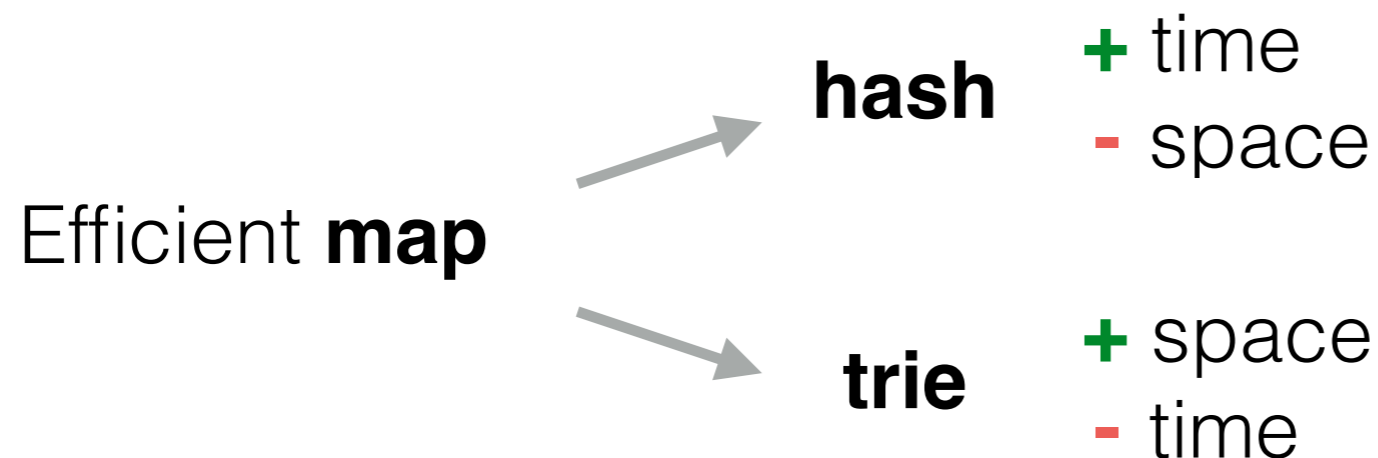
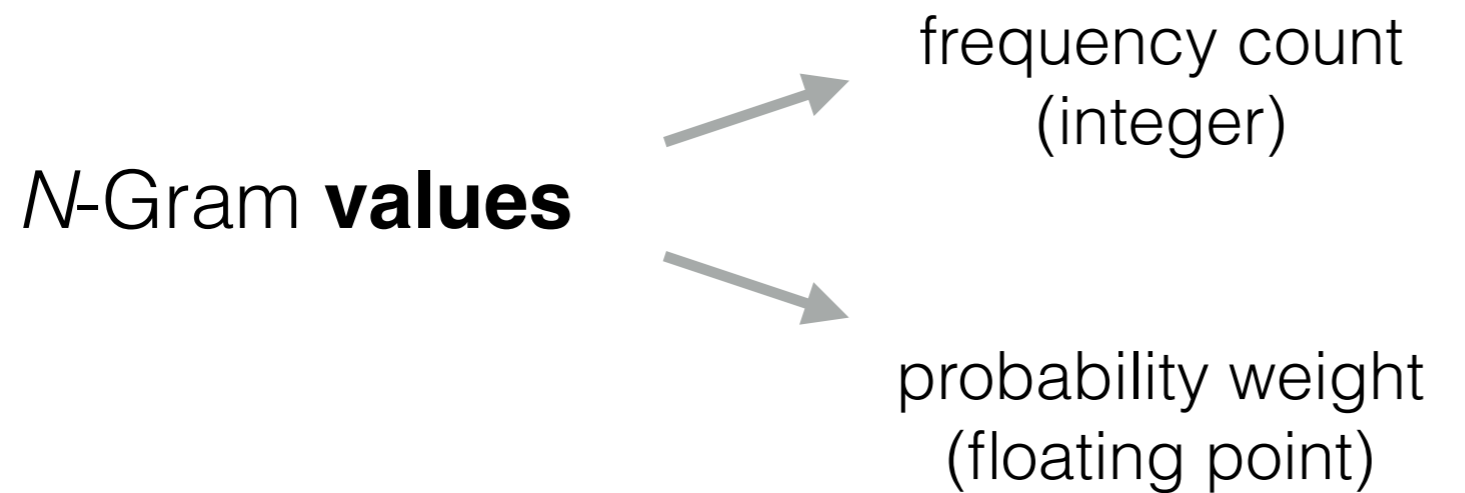
# N-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



# N-grams - Challenge

Store massive *N*-grams datasets in **compressed space** such that given a pattern, we can **return its value efficiently**.



Active field of research  
Many software libraries

- **KenLM** [Heafield, WMT 2011]
- **BerkeleyLM** [Pauls and Klein, ACL 2011]
- **ExpGram** [Watanabe et al., IJCNLP 2009]
- **IRSTLM** [Federico et al., ACL 2008]
- **RandLM** [Talbot and Osborne, ACL 2007]
- **SRILM** [Stolcke, INTERSPEECH 2002]

# Trie Indexing



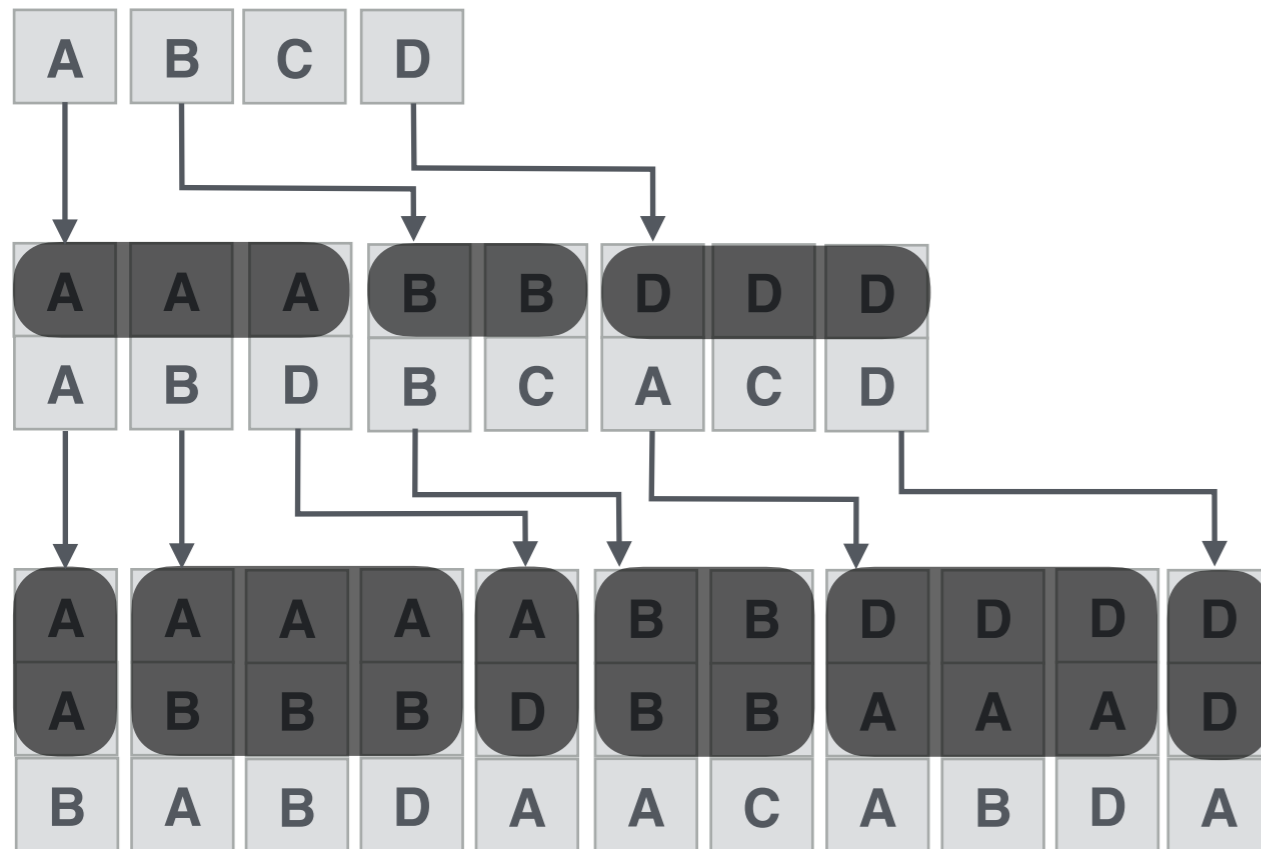
# Trie Indexing

A	B	C	D
---	---	---	---

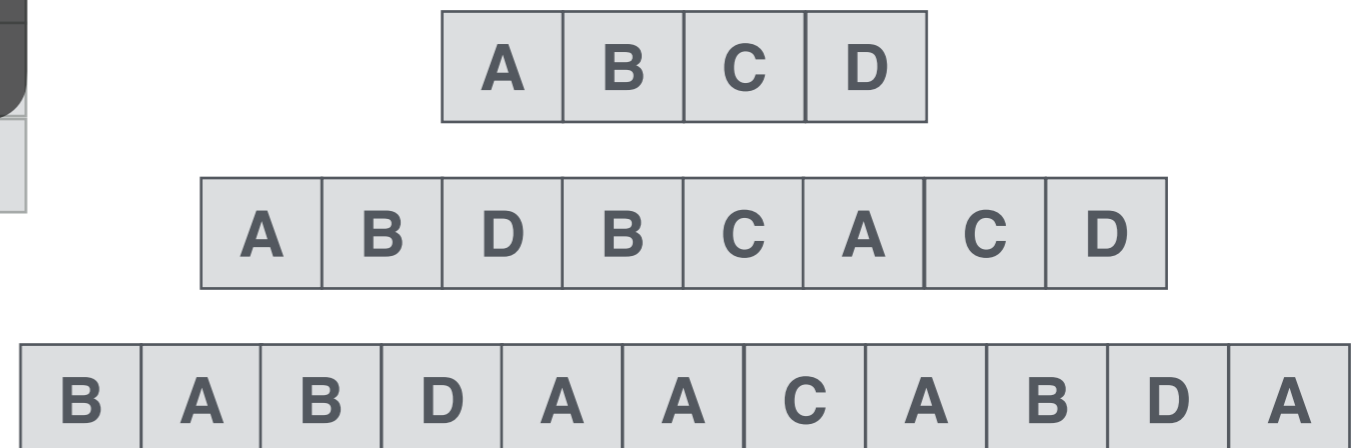
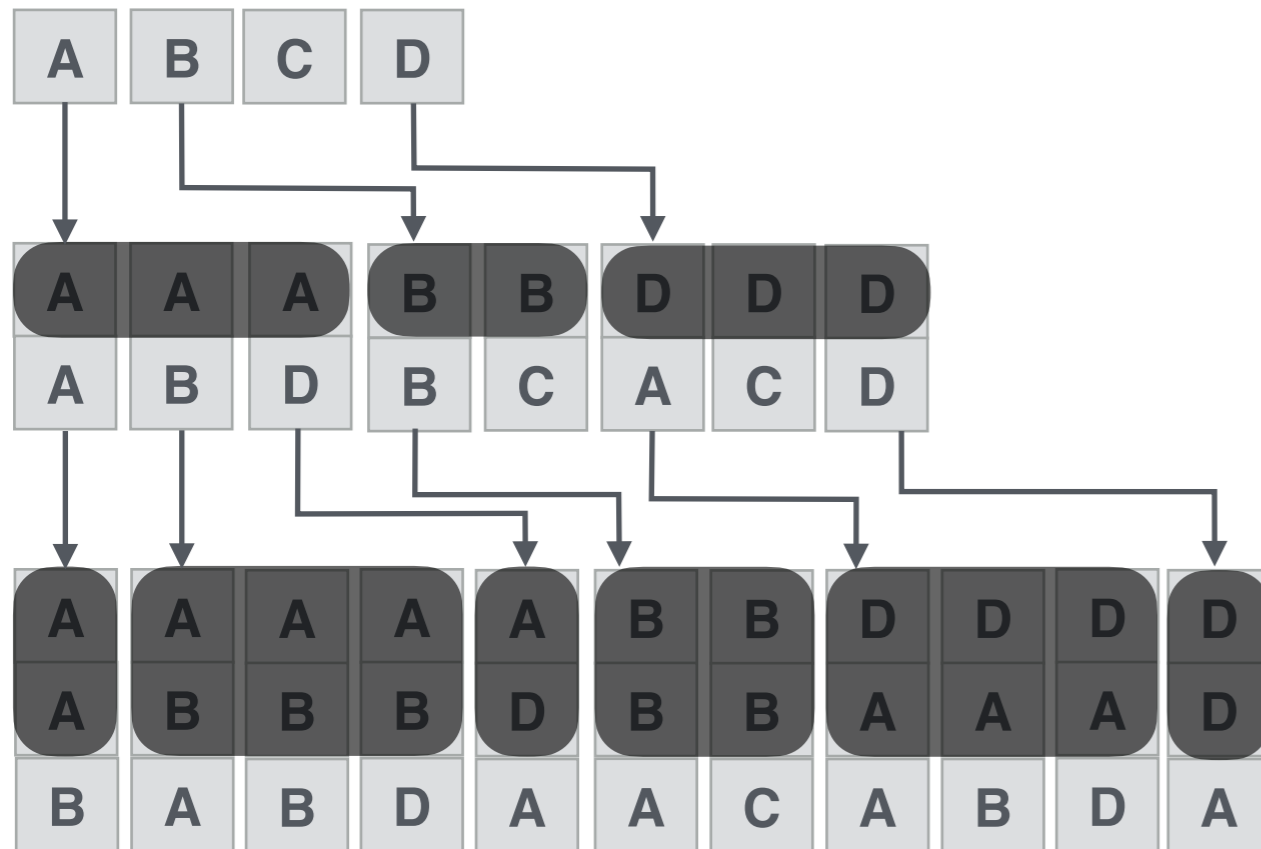
A	A	A	B	B	D	D	D
A	B	D	B	C	A	C	D

A	A	A	A	A	B	B	D	D	D	D
A	B	B	B	D	B	B	A	A	A	D
B	A	B	D	A	A	C	A	B	D	A

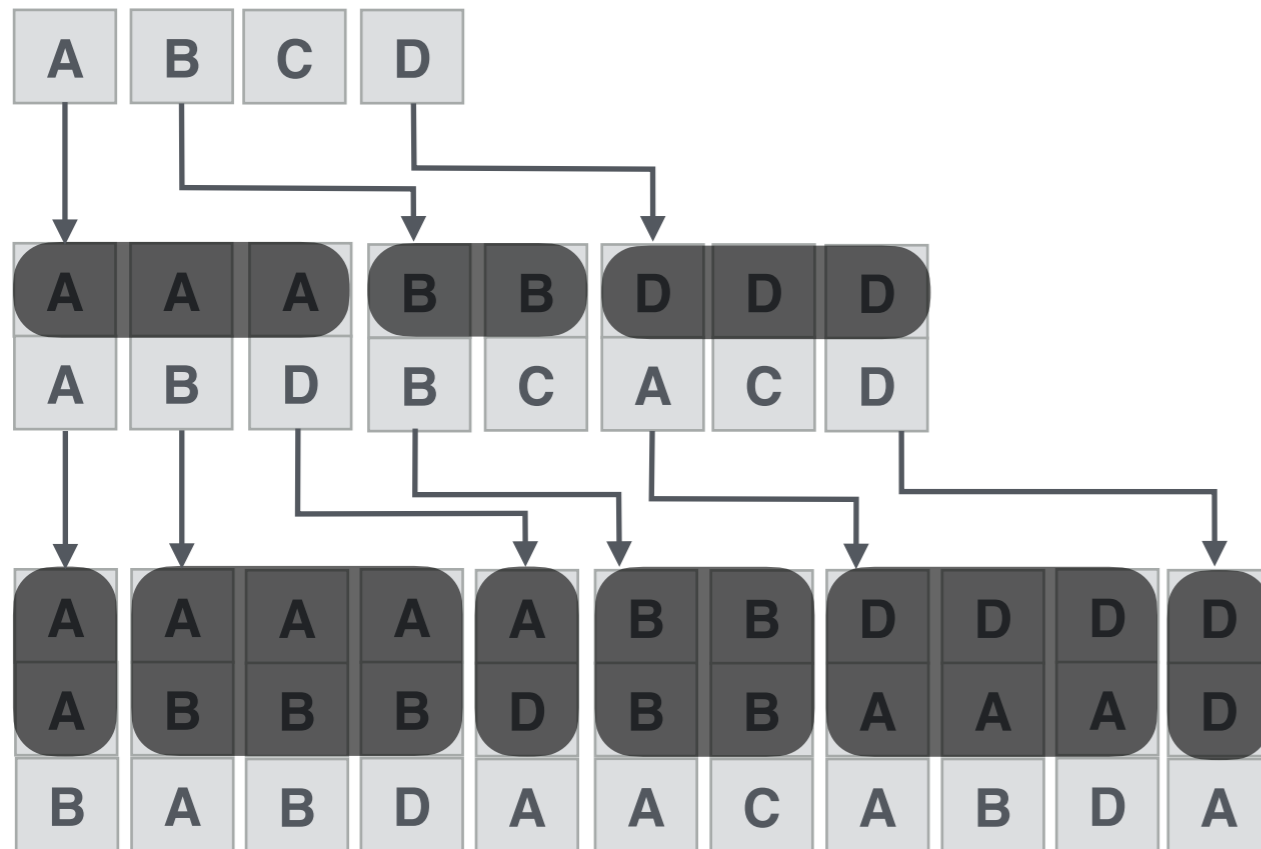
# Trie Indexing



# Trie Indexing



# Trie Indexing



**hash**  
vocabulary

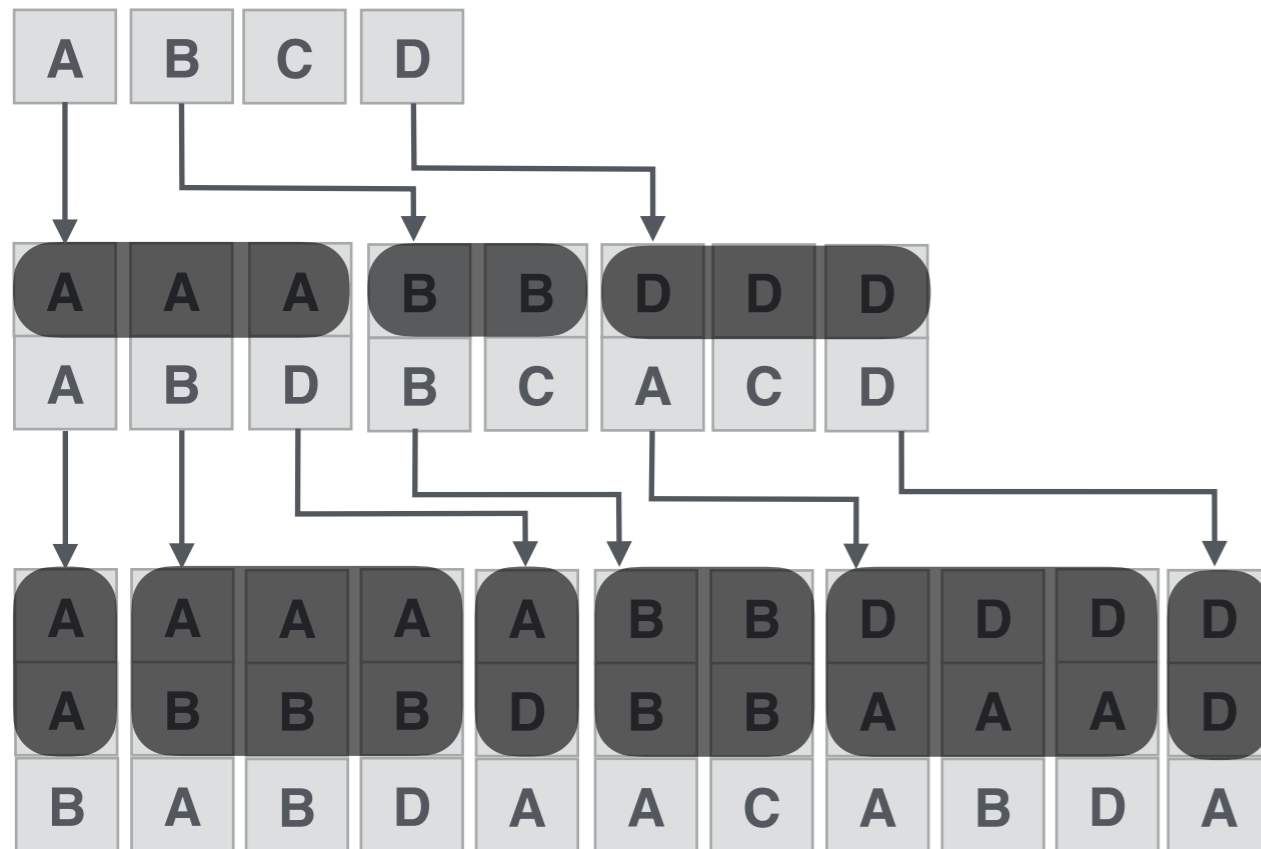
A	→	0
B	→	1
C	→	2
D	→	3

A B C D

A B D B C A C D

B A B D A A C A B D A

# Trie Indexing



**hash**  
vocabulary

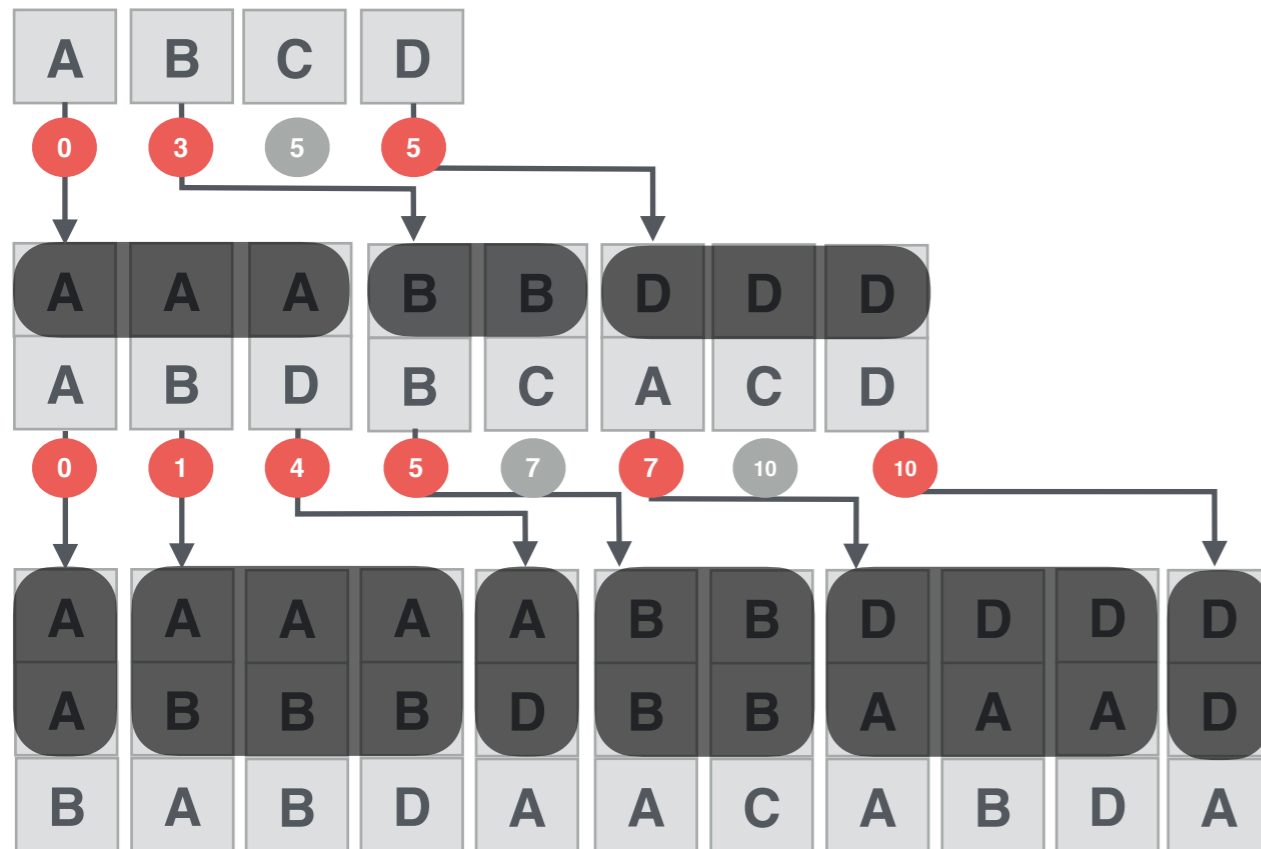
A	→	0
B	→	1
C	→	2
D	→	3

0	1	2	3
---	---	---	---

0	1	3	1	2	0	2	3
---	---	---	---	---	---	---	---

1	0	1	3	0	0	2	0	1	3	0
---	---	---	---	---	---	---	---	---	---	---

# Trie Indexing



hash  
vocabulary

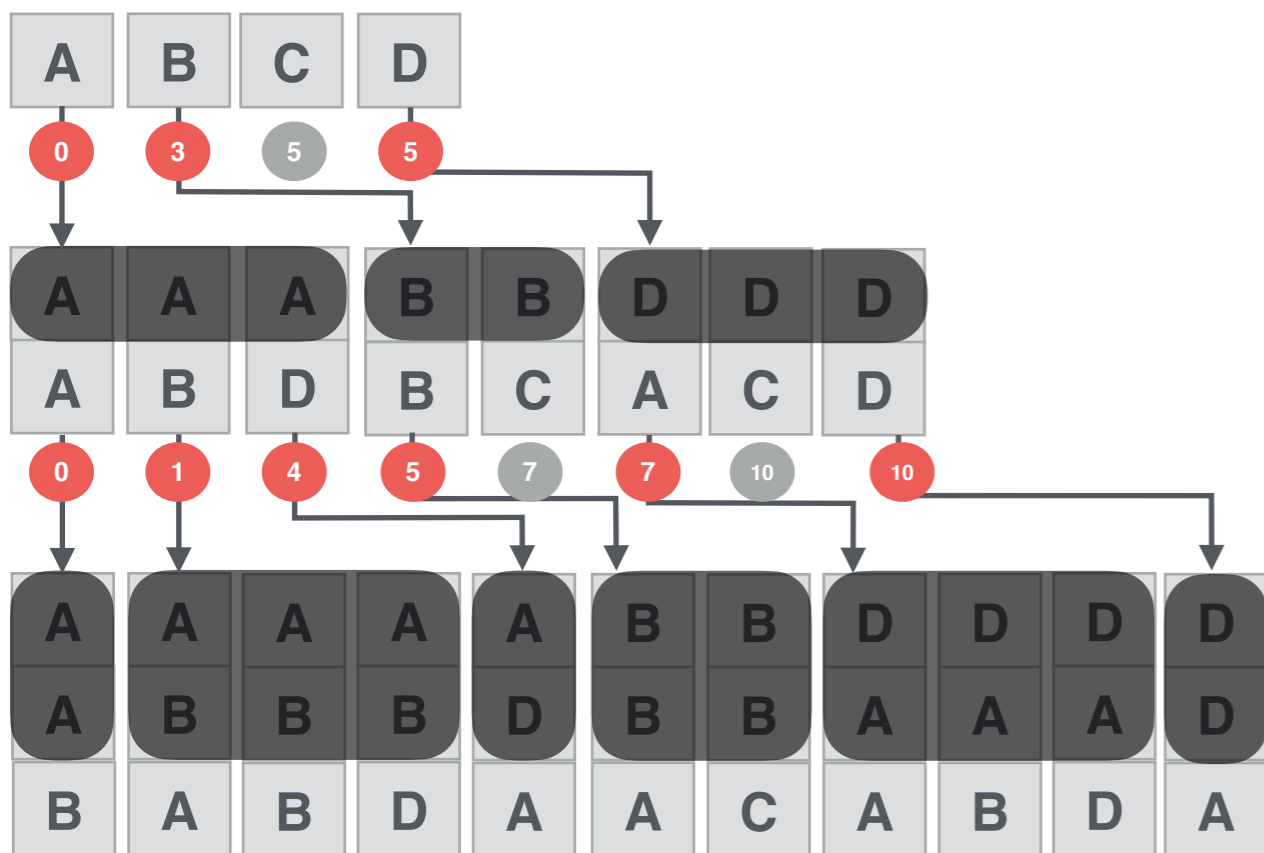
A	→	0
B	→	1
C	→	2
D	→	3

0	1	2	3
---	---	---	---

0	1	3	1	2	0	2	3
---	---	---	---	---	---	---	---

1	0	1	3	0	0	2	0	1	3	0
---	---	---	---	---	---	---	---	---	---	---

# Trie Indexing



**hash**  
vocabulary

A	→	0
B	→	1
C	→	2
D	→	3



# Trie Indexing

**hash**  
vocabulary

A	→	0
B	→	1
C	→	2
D	→	3





# Trie Indexing

We need an encoder for integer sequences, supporting fast **random Access**.

**hash**  
vocabulary

A	→	0
B	→	1
C	→	2
D	→	3



# Trie Indexing

We need an encoder for integer sequences, supporting fast **random Access**.

Take *range-wise* prefix sums on gram-ID sequences.

<b>hash</b> vocabulary	A	→	0
	B	→	1
	C	→	2
	D	→	3

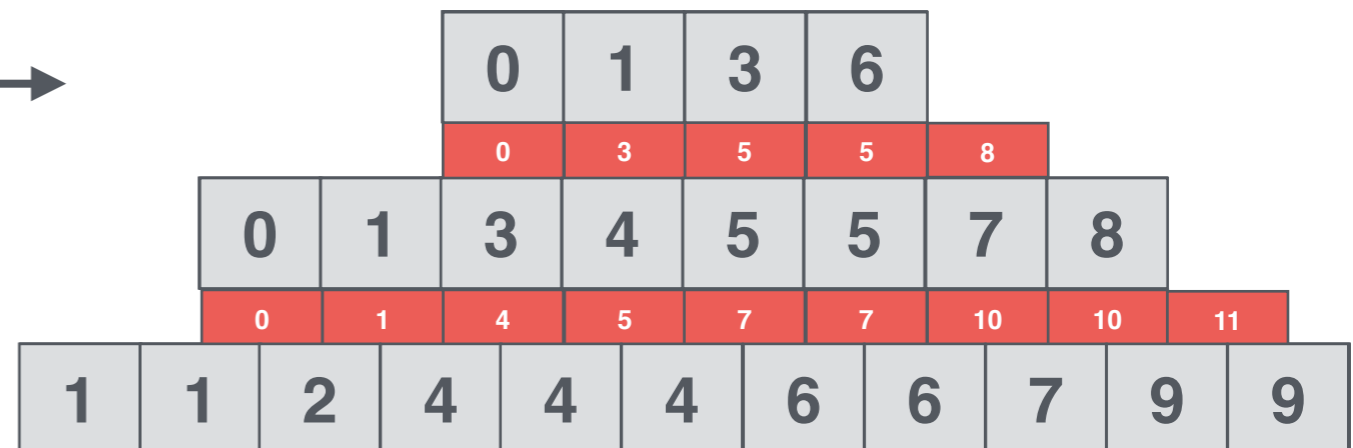


# Trie Indexing

We need an encoder for integer sequences, supporting fast **random Access**.

Take *range-wise* prefix sums on gram-ID sequences.

<b>hash</b> vocabulary	A	→	0
	B	→	1
	C	→	2
	D	→	3



# Trie Indexing

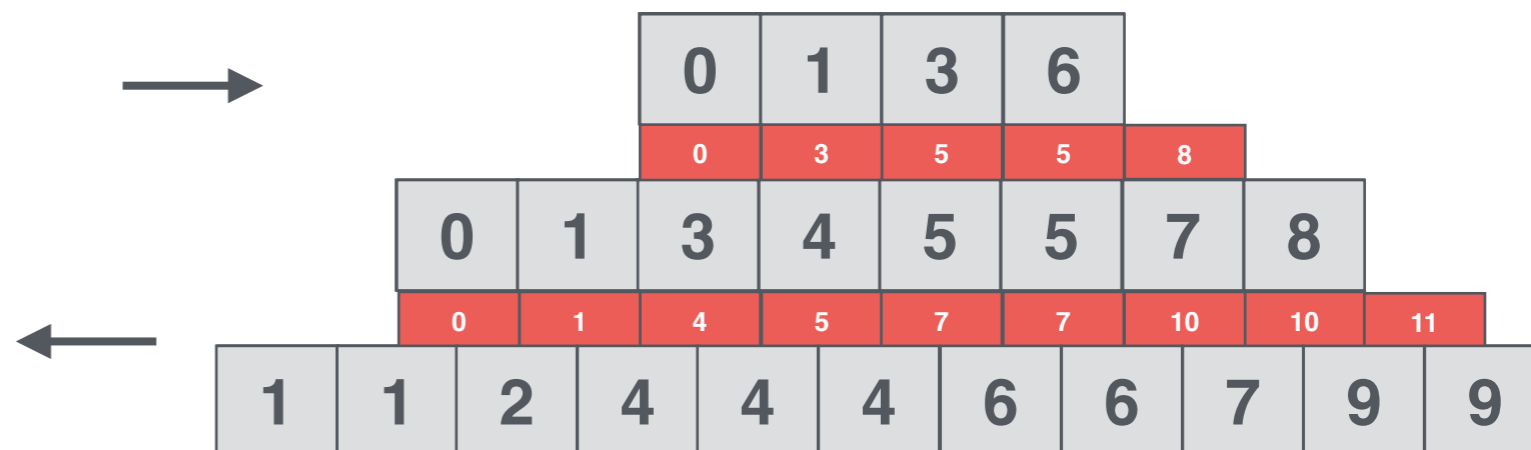
We need an encoder for integer sequences, supporting fast **random Access**.

<b>hash</b> vocabulary	A	→	0
	B	→	1
	C	→	2
	D	→	3

Take *range-wise* prefix sums on gram-ID sequences.

## Elias-Fano Tries

One **Successor** query per level  
Constant-time random **Access**



# Trie Indexing

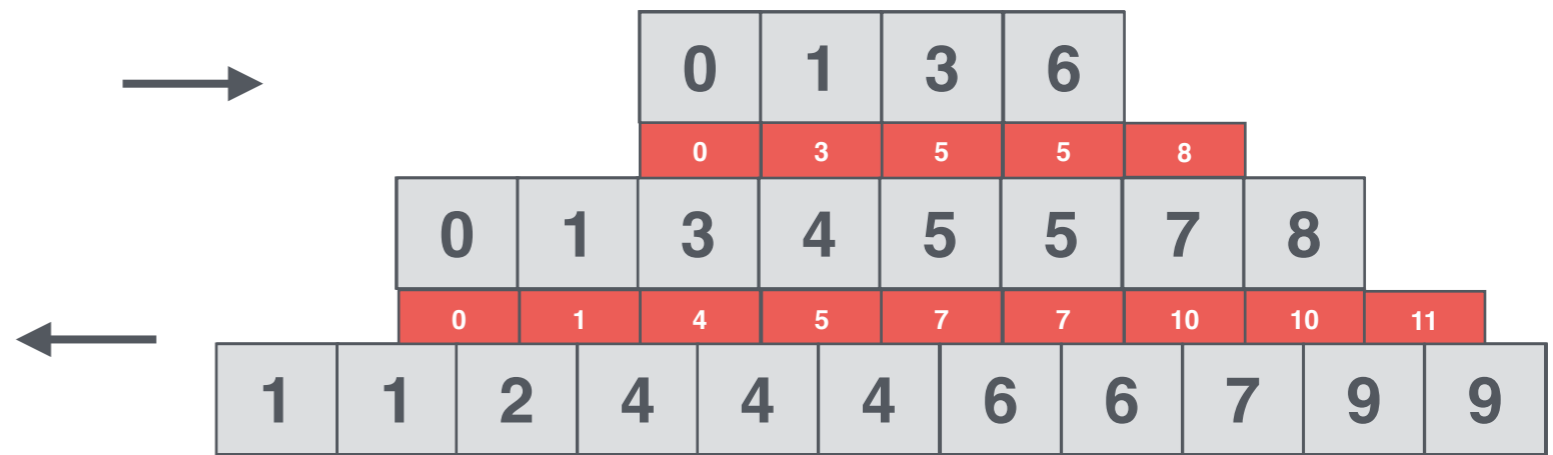
We need an encoder for integer sequences, supporting fast **random Access**.

<b>hash</b> vocabulary	A	→	0
	B	→	1
	C	→	2
	D	→	3

Take *range-wise* prefix sums on gram-ID sequences.

## Elias-Fano Tries

One **Successor** query per level  
Constant-time random **Access**



Remember:  
Elias-Fano takes  
 **$\log(u/n) + 2$**  bits  
per integer

# Context-based ID Remapping

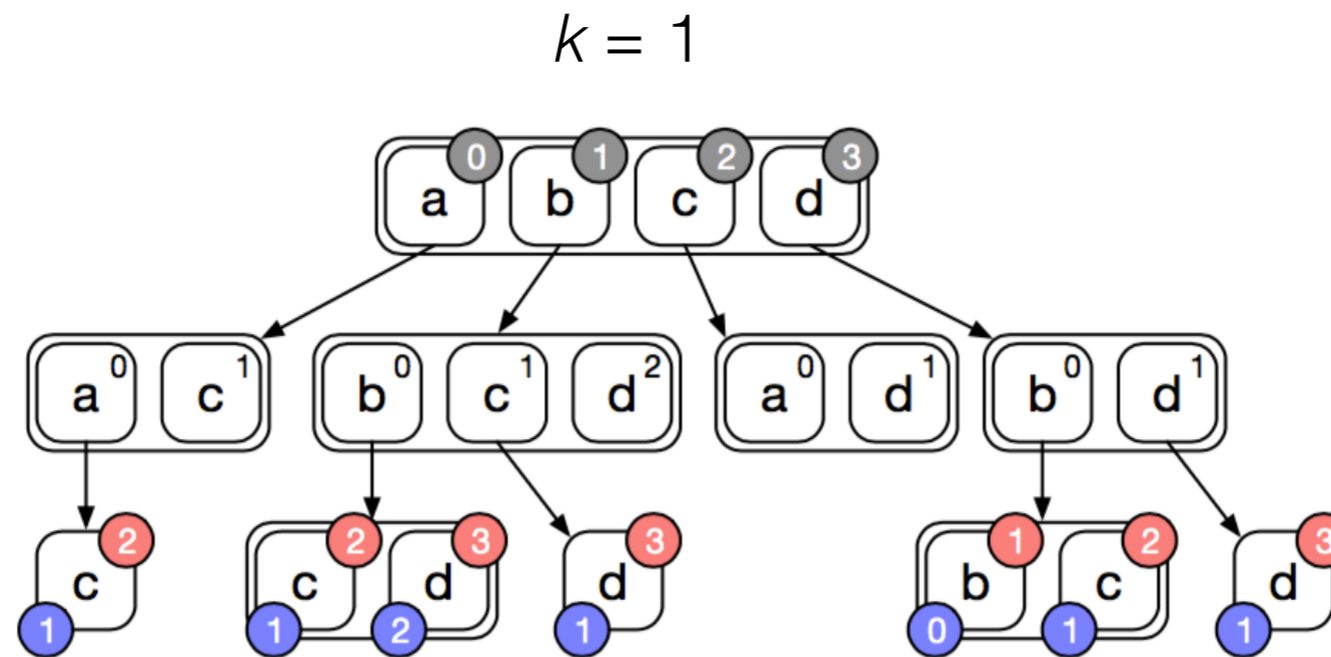
Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

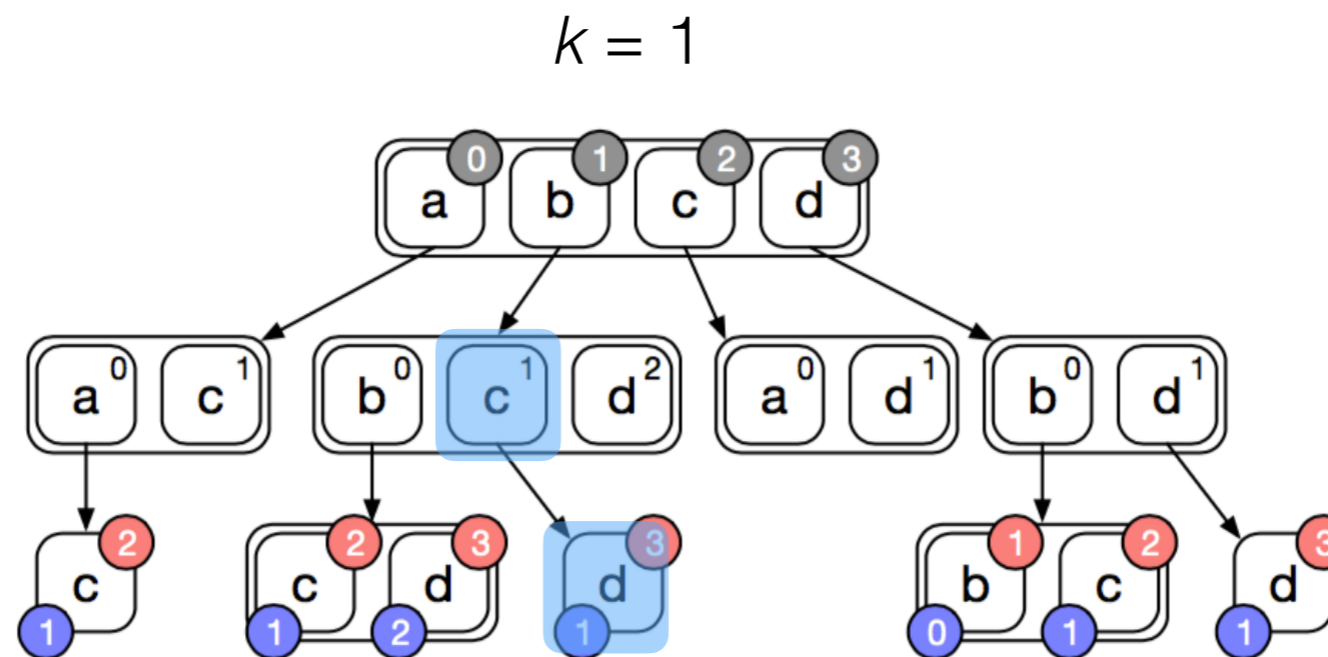
High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).



# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

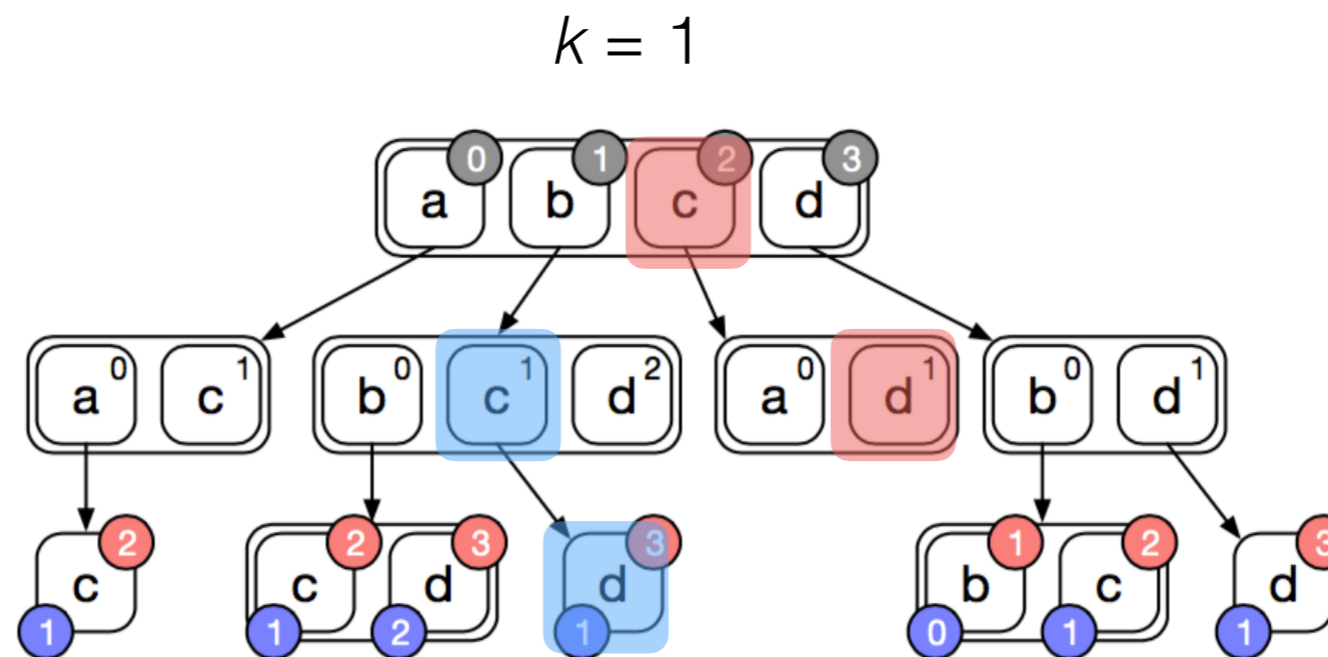




# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

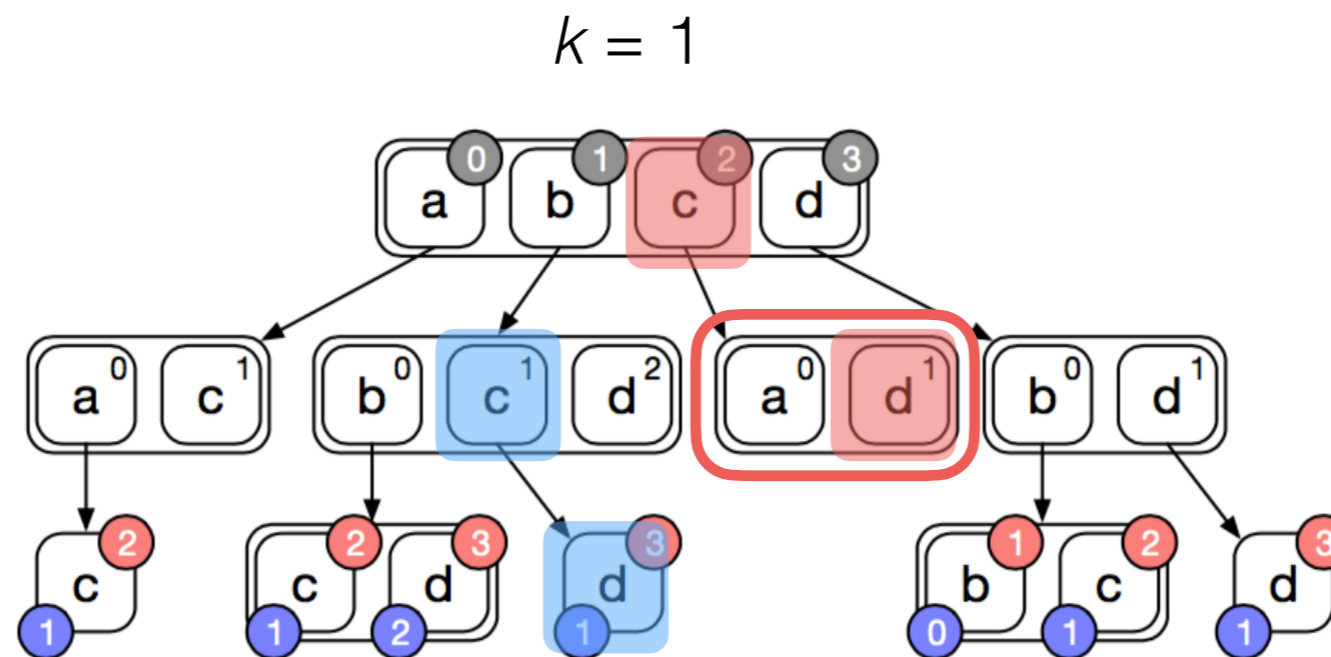
High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).



# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).



# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

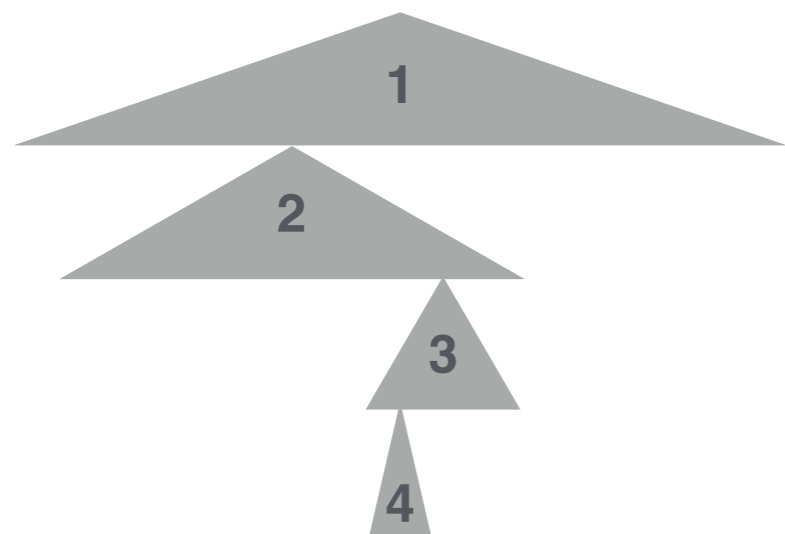
- **Millions** of unigrams.
- Height 5: **longer** contexts.
- The number of siblings has a **funnel**-shaped distribution.

# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

- **Millions** of unigrams.
- Height 5: **longer** contexts.
- The number of siblings has a **funnel-shaped** distribution.

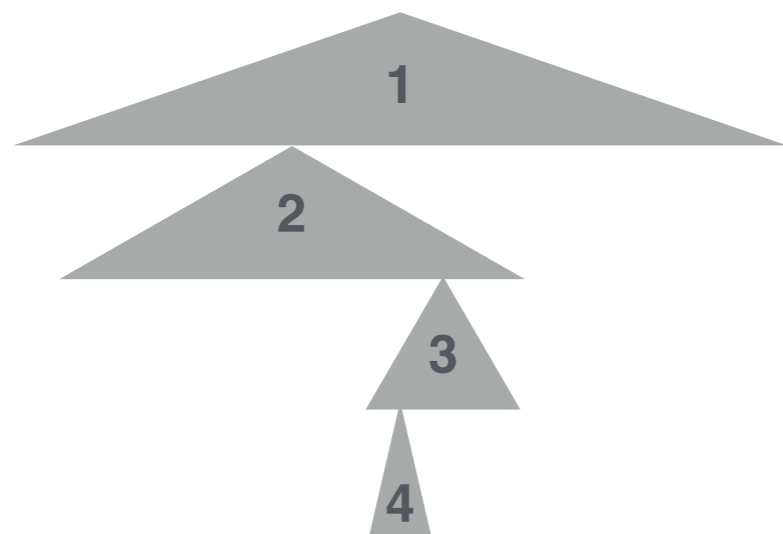


# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

- **Millions** of unigrams.
- Height 5: **longer** contexts.
- The number of siblings has a **funnel-shaped** distribution.



u/n by varying context-length  $k$

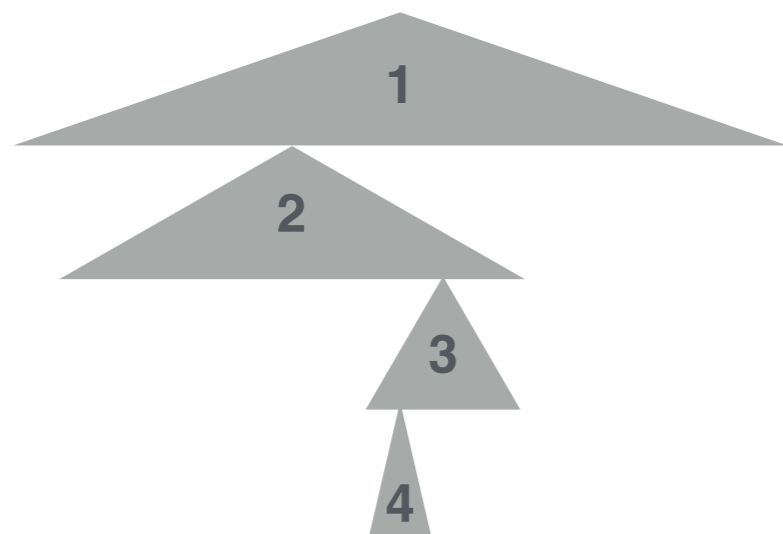
	$k$	3-grams	4-grams	5-grams
Europarl	0	2404	2782	2920
	1	213 ( $\times 11.28$ )	480 ( $\times 5.79$ )	646 ( $\times 4.52$ )
	2	2404	48 ( $\times 57.95$ )	101 ( $\times 28.91$ )
YahooV2	0	7350	7197	7417
	1	753 ( $\times 9.76$ )	1461 ( $\times 4.93$ )	1963 ( $\times 3.78$ )
	2	7350	104 ( $\times 69.20$ )	249 ( $\times 29.79$ )
GoogleV2	0	4050	6631	6793
	1	1025 ( $\times 3.95$ )	2192 ( $\times 3.03$ )	2772 ( $\times 2.45$ )
	2	4050	221 ( $\times 30.00$ )	503 ( $\times 13.50$ )

# Context-based ID Remapping

Observation: the number of words following a given context is **small**.

High-level idea: map a word ID to the **position** it takes within its *sibling* IDs (the IDs following a context of fixed length  $k$ ).

- **Millions** of unigrams.
- Height 5: **longer** contexts.
- The number of siblings has a **funnel-shaped** distribution.



u/n by varying context-length  $k$

	$k$	3-grams	4-grams	5-grams
Europarl	0	2404	2782	2920
	1	213 ( $\times 11.28$ )	480 ( $\times 5.79$ )	646 ( $\times 4.52$ )
	2	2404	48 ( $\times 57.95$ )	101 ( $\times 28.91$ )
YahooV2	0	7350	7197	7417
	1	753 ( $\times 9.76$ )	1461 ( $\times 4.93$ )	1963 ( $\times 3.78$ )
	2	7350	104 ( $\times 69.20$ )	249 ( $\times 29.79$ )
GoogleV2	0	4050	6631	6793
	1	1025 ( $\times 3.95$ )	2192 ( $\times 3.03$ )	2772 ( $\times 2.45$ )
	2	4050	221 ( $\times 30.00$ )	503 ( $\times 13.50$ )

# Experimental Analysis - EF/PEF (R)Trie

$N$	Europarl	YahooV2	GoogleV2
	$n$	$n$	$n$
1	304 579	3 475 482	24 357 349
2	5 192 260	53 844 927	665 752 080
3	18 908 249	187 639 522	7 384 478 110
4	33 862 651	287 562 409	1 642 783 634
5	43 160 518	295 701 337	1 413 870 914
Total	101 428 257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine  
Intel Xeon E5-2630 v3, 2.4 GHz  
193 GB of RAM, Linux 64 bits

**C++** implementation  
gcc 5.4.1 with the highest  
optimization setting



# Experimental Analysis - EF/PEF (R)Trie

$N$	Europarl	YahooV2	GoogleV2
	$n$	$n$	$n$
1	304 579	3 475 482	24 357 349
2	5 192 260	53 844 927	665 752 080
3	18 908 249	187 639 522	7 384 478 110
4	33 862 651	287 562 409	1 642 783 634
5	43 160 518	295 701 337	1 413 870 914
Total	101 428 257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine  
Intel Xeon E5-2630 v3, 2.4 GHz  
193 GB of RAM, Linux 64 bits

**C++** implementation  
gcc 5.4.1 with the highest  
optimization setting

		Europarl		YahooV2		GoogleV2		
		bpg	$\mu s \times \text{query}$	bpg	$\mu s \times \text{query}$	bpg	$\mu s \times \text{query}$	
CONTEXT-BASED ID REMAPPING	$k = 1$	EF	1.97	1.28	2.17	1.60	2.13	2.09
		PEF	1.87 (-4.99%)	1.35 (+5.93%)	1.91 (-12.03%)	1.73 (+8.00%)	1.52 (-28.60%)	1.91 (-8.79%)
	$k = 2$	EF	1.67 (-15.30%)	1.58 (+23.86%)	1.89 (-12.92%)	2.05 (+28.07%)	1.91 (-10.24%)	3.03 (+44.61%)
		PEF	1.53 (-22.36%)	1.61 (+25.89%)	1.63 (-24.91%)	2.16 (+35.22%)	1.31 (-38.71%)	2.30 (+9.88%)
	$k = 2$	EF	1.46 (-25.62%)	1.60 (+25.17%)	1.68 (-22.32%)	2.08 (+30.23%)	—	—
		PEF	1.28 (-34.87%)	1.64 (+28.12%)	1.38 (-36.15%)	2.15 (+34.81%)	—	—

# Experimental Analysis - EF/PEF (R)Trie

$N$	Europarl	YahooV2	GoogleV2
	$n$	$n$	$n$
1	304 579	3 475 482	24 357 349
2	5 192 260	53 844 927	665 752 080
3	18 908 249	187 639 522	7 384 478 110
4	33 862 651	287 562 409	1 642 783 634
5	43 160 518	295 701 337	1 413 870 914
Total	101 428 257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine  
Intel Xeon E5-2630 v3, 2.4 GHz  
193 GB of RAM, Linux 64 bits

**C++** implementation  
gcc 5.4.1 with the highest  
optimization setting

		Europarl		YahooV2		GoogleV2		
		bpg	$\mu s \times \text{query}$	bpg	$\mu s \times \text{query}$	bpg	$\mu s \times \text{query}$	
CONTEXT-BASED ID REMAPPING	$k = 1$	EF	1.97	1.28	2.17	1.60	2.13	2.09
		PEF	1.87 (-4.99%)	1.35 (+5.93%)	1.91 (-12.03%)	1.73 (+8.00%)	1.52 (-28.60%)	1.91 (-8.79%)
	$k = 2$	EF	1.67 (-15.30%)	1.58 (+23.86%)	1.89 (-12.92%)	2.05 (+28.07%)	1.91 (-10.24%)	3.03 (+44.61%)
		PEF	1.53 (-22.36%)	1.61 (+25.89%)	1.63 (-24.91%)	2.16 (+35.22%)	1.31 (-38.71%)	2.30 (+9.88%)
	$k = 2$	EF	1.46 (-25.62%)	1.60 (+25.17%)	1.68 (-22.32%)	2.08 (+30.23%)	—	—
		PEF	1.28 (-34.87%)	1.64 (+28.12%)	1.38 (-36.15%)	2.15 (+34.81%)	—	—

# Experimental Analysis - EF/PEF (R)Trie

$N$	Europarl	YahooV2	GoogleV2
	$n$	$n$	$n$
1	304 579	3 475 482	24 357 349
2	5 192 260	53 844 927	665 752 080
3	18 908 249	187 639 522	7 384 478 110
4	33 862 651	287 562 409	1 642 783 634
5	43 160 518	295 701 337	1 413 870 914
Total	101 428 257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine  
Intel Xeon E5-2630 v3, 2.4 GHz  
193 GB of RAM, Linux 64 bits

**C++** implementation  
gcc 5.4.1 with the highest  
optimization setting

		Europarl		YahooV2		GoogleV2	
		bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
EF		1.97	1.28	2.17	1.60	2.13	2.09
PEF		1.87 (-4.99%)	1.35 (+5.93%)	1.91 (-12.03%)	1.73 (+8.00%)	1.52 (-28.60%)	1.91 (-8.79%)
CONTEXT-BASED ID REMAPPING	$k = 1$						
	EF	1.67 (-15.30%)	1.58 (+23.86%)	1.89 (-12.92%)	2.05 (+28.07%)	1.91 (-10.24%)	3.03 (+44.61%)
	PEF	1.53 (-22.36%)	1.61 (+25.89%)	1.63 (-24.91%)	2.16 (+35.22%)	1.31 (-38.71%)	2.30 (+9.88%)
	$k = 2$						
	EF	1.46 (-25.62%)	1.60 (+25.17%)	1.68 (-22.32%)	2.08 (+30.23%)	—	—
	PEF	1.28 (-34.87%)	1.64 (+28.12%)	1.38 (-36.15%)	2.15 (+34.81%)	—	—

## Context-based ID Remapping

- reduces space by more than **36%** on average  $\longrightarrow$  **you will** notice this!

# Experimental Analysis - EF/PEF (R)Trie

$N$	Europarl	YahooV2	GoogleV2
	$n$	$n$	$n$
1	304 579	3 475 482	24 357 349
2	5 192 260	53 844 927	665 752 080
3	18 908 249	187 639 522	7 384 478 110
4	33 862 651	287 562 409	1 642 783 634
5	43 160 518	295 701 337	1 413 870 914
Total	101 428 257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine  
Intel Xeon E5-2630 v3, 2.4 GHz  
193 GB of RAM, Linux 64 bits

**C++** implementation  
gcc 5.4.1 with the highest  
optimization setting

		Europarl		YahooV2		GoogleV2	
		bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
EF		1.97	1.28	2.17	1.60	2.13	2.09
PEF		1.87 (-4.99%)	1.35 (+5.93%)	1.91 (-12.03%)	1.73 (+8.00%)	1.52 (-28.60%)	1.91 (-8.79%)
CONTEXT-BASED ID REMAPPING	$k = 1$						
	EF	1.67 (-15.30%)	1.58 (+23.86%)	1.89 (-12.92%)	2.05 (+28.07%)	1.91 (-10.24%)	3.03 (+44.61%)
	PEF	1.53 (-22.36%)	1.61 (+25.89%)	1.63 (-24.91%)	2.16 (+35.22%)	1.31 (-38.71%)	2.30 (+9.88%)
	$k = 2$						
EF	1.46 (-25.62%)	1.60 (+25.17%)	1.68 (-22.32%)	2.08 (+30.23%)	—	—	
PEF	1.28 (-34.87%)	1.64 (+28.12%)	1.38 (-36.15%)	2.15 (+34.81%)	—	—	

## Context-based ID Remapping

- reduces space by more than **36%** on average  $\longrightarrow$  **you will** notice this!

# Experimental Analysis - EF/PEF (R)Trie

$N$	Europarl	YahooV2	GoogleV2
	$n$	$n$	$n$
1	304 579	3 475 482	24 357 349
2	5 192 260	53 844 927	665 752 080
3	18 908 249	187 639 522	7 384 478 110
4	33 862 651	287 562 409	1 642 783 634
5	43 160 518	295 701 337	1 413 870 914
Total	101 428 257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine  
Intel Xeon E5-2630 v3, 2.4 GHz  
193 GB of RAM, Linux 64 bits

**C++** implementation  
gcc 5.4.1 with the highest  
optimization setting

		Europarl		YahooV2		GoogleV2	
		bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
EF		1.97	1.28	2.17	1.60	2.13	2.09
PEF		1.87 (-4.99%)	1.35 (+5.93%)	1.91 (-12.03%)	1.73 (+8.00%)	1.52 (-28.60%)	1.91 (-8.79%)
CONTEXT-BASED ID REMAPPING	$k = 1$						
	EF	1.67 (-15.30%)	1.58 (+23.86%)	1.89 (-12.92%)	2.05 (+28.07%)	1.91 (-10.24%)	3.03 (+44.61%)
	PEF	1.53 (-22.36%)	1.61 (+25.89%)	1.63 (-24.91%)	2.16 (+35.22%)	1.31 (-38.71%)	2.30 (+9.88%)
	$k = 2$						
EF	1.46 (-25.62%)	1.60 (+25.17%)	1.68 (-22.32%)	2.08 (+30.23%)	—	—	
PEF	1.28 (-34.87%)	1.64 (+28.12%)	1.38 (-36.15%)	2.15 (+34.81%)	—	—	

## Context-based ID Remapping

- reduces space by more than **36%** on average → **you will** notice this!
- brings approximately **30%** more time → **will you** notice this?

# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$
PEF-Trie	1.87	1.35	1.91	1.73	1.52	1.91
PEF-RTrie	1.28	1.64	1.38	2.15	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	2.83 (+108.88%) (+72.70%)	1.69 (-11.41%) (+22.04%)	3.48 (+101.84%) (+61.70%)	1.45 (-4.87%) (+10.83%)	4.13 (+116.57%) (+79.76%)
BerkeleyLM H.3	6.70 (+258.81%) (+423.40%)	0.97 (-28.46%) (-40.85%)	7.82 (+310.38%) (+465.36%)	1.13 (-34.35%) (-47.41%)	9.24 (+507.79%) (+608.07%)	2.18 (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+326.03%) (+521.45%)	0.97 (-28.49%) (-40.88%)	9.37 (+391.32%) (+576.87%)	0.96 (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	2.80 (+106.61%) (+70.82%)	2.24 (+17.36%) (+61.68%)	9.23 (+435.33%) (+328.87%)	—	—
KenLM T.	2.99 (+60.11%) (+133.56%)	1.28 (-5.47%) (-21.84%)	3.44 (+80.39%) (+148.52%)	1.94 (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+93.09%) (+181.66%)	2.06 (+52.00%) (+25.67%)	3.81 (+99.60%) (+174.98%)	3.24 (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	4.39 (+224.20%) (+168.04%)	2.02 (+6.18%) (+46.29%)	5.08 (+194.35%) (+135.82%)	2.60 (+70.73%) (+98.90%)	9.25 (+384.54%) (+302.19%)

# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$
PEF-Trie	1.87	1.35	1.91	1.73	1.52	1.91
PEF-RTrie	1.28	1.64	1.38	2.15	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	2.83 (+108.88%) (+72.70%)	1.69 (-11.41%) (+22.04%)	3.48 (+101.84%) (+61.70%)	1.45 (-4.87%) (+10.83%)	4.13 (+116.57%) (+79.76%)
BerkeleyLM H.3	6.70 (+258.81%) (+423.40%)	0.97 (-28.46%) (-40.85%)	7.82 (+310.38%) (+465.36%)	1.13 (-34.35%) (-47.41%)	9.24 (+507.79%) (+608.07%)	2.18 (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+326.03%) (+521.45%)	0.97 (-28.49%) (-40.88%)	9.37 (+391.32%) (+576.87%)	0.96 (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	2.80 (+106.61%) (+70.82%)	2.24 (+17.36%) (+61.68%)	9.23 (+435.33%) (+328.87%)	—	—
KenLM T.	2.99 (+60.11%) (+133.56%)	1.28 (-5.47%) (-21.84%)	3.44 (+80.39%) (+148.52%)	1.94 (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+93.09%) (+181.66%)	2.06 (+52.00%) (+25.67%)	3.81 (+99.60%) (+174.98%)	3.24 (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	4.39 (+224.20%) (+168.04%)	2.02 (+6.18%) (+46.29%)	5.08 (+194.35%) (+135.82%)	2.60 (+70.73%) (+98.90%)	9.25 (+384.54%) (+302.19%)

# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$
PEF-Trie	1.87	1.35	1.91	1.73	1.52	1.91
PEF-RTrie	1.28	1.64	1.38	2.15	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	2.83 (+108.88%) (+72.70%)	1.69 (-11.41%) (+22.04%)	3.48 (+101.84%) (+61.70%)	1.45 (-4.87%) (+10.83%)	4.13 (+116.57%) (+79.76%)
BerkeleyLM H.3	6.70 (+258.81%) (+423.40%)	0.97 (-28.46%) (-40.85%)	7.82 (+310.38%) (+465.36%)	1.13 (-34.35%) (-47.41%)	9.24 (+507.79%) (+608.07%)	2.18 (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+326.03%) (+521.45%)	0.97 (-28.49%) (-40.88%)	9.37 (+391.32%) (+576.87%)	0.96 (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	2.80 (+106.61%) (+70.82%)	2.24 (+17.36%) (+61.68%)	9.23 (+435.33%) (+328.87%)	—	—
KenLM T.	2.99 (+133.56%)	1.28 (-5.47%) (-21.84%)	3.44 (+148.52%)	1.94 (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+93.09%) (+181.66%)	2.06 (+52.00%) (+25.67%)	3.81 (+99.60%) (+174.98%)	3.24 (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	4.39 (+224.20%) (+168.04%)	2.02 (+6.18%) (+46.29%)	5.08 (+194.35%) (+135.82%)	2.60 (+70.73%) (+98.90%)	9.25 (+384.54%) (+302.19%)



# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie	1.87	<b>1.35</b>	1.91	<b>1.73</b>	1.52	1.91
PEF-RTrie	1.28	<b>1.64</b>	1.38	<b>2.15</b>	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	<b>2.83</b> (+108.88%) (+72.70%)	1.69 (-11.41%) (+22.04%)	<b>3.48</b> (+101.84%) (+61.70%)	1.45 (-4.87%) (+10.83%)	<b>4.13</b> (+116.57%) (+79.76%)
BerkeleyLM H.3	6.70 (+258.81%) (+423.40%)	<b>0.97</b> (-28.46%) (-40.85%)	7.82 (+310.38%) (+465.36%)	<b>1.13</b> (-34.35%) (-47.41%)	9.24 (+507.79%) (+608.07%)	<b>2.18</b> (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+326.03%) (+521.45%)	<b>0.97</b> (-28.49%) (-40.88%)	9.37 (+391.32%) (+576.87%)	<b>0.96</b> (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	<b>2.80</b> (+106.61%) (+70.82%)	2.24 (+17.36%) (+61.68%)	<b>9.23</b> (+435.33%) (+328.87%)	—	—
KenLM T.	2.99 (+133.56%)	<b>1.28</b> (-5.47%) (-21.84%)	3.44 (+148.52%)	<b>1.94</b> (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+93.09%) (+181.66%)	<b>2.06</b> (+52.00%) (+25.67%)	3.81 (+99.60%) (+174.98%)	<b>3.24</b> (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	<b>4.39</b> (+224.20%) (+168.04%)	2.02 (+6.18%) (+46.29%)	<b>5.08</b> (+194.35%) (+135.82%)	2.60 (+70.73%) (+98.90%)	<b>9.25</b> (+384.54%) (+302.19%)

# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie	1.87	1.35	1.91	1.73	1.52	1.91
PEF-RTrie	1.28	1.64	1.38	2.15	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	2.83 (+105.78%) (+72.70%) <b>2X</b>	1.69 (-11.41%) (+22.04%)	3.48 (+100.54%) (+61.70%) <b>2X</b>	1.45 (-4.87%) (+10.83%)	4.13 (+170.77%) (+79.76%) <b>2X</b>
BerkeleyLM H.3	6.70 (+258.81%) (+105.10%) <b>2.5÷</b>	0.97 (-28.46%) (-40.85%)	7.82 (+310.38%) (+105.35%) <b>3.1÷</b>	1.13 (-34.35%) (-47.41%)	9.24 (+597.79%) (+608.07%) <b>5.5X</b>	2.18 (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+292.32%) (+521.45%) <b>5.2X</b>	0.97 (-28.49%) (-40.88%)	9.37 (+345.33%) (+576.87%) <b>5.8X</b>	0.96 (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	2.80 (+105.71%) (+70.82%) <b>2X</b>	2.24 (+17.36%) (+61.68%)	9.23 (+335.97%) (+328.87%) <b>3.5X</b>	—	—
KenLM T.	2.99 (+108.11%) (+133.56%) <b>2.3X</b>	1.28 (-5.47%) (-21.84%)	3.44 (+120.33%) (+148.52%) <b>2.5X</b>	1.94 (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+193.09%) (+161.65%) <b>2.8X</b>	2.06 (+52.00%) (+25.67%)	3.81 (+199.60%) (+114.95%) <b>2.7X</b>	3.24 (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	4.39 (+222.77%) (+168.04%) <b>2.5X</b>	2.02 (+6.18%) (+46.29%)	5.08 (+193.39%) (+135.82%) <b>2.5X</b>	2.60 (+70.73%) (+98.90%)	9.25 (+302.54%) (+302.19%) <b>3X</b>

# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie	1.87	1.35	1.91	1.73	1.52	1.91
PEF-RTrie	1.28	1.64	1.38	2.15	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	2.83 (+105.78%) (+72.70%) <b>2X</b>	1.69 (-11.41%) (+22.04%)	3.48 (+104.44%) (+61.70%) <b>2X</b>	1.45 (-4.87%) (+10.83%)	4.13 (+170.77%) (+79.76%) <b>2X</b>
BerkeleyLM H.3	6.70 (+258.81%) (+105.10%) <b>2.5X</b>	0.97 (-28.46%) (-40.85%)	7.82 (+310.38%) (+105.35%) <b>3.1X</b>	1.13 (-34.35%) (-47.41%)	9.24 (+597.76%) (+608.07%) <b>5.5X</b>	2.18 (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+292.03%) (+521.45%) <b>5.2X</b>	0.97 (-28.49%) (-40.88%)	9.37 (+349.33%) (+576.87%) <b>5.8X</b>	0.96 (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	2.80 (+105.71%) (+70.82%) <b>2X</b>	2.24 (+17.36%) (+61.68%)	9.23 (+355.97%) (+328.87%) <b>3.5X</b>	—	—
KenLM T.	2.99 (+108.11%) (+133.56%) <b>2.3X</b>	1.28 (-5.47%) (-21.84%)	3.44 (+180.33%) (+148.52%) <b>2.5X</b>	1.94 (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+93.09%) (+161.65%) <b>2.8X</b>	2.06 (+52.00%) (+25.67%)	3.81 (+99.60%) (+114.95%) <b>2.7X</b>	3.24 (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	4.39 (+225.77%) (+168.04%) <b>2.5X</b>	2.02 (+6.18%) (+46.29%)	5.08 (+193.39%) (+135.82%) <b>2.5X</b>	2.60 (+70.73%) (+98.90%)	9.25 (+385.54%) (+302.19%) <b>3X</b>

# Experimental Analysis - Overall comparison

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$	bpg	$\mu\text{s} \times \text{query}$
PEF-Trie	1.87	1.35	1.91	1.73	1.52	1.91
PEF-RTrie	1.28	1.64	1.38	2.15	1.31	2.30
BerkeleyLM C.	1.70 (-8.89%) (+32.90%)	2.83 (+108.8%) (+72.70%) <b>2X</b>	1.69 (-11.41%) (+22.04%)	3.48 (+104.4%) (+61.70%) <b>2X</b>	1.45 (-4.87%) (+10.83%)	4.13 (+113.7%) (+79.76%) <b>2X</b>
BerkeleyLM H.3	6.70 (+258.81%) (+105.00%) <b>2.5X</b>	0.97 (-28.46%) (-40.85%)	7.82 (+310.38%) (+105.30%) <b>3.1X</b>	1.13 (-34.35%) (-47.41%)	9.24 (+597.70%) (+608.07%) <b>5.5X</b>	2.18 (+13.95%) (-5.42%)
BerkeleyLM H.50	7.96 (+292.00%) (+521.45%) <b>5.2X</b>	0.97 (-28.49%) (-40.88%)	9.37 (+395.30%) (+576.87%) <b>5.8X</b>	0.96 (-44.27%) (-55.35%)	—	—
Expgram	2.06 (+10.18%) (+60.73%)	2.80 (+105.1%) (+70.82%) <b>2X</b>	2.24 (+17.36%) (+61.68%)	9.23 (+355.10%) (+328.87%) <b>3.5X</b>	—	—
KenLM T.	2.99 (+86.10%) (+133.56%) <b>2.3X</b>	1.28 (-5.47%) (-21.84%)	3.44 (+86.30%) (+148.52%) <b>2.5X</b>	1.94 (+12.32%) (-10.01%)	—	—
Marisa	3.61 (+93.09%) (+161.60%) <b>2.8X</b>	2.06 (+52.00%) (+25.67%)	3.81 (+99.60%) (+114.95%) <b>2.7X</b>	3.24 (+87.96%) (+50.58%)	—	—
RandLM	1.81 (-3.06%) (+41.41%)	4.39 (+225.70%) (+168.04%) <b>2.5X</b>	2.02 (+6.18%) (+46.29%)	5.08 (+193.30%) (+135.82%) <b>2.5X</b>	2.60 (+70.73%) (+98.90%)	9.25 (+302.54%) (+302.19%) <b>3X</b>

- Elias-Fano Tries substantially **outperform ALL** previous solutions in **both space and time**.
- As fast as the state-of-the-art (KenLM) but more than twice smaller.

# Summary

Elias-Fano encodes *monotone integer sequences* in *space close to the information theoretic minimum*, while allowing *powerful search operations*, namely **Predecessor/Successor** queries and random **Access**.

Successfully applied to crucial problems, such as *inverted indexes*, *social graphs* and *tries* representation.

Several *optimized* software implementations are available.

- [Fano-1971] Robert Mario Fano. *On the number of bits required to implement an associative memory*. Memorandum 61, Computer Structures Group, MIT (1971).
- [Elias-1974] Peter Elias. *Efficient Storage and Retrieval by Content and Address of Static Files*. Journal of the ACM (JACM) 21, 2, 246–260 (1974).
- [Jacobson-1989] Guy Jacobson. *Succinct Static Data Structures*. Ph.D. Thesis, Carnegie Mellon University (1989).
- [Clark-1996] David Clark. *Compact Pat Trees*. Ph.D. Thesis, University of Waterloo (1996).
- [Moffat and Stuiver-2000] Alistair Moffat and Lang Stuiver. *Binary Interpolative Coding for Effective Index Compression*. Information Retrieval Journal 3, 1, 25–47 (2000).
- [Anh and Moffat-2005] Vo Ngoc Anh and Alistair Moffat. *Inverted Index Compression Using Word-Aligned Binary Codes*. Information Retrieval Journal 8, 1, 151–166 (2005).
- [Salomon-2007] David Salomon. *Variable-length Codes for Data Compression*. Springer (2007).
- [Vigna-2008] Sebastiano Vigna. *Broadword implementation of rank/select queries*. In Workshop in Experimental Algorithms (WEA), 154–168 (2008).

- [Yan *et al.*-2009] Hao Yan, Shuai Ding, and Torsten Suel. *Inverted index compression and query processing with optimized document ordering*. In Proceedings of the 18th International Conference on World Wide Web (WWW). 401–410 (2009).
- [Anh and Moffat-2010] Vo Ngoc Anh and Alistair Moffat. *Index compression using 64-bit words*. In Software: Practice and Experience 40, 2, 131–147 (2010).
- [Zukowski *et al.*-2010] Marcin Zukowski, Sandor Hèman, Niels Nes, and Peter Boncz. *Super-Scalar RAM-CPU Cache Compression*. In Proceedings of the 22nd International Conference on Data Engineering (ICDE). 59–70 (2006).
- [Stepanov *et al.*-2011] Alexander Stepanov, Anil Gangolli, Daniel Rose, Ryan Ernst, and Paramjit Oberoi. *SIMD-based decoding of posting lists*. In Proceedings of the 20th ACM International Conference on Information and Knowledge Management (CIKM). 317–326 (2011).
- [Zhou *et al.*-2013] Dong Zhou, David Andersen, Michael Kaminsky. *Space-Efficient, High-Performance Rank and Select Structures on Uncompressed Bit Sequences*. In Proceedings of the 12-nd International Symposium on Experimental Algorithms (SEA), 151-163 (2013).
- [Vigna-2013] Sebastiano Vigna. *Quasi-succinct indices*. In Proceedings of the 6-th ACM International Conference on Web Search and Data Mining (WSDM), 83-92 (2013).
- [Curtiss *et al.*-2013] Michael Curtiss *et al.* *Unicorn: A System for Searching the Social Graph*. In Proceedings of the Very Large Database Endowment (PVLDB), 1150-1161 (2013).
- [Ottaviano and Venturini-2014] Giuseppe Ottaviano, Rossano Venturini. *Partitioned Elias-Fano Indexes*. In Proceedings of the 37-th ACM International Conference on Research and Development in Information Retrieval (SIGIR), 273-282 (2014).

Thanks for your attention,  
time, patience!

Any questions?



# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$

1 2 3 4 5 6 7 8

$H = 1110111010001000$

$L = 011100111101110111101011$

# successor example

S = [3, 4, 7, 13, 14, 15, 21, 43]  
1 2 3 4 5 6 7 8

successor(12) = ?

H = 1110111010001000

L = 011100111101110111101011

# successor example

S = [3, 4, 7, 13, 14, 15, 21, 43]  
1 2 3 4 5 6 7 8

successor(12) = ?

001100

H = 1110111010001000

L = 011100111101110111101011

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

successor(12) = ?

$h_{12} = 001100$

$H = 1110111010001000$

$L = 011100111101110111101011$

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

successor(12) = ?

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$
$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = 1110111010001000$

$L = 011100111101110111101011$

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

successor(12) = ?

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$

$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = \underline{1110}111010001000$

$L = 011100111101110111101011$

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

successor(12) = ?

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$
$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = \underline{1110}1110\underline{10001000}$

$L = 011100111101110111101011$

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

successor(12) = ?

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$

$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = \underline{1110}1110\underline{10001000}$

$L = \underline{011100111}101110111\underline{101011}$



$p_1$



$p_2$



# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

successor(12) = ?

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$
$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = \underline{1110}1110\underline{10001000}$

$L = \underline{011100111}101110111\underline{101011}$

↑  
 $p_1$

↑  
 $p_2$



*binary search*  
in  $[p_1, p_2)$

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{successor}(12) = 13$

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$
$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = \underline{1110}1110\underline{10001000}$

$L = \underline{011100111}101110111\underline{101011}$

$\uparrow$   
 $p_1$

$\uparrow$   
 $p_2$



*binary search*  
in  $[p_1, p_2)$

# successor example

$S = [3, 4, 7, 13, 14, 15, 21, 43]$   
1 2 3 4 5 6 7 8

$\text{successor}(12) = 13$

$h_{12} = 001100$

$$p_1 = \text{select}_0(h_x) - h_x$$
$$p_2 = \text{select}_0(h_{x+1}) - h_x - 1$$

$H = \underline{1110}1110\underline{10001000}$   
 $L = \underline{011100111}101110111\underline{101011}$   
                                  ↑                                  ↑  
                                   $p_1$                                    $p_2$



*binary search*  
in  $[p_1, p_2)$

$$\text{Complexity: } O\left(\log \frac{u}{n}\right)$$