

On Elias-Fano for Rank Queries in FM-Indexes

Danyang Ma, Simon J. Puglisi, Rajeev Raman, and Bella Zhukova

DCC 2021

Definitions: BWT

$T = \text{mississippi\$}$

$\Sigma = \{\$, i, m, p, s\}$

| | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|
| m | i | s | s | i | s | s | i | p | p | i | \$ |
| i | s | s | i | s | s | i | p | p | i | \$ | m |
| s | s | i | s | s | i | p | p | i | \$ | m | i |
| s | i | s | s | i | p | p | i | \$ | m | i | s |
| i | s | s | i | p | p | i | \$ | m | i | s | s |
| s | s | i | p | p | i | \$ | m | i | s | s | i |
| s | i | p | p | i | \$ | m | i | s | s | i | s |
| i | p | p | i | \$ | m | i | s | s | i | s | s |
| p | p | i | \$ | m | i | s | s | i | s | s | i |
| p | i | \$ | m | i | s | s | i | s | s | i | p |
| i | \$ | m | i | s | s | i | s | s | i | p | p |
| \$ | m | i | s | s | i | s | s | i | p | p | i |

Definitions: BWT

$T = \text{mississippi\$}$

$\Sigma = \{\$, i, m, p, s\}$

| | | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 5 | m | i | s | s | i | s | s | i | p | p | i | \$ |
| 4 | i | s | s | i | s | s | i | p | p | i | \$ | m |
| 11 | s | s | i | s | s | i | p | p | i | \$ | m | i |
| 9 | s | i | s | s | i | p | p | i | \$ | m | i | s |
| 3 | i | s | s | i | p | p | i | \$ | m | i | s | s |
| 10 | s | s | i | p | p | i | \$ | m | i | s | s | i |
| 8 | s | i | p | p | i | \$ | m | i | s | s | i | s |
| 2 | i | p | p | i | \$ | m | i | s | s | i | s | s |
| 7 | p | p | i | \$ | m | i | s | s | i | s | s | i |
| 6 | p | i | \$ | m | i | s | s | i | s | s | i | p |
| 1 | i | \$ | m | i | s | s | i | s | s | i | p | p |
| 0 | \$ | m | i | s | s | i | s | s | i | p | p | i |

Definitions: BWT

$T = \text{mississippi}\$$

$\Sigma = \{\$, i, m, p, s\}$

| | | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0 | \$ | m | i | s | s | i | s | s | i | p | p | i |
| 1 | i | \$ | m | i | s | s | i | s | s | i | p | p |
| 2 | i | p | p | i | \$ | m | i | s | s | i | s | s |
| 3 | i | s | s | i | p | p | i | \$ | m | i | s | s |
| 4 | i | s | s | i | s | s | i | p | p | i | \$ | m |
| 5 | m | i | s | s | i | s | s | i | p | p | i | \$ |
| 6 | p | i | \$ | m | i | s | s | i | s | s | i | p |
| 7 | p | p | i | \$ | m | i | s | s | i | s | s | i |
| 8 | s | i | p | p | i | \$ | m | i | s | s | i | s |
| 9 | s | i | s | s | i | p | p | i | \$ | m | i | s |
| 10 | s | s | i | p | p | i | \$ | m | i | s | s | i |
| 11 | s | s | i | s | s | i | p | p | i | \$ | m | i |

Definitions: BWT

$T = \text{mississippi}\$$

$\Sigma = \{\$, i, m, p, s\}$

F

| | | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0 | \$ | m | i | s | s | i | s | s | i | p | p | i |
| 1 | i | \$ | m | i | s | s | i | s | s | i | p | p |
| 2 | i | p | p | i | \$ | m | i | s | s | i | s | s |
| 3 | i | s | s | i | p | p | i | \$ | m | i | s | s |
| 4 | i | s | s | i | s | s | i | p | p | i | \$ | m |
| 5 | m | i | s | s | i | s | s | i | p | p | i | \$ |
| 6 | p | i | \$ | m | i | s | s | i | s | s | i | p |
| 7 | p | p | i | \$ | m | i | s | s | i | s | s | i |
| 8 | s | i | p | p | i | \$ | m | i | s | s | i | s |
| 9 | s | i | s | s | i | p | p | i | \$ | m | i | s |
| 10 | s | s | i | p | p | i | \$ | m | i | s | s | i |
| 11 | s | s | i | s | s | i | p | p | i | \$ | m | i |

Definitions: BWT

$T = \text{mississippi}\$$

$F \quad \Sigma = \{\$, i, m, p, s\}$

L, BWT

| | | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0 | \$ | m | i | s | s | i | s | s | i | p | p | i |
| 1 | i | \$ | m | i | s | s | i | s | s | i | p | p |
| 2 | i | p | p | i | \$ | m | i | s | s | i | s | s |
| 3 | i | s | s | i | p | p | i | \$ | m | i | s | s |
| 4 | i | s | s | i | s | s | i | p | p | i | \$ | m |
| 5 | m | i | s | s | i | s | s | i | p | p | i | \$ |
| 6 | p | i | \$ | m | i | s | s | i | s | s | i | p |
| 7 | p | p | i | \$ | m | i | s | s | i | s | s | i |
| 8 | s | i | p | p | i | \$ | m | i | s | s | i | s |
| 9 | s | i | s | s | i | p | p | i | \$ | m | i | s |
| 10 | s | s | i | p | p | i | \$ | m | i | s | s | i |
| 11 | s | s | i | s | s | i | p | p | i | \$ | m | i |

Problem (Counting Queries)

For a pattern P and text T , $count(P)$ returns the number of occurrences of P in T

Definitions: *rank* and *select*

$rank_X(i, c)$ on string X

is the number of occurrences of symbol c in prefix $X[0, i-1]$

$select_X(i, c)$ on string X

is the position of the i -th occurrence of symbol c

Definitions: BWT

$T = \text{mississippi\$}$

$\Sigma = \{\$, i, m, p, s\}$

$P = \text{isi}$

F L, BWT

| | | |
|----|----|----|
| 0 | \$ | i |
| 1 | i | p |
| 2 | i | s |
| 3 | i | s |
| 4 | i | m |
| 5 | m | \$ |
| 6 | p | p |
| 7 | p | i |
| 8 | s | s |
| 9 | s | s |
| 10 | s | i |
| 11 | s | i |

| | | | | | | | | | | | | |
|-----|---|---|---|---|---|---|---|---|---|---|----|----|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| \$: | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| i: | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |
| m: | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| p: | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| s: | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |

Previous Work

- Many authors at this bit-vector representation on BWT (Mäkinen, Navarro, Huo, Sirén, Grabowski,...)
- We re-examine the use of Elias-Fano for representing bit-vectors

Our Contribution

- Optimizations on Elias-Fano that quicken rank query
 - may be useful in other problems/contexts, not only for counting queries
- Two Elias-Fano-like variants independent of `select`

Elias-Fano

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|-----|---|---|---|---|---|---|---|---|---|---|----|----|
| \$: | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| i: | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |
| m: | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| p: | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| s: | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |

Elias-Fano

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|-----|---|---|---|---|---|---|---|---|---|---|----|----|
| \$: | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| i: | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |
| m: | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| p: | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| s: | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |

$B = 000001000000 \ 100000010011 \ 000010000000 \ 010000100000 \ 001100001100$

Elias-Fano

$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 001100001100$

5 12 19 22 23 28 37 42 50 51 56 57

Elias-Fano

$$m = 57, n = 12, l = 2$$

5 12 19 22 23 28 37 42 50 51 56 57

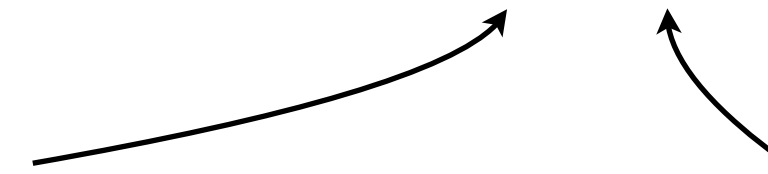
Elias-Fano

$$m = 57, n = 12, l = 2$$

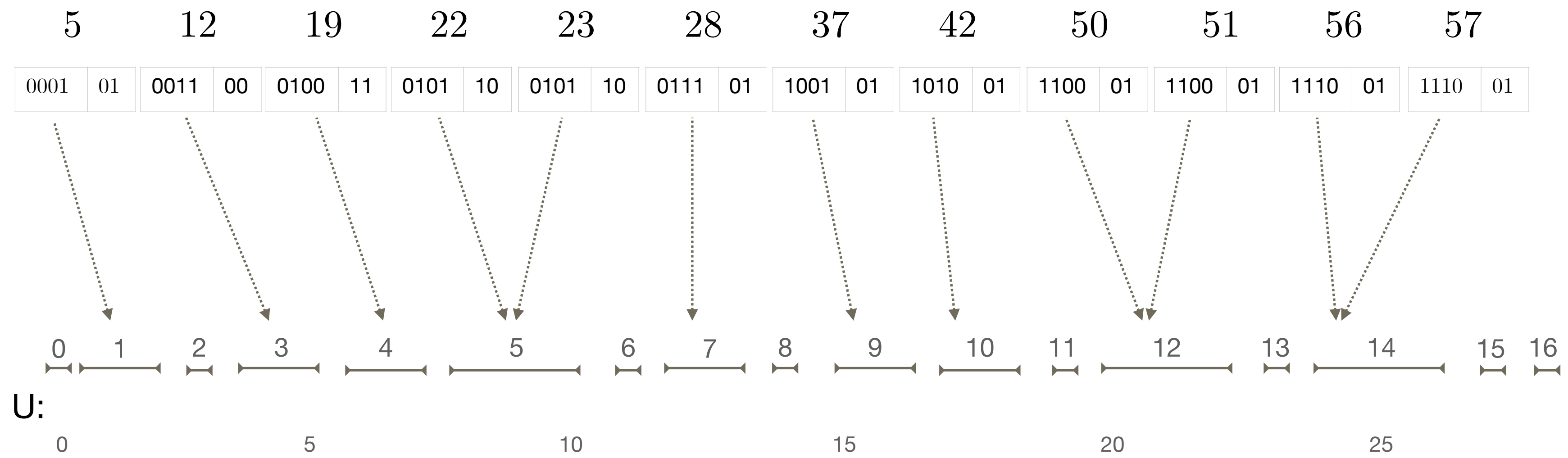
| | | | | | | | | | | | |
|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| 5 | 12 | 19 | 22 | 23 | 28 | 37 | 42 | 50 | 51 | 56 | 57 |
| 0001 01 | 0011 00 | 0100 11 | 0101 10 | 0101 10 | 0111 01 | 1001 01 | 1010 01 | 1100 01 | 1100 01 | 1110 01 | 1110 01 |

quotient: $\lfloor i/2^l \rfloor$, $\lceil \log_2 m \rceil - l$ bits

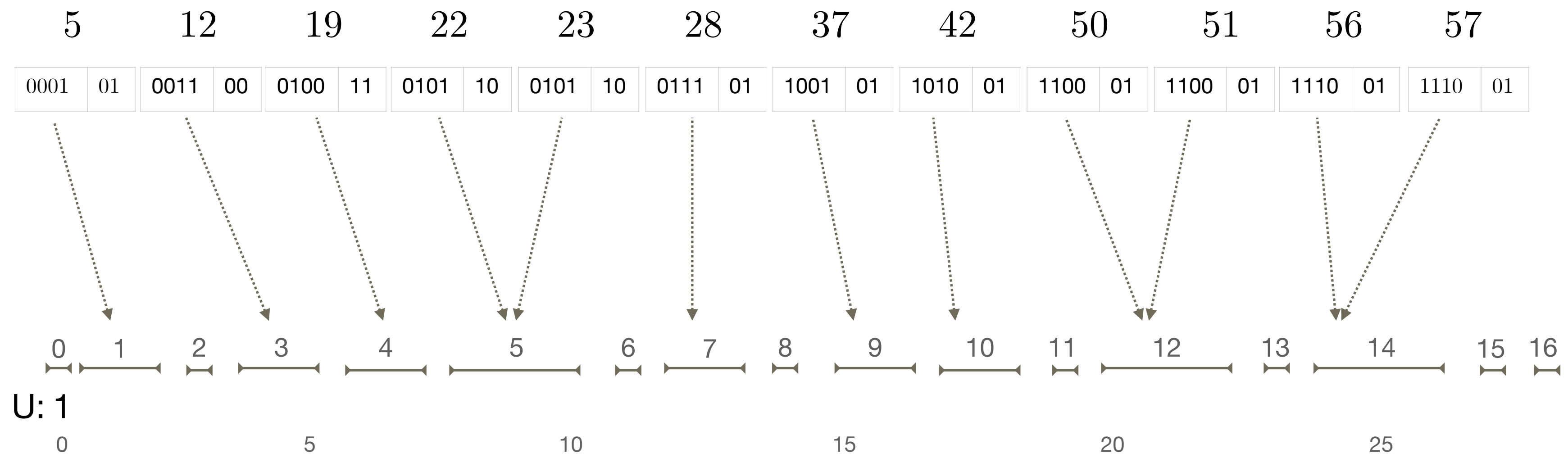
remainder: $i \bmod 2^l$, l bits



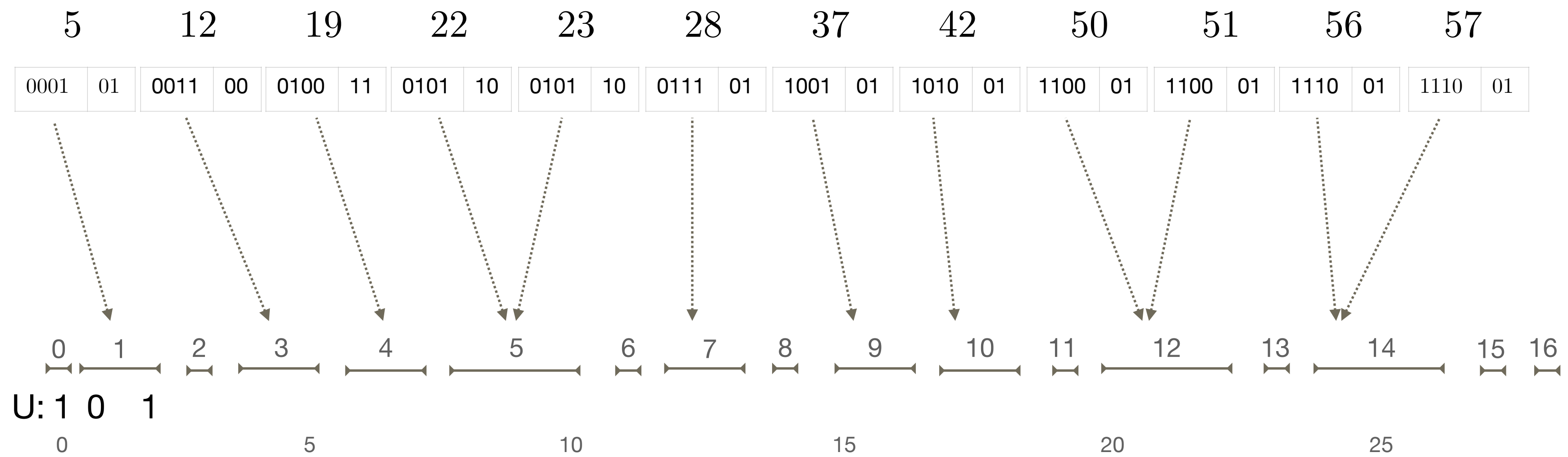
Elias-Fano



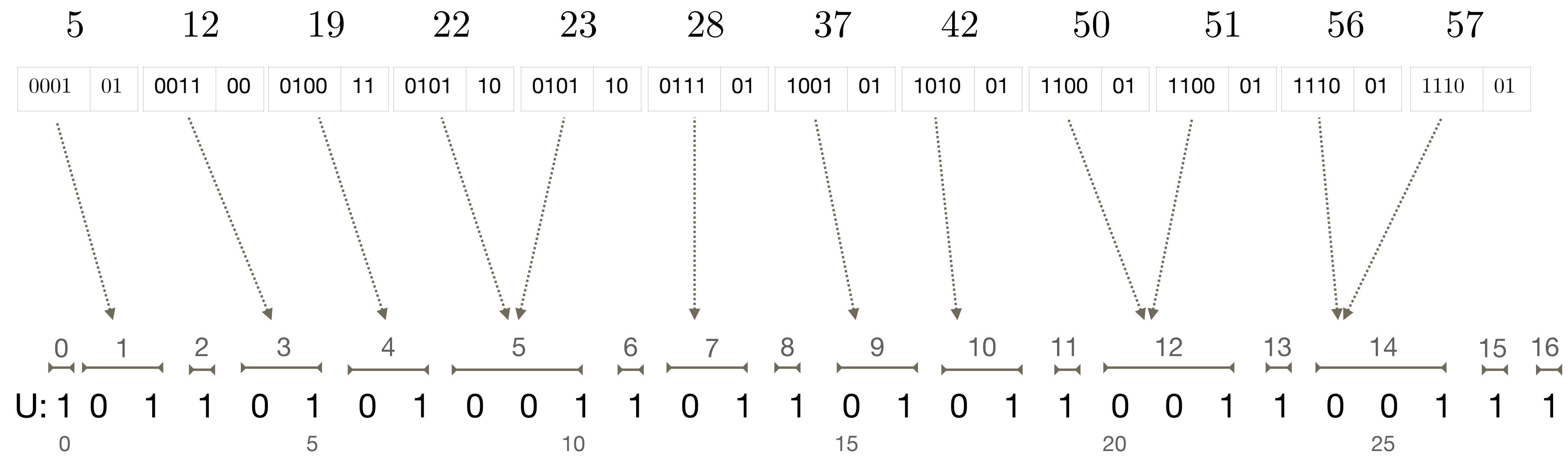
Elias-Fano



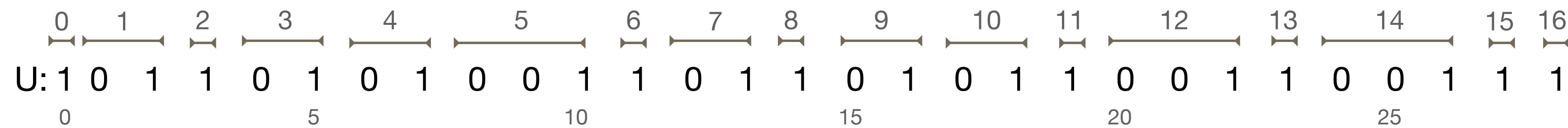
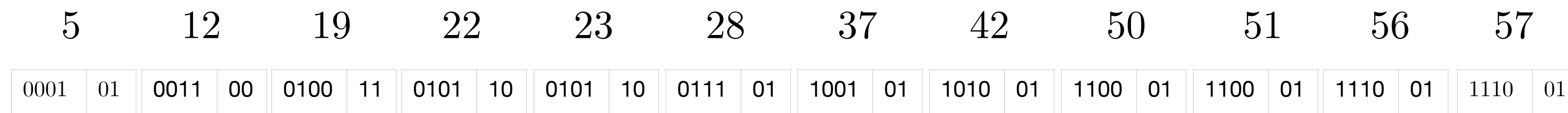
Elias-Fano



Elias-Fano



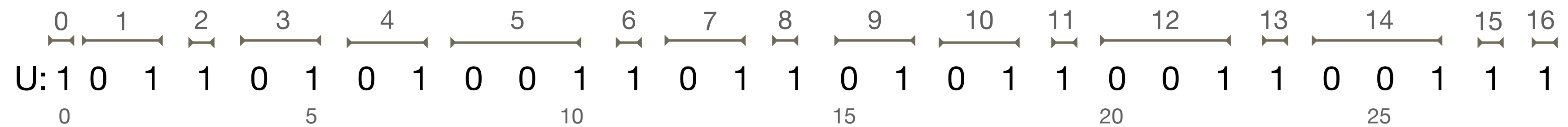
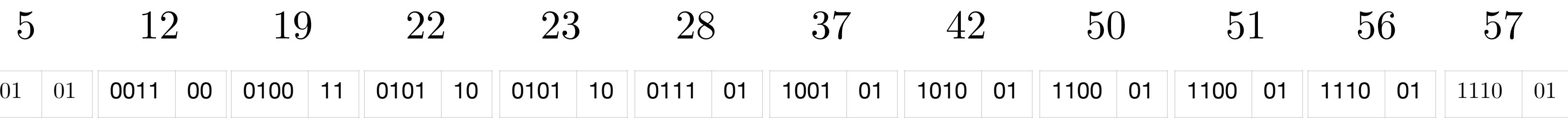
Elias-Fano



Elias-Fano



$B = 000001000000\ 100000010011\ 000010000000\ 0100000100000\ 0011000001100$



$38 = 1001\ 10, 1001 \implies 9$

Experiment

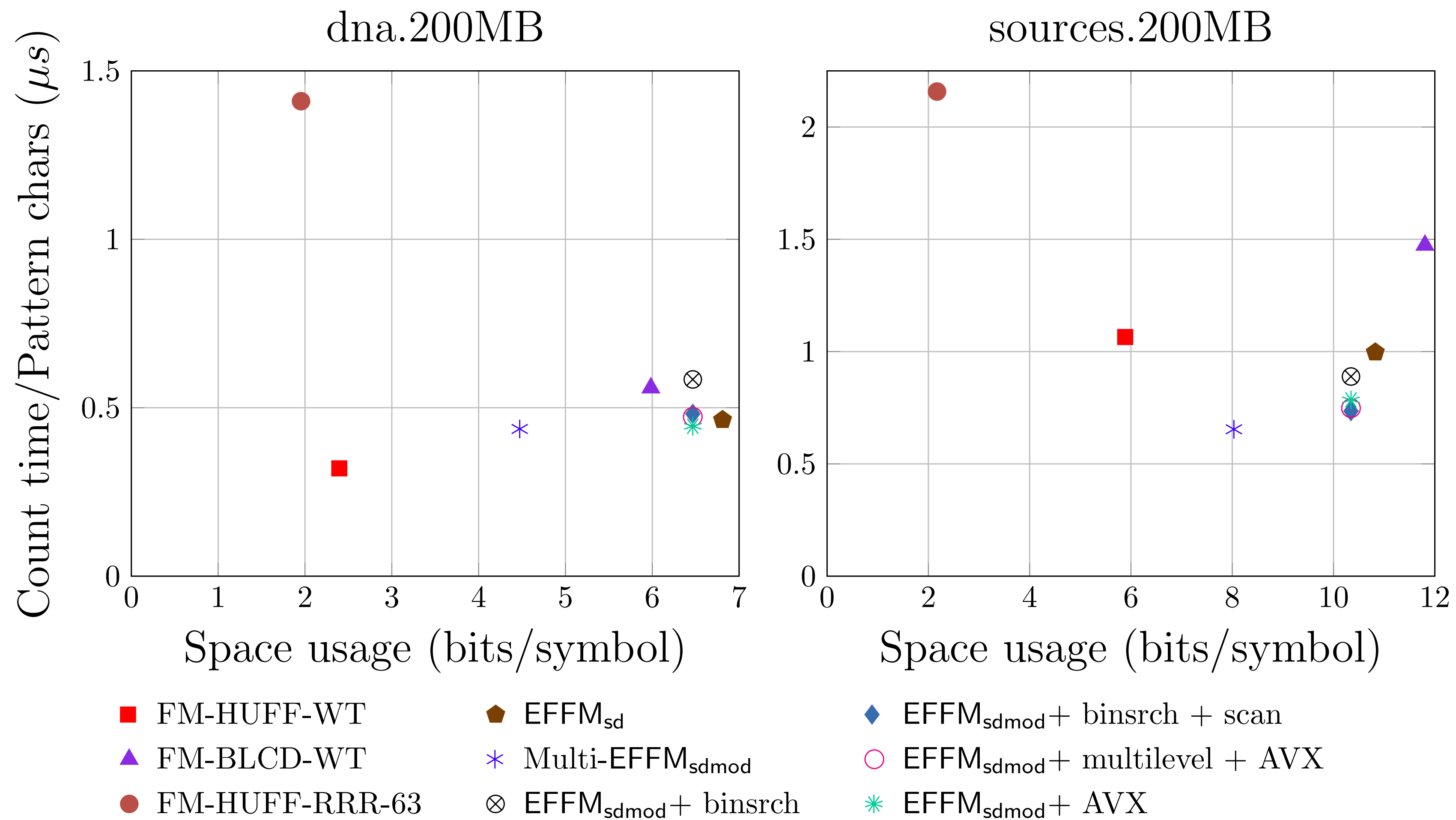
Datasets, 200MB, from the [Pizza&Chili corpus](#):

- DNA
- english
- sources
- XML

Search queries:

- 50 000 patterns
- length = 20

Elias-Fano



New Approaches

- $EFFM_{zs}$
 - Zero-Suppressed
- $EFFM_{pu}$
 - Partitioned-Upper

New Approaches: EFFM_{ZS}

$B = 000001000000 \ 100000010011 \ 000010000000 \ 010000100000 \ 001100001100 \quad \Sigma = \{\$, i, m, p, s\}$

let $l = 2$

U bit vector of length $(n\sigma)/2^\ell$, and $U[i] = 1$ iff the i -th bucket is not all zeros

L $2^\ell \cdot$ number of non-zero buckets

New Approaches: E_zFFM_zS

$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 001100001100$ $\Sigma = \{\$, i, m, p, s\}$

let $l = 2$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

U



New Approaches: EFFM_{ZS}

$B = 000001000000 100000010011 000010000000 010000100000 001100001100$ $\Sigma = \{\$, i, m, p, s\}$
let $l = 2$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

U



New Approaches: EFFM_{ZS}

$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 001100001100$ $\Sigma = \{\$, i, m, p, s\}$
let $l = 2$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

U

| |
|---|
| 0 |
|---|

New Approaches: EFFM_{ZS}

$B = 0000$ 0100 $0000 100000010011 000010000000 010000100000 001100001100$ $\Sigma = \{\$, i, m, p, s\}$
let $l = 2$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

U

| |
|---|
| 0 |
|---|

New Approaches: EFFM_{ZS}

$B = 0000$ 0100 $0000 100000010011 000010000000 010000100000 001100001100$ $\Sigma = \{\$, i, m, p, s\}$
let $l = 2$



New Approaches: EFFM_{ZS}

$B = 000001000000\ 1000\boxed{0001}0011\ 000010000000\ 010000100000\ 001100001100$ $\Sigma = \{\$, i, m, p, s\}$
let $l = 2$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

U

| | | | | |
|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 |
|---|---|---|---|---|

New Approaches: E_zFFM_zS

$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 001100001100$ $\Sigma = \{\$, i, m, p, s\}$

let $l = 2$

| | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| U | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |

New Approaches: E_{ZS}FFM

$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 001100001100$ $\Sigma = \{\$, i, m, p, s\}$

let $l = 2$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35

U

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

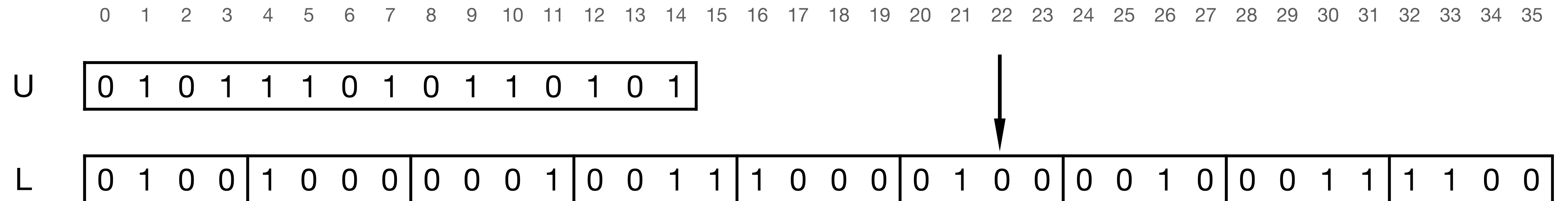
L

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

New Approaches: EFFM_{ZS}

$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 0011000001100$ $\Sigma = \{\$, i, m, p, s\}$

let $l = 2$



$$\begin{aligned} \text{rank}_B(i, 1): & \quad j = \text{rank}_U(\lfloor i/2^\ell \rfloor, 1) \\ \text{if } U[\lfloor i/2^\ell \rfloor] = 0 & \quad \text{rank}_B(i, 1) = \text{rank}_L(j \cdot 2^\ell, 1) \\ \text{if } U[\lfloor i/2^\ell \rfloor] = 1 & \quad \text{rank}_B(i, 1) = \text{rank}_L(j \cdot 2^\ell + i \bmod 2^\ell, 1) \end{aligned}$$

New Approaches: EFFM_{pu}



$B = 000001000000\ 100000010011\ 000010000000\ 010000100000\ 001100001100$ $\Sigma = \{\$, i, m, p, s\}$
 $l = 2$

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|----------|------------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| U_{nz} | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
| U_{sz} | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | | | |
| L | as in sdsI | | | | | | | | | | | | | | |

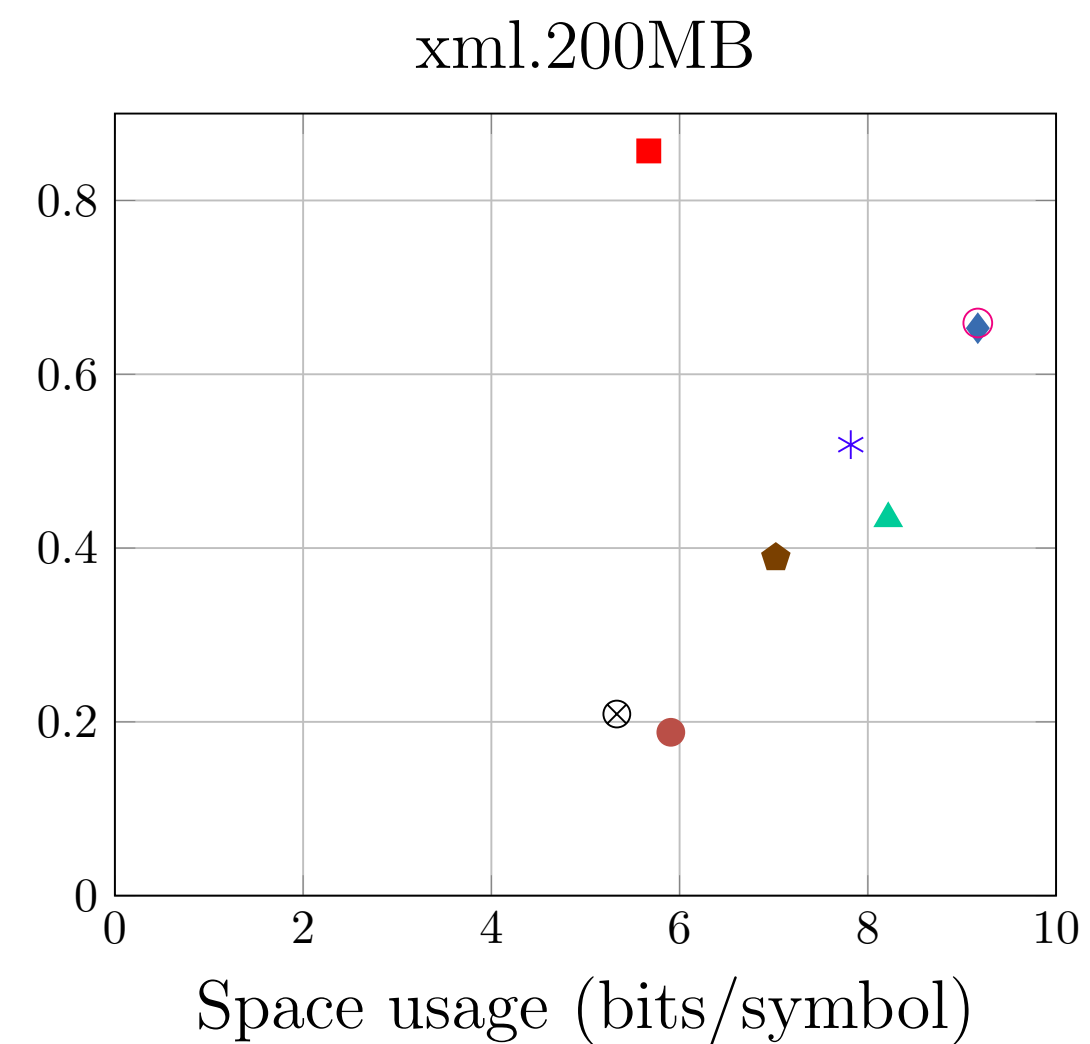
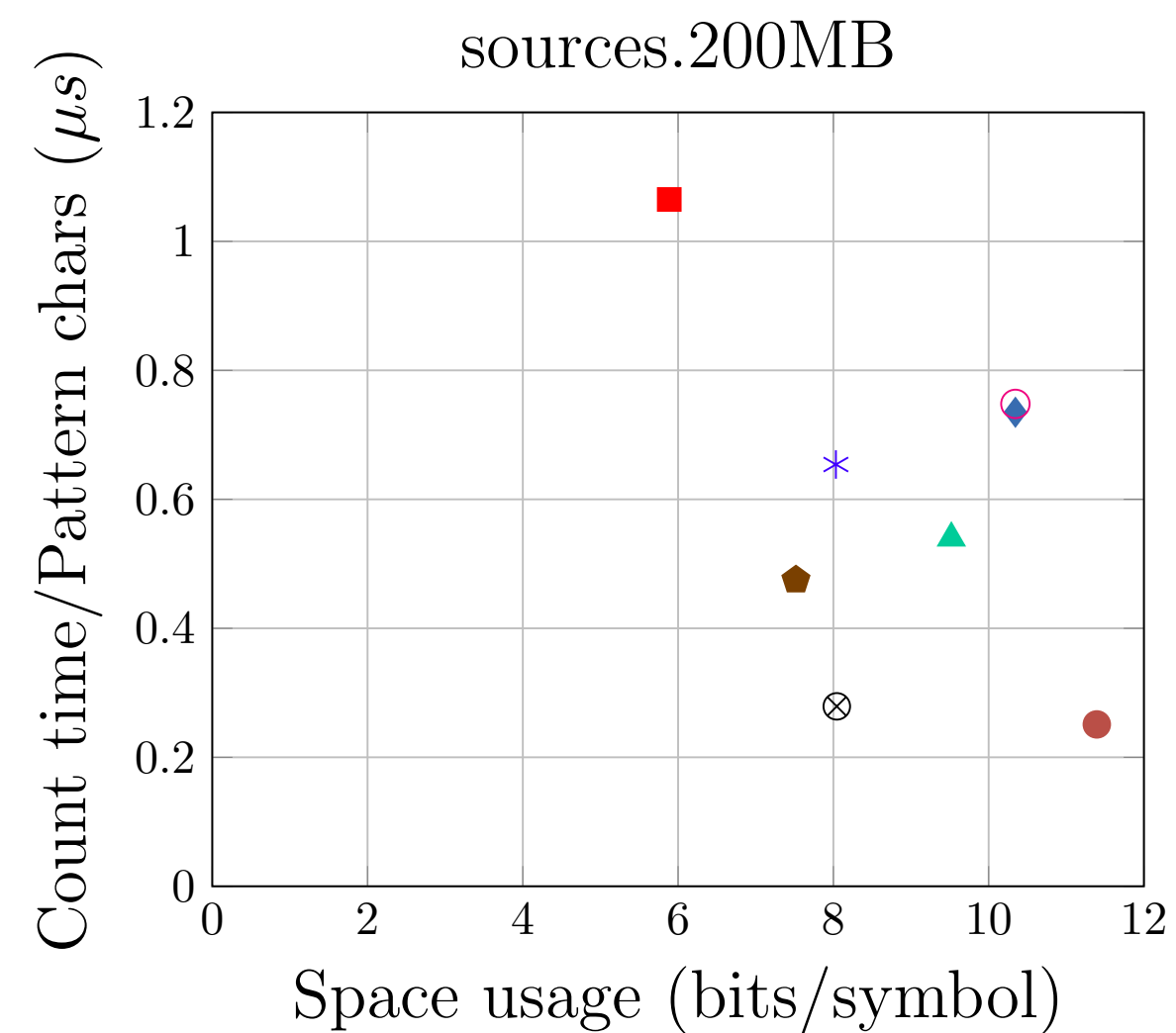
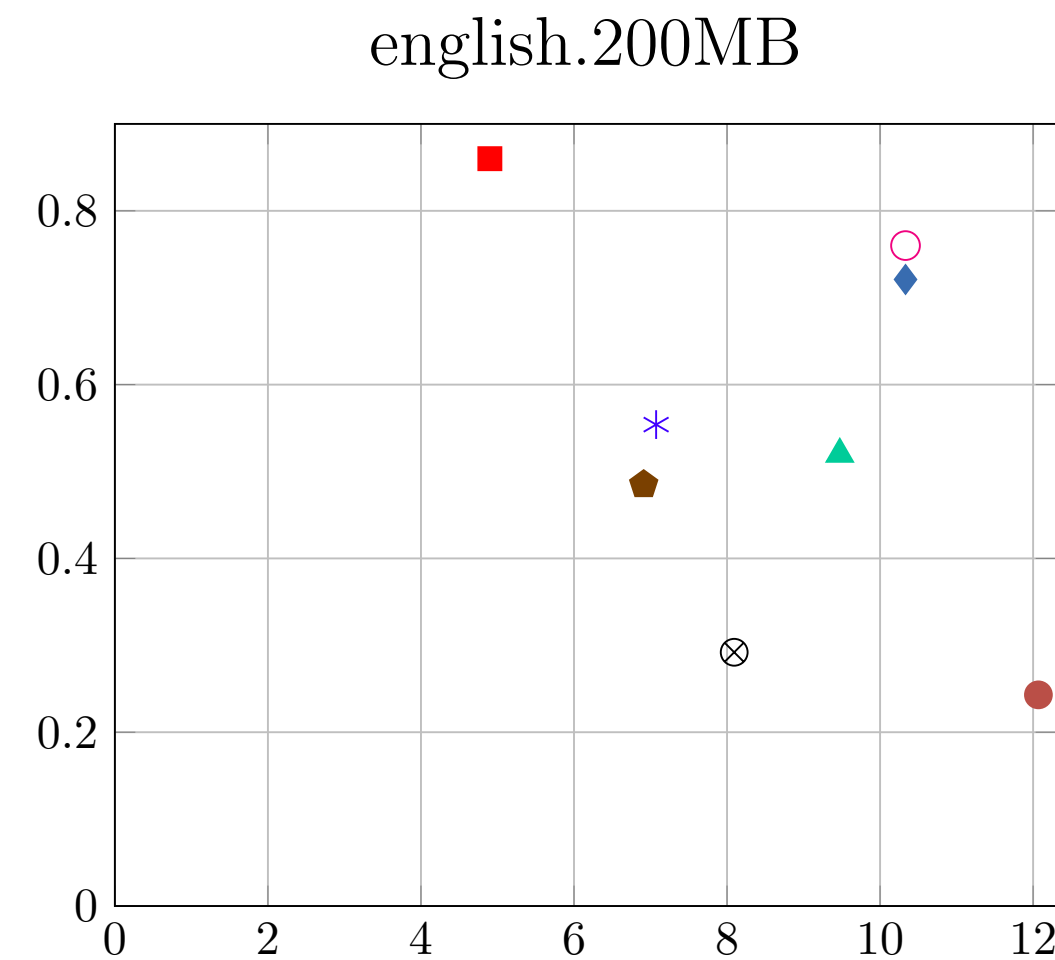
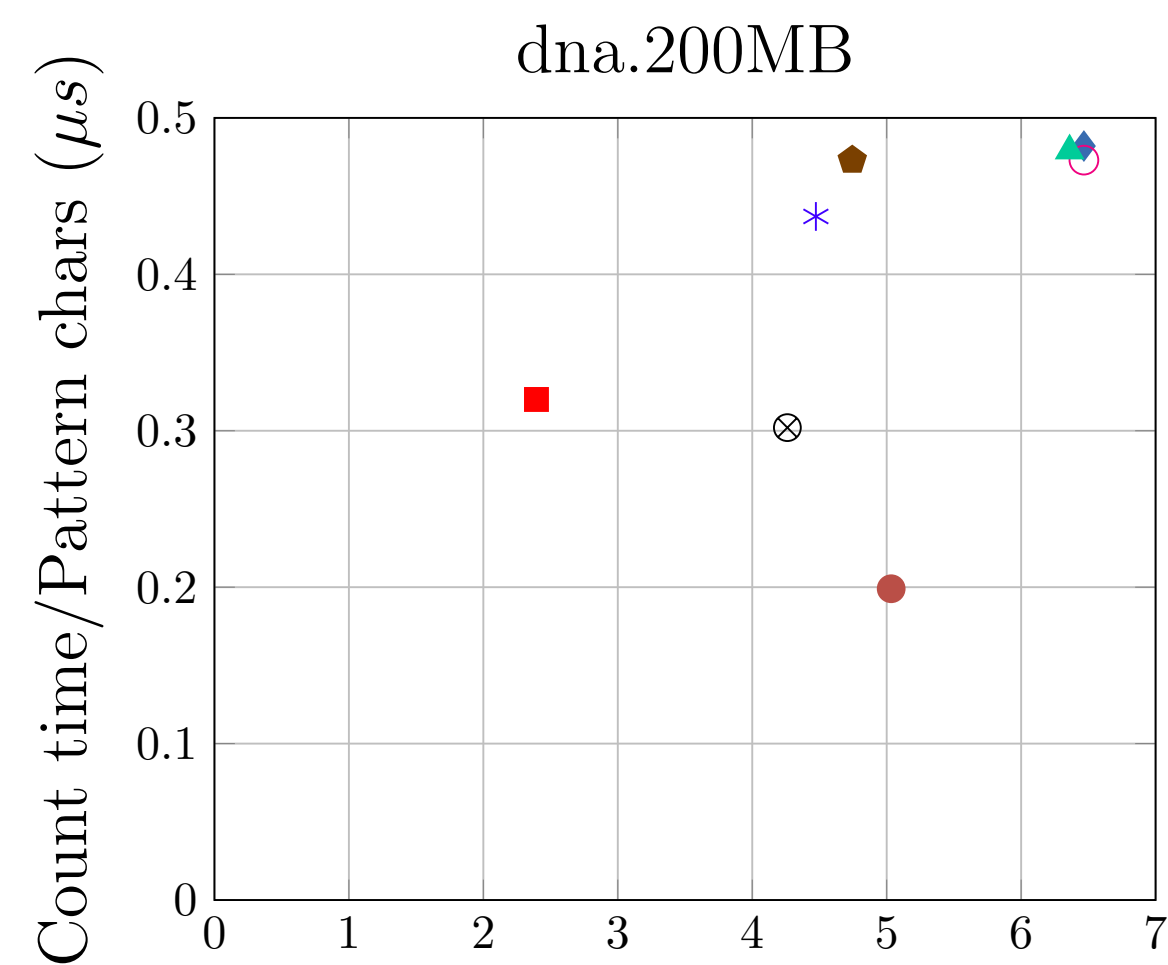
$$\text{select}_U(\lfloor i/2^\ell \rfloor, 1) = \text{select}_{U_{sz}}(j, 1), \text{ where } j = \text{rank}_{U_{nz}}(\lfloor i/2^\ell \rfloor, 1)$$

New Approaches: EFFM_{pu}

$B = 000001000000 \ 100000010011 \ 000010000000 \ 010000100000 \ 0011000001100 \quad \Sigma = \{\$, i, m, p, s\}$
 $l = 2$

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|----------|---|---|---|---|---|---|---|---|---|----------------------|----|----|----|----|----|
| U_{nz} | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
| U_{sz} | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | | | |
| | 1 | 1 | 1 | 2 | 1 | 1 | 1 | 2 | 2 | → l bits per entry | | | | | |

New Approaches



- FM-HUFF-WT
- ◆ EFFM_{sdmod} + binsrch + scan
- EFFM_{sdmod} + multilevel + AVX
- * Multi-EFFM_{sdmod}
- EFFM_{zs}
- ⊗ Multi-EFFM_{zs}
- ▲ EFFM_{pu}
- ◆ Multi-EFFM_{pu}

Future Work

- Explore the use of our indexes for querying labelled trees and automata
- Exploit AVX instructions more
 - AVX still haven't been used much in succinct data structures
- Using our new EF implementations in other applications where EF is currently used

Thank you!