# Burrows-Wheeler Transform

CMSC 423

# Motivation - Short Read Mapping

A Cow Genome

Sequencing technologies produce millions of "reads" = a random, short substring of the genome

If we already know the genome of one cow, we can get reads from a 2nd cow and map them onto the known cow genome.

Need to do millions of string searches in a long string.

# Bowtie

**Software**

## Ultrafast and memory-efficient alignment of short DNA sequences to the human genome

**Ben Langmead**[*], **Cole Trapnell**, **Mihai Pop** and **Steven L Salzberg**

* Corresponding author: Ben Langmead langmead@cs.umd.edu

▼ Author Affiliations

Center for Bioinformatics and Computational Biology, Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA
For all author emails, please log on.

The electronic version of this article is the complete one and can be found online at:
http://genomebiology.com/2009/10/3/R25

# BWA

## Fast and accurate short read alignment with Burrows–Wheeler transform

→

**Heng Li** and **Richard Durbin**[*]

+ Author Affiliations

* To whom correspondence should be addressed.

# Bowtie Performance

**Varying read length using Bowtie, Maq and SOAP**

| Length | Program | CPU time | Wall clock time | Peak virtual memory footprint (megabytes) | Bowtie speed-up | Reads aligned (%) |
|--------|---------|----------|-----------------|-------------------------------------------|-----------------|-------------------|
| 36 bp | Bowtie | 6 m 15 s | 6 m 21 s | 1,305 | - | 62.2 |
| | Maq | 3 h 52 m 26 s | 3 h 52 m 54 s | 804 | 36.7× | 65.0 |
| | Bowtie -v 2 | 4 m 55 s | 5 m 00 s | 1,138 | - | 55.0 |
| | SOAP | 16 h 44 m 3 s | 18 h 1 m 38 s | 13,619 | 216× | 55.1 |
| 50 bp | Bowtie | 7 m 11 s | 7 m 20 s | 1,310 | - | 67.5 |
| | Maq | 2 h 39 m 56 s | 2 h 40 m 9 s | 804 | 21.8× | 67.9 |
| | Bowtie -v 2 | 5 m 32 s | 5 m 46 s | 1,138 | - | 56.2 |
| | SOAP | 48 h 42 m 4 s | 66 h 26 m 53 s | 13,619 | 691× | 56.2 |
| 76 bp | Bowtie | 18 m 58 s | 19 m 6 s | 1,323 | - | 44.5 |
| | Maq 0.7.1 | 4 h 45 m 7 s | 4 h 45 m 17 s | 1,155 | 14.9× | 44.9 |
| | Bowtie -v 2 | 7 m 35 s | 7 m 40 s | 1,138 | - | 31.7 |

Maq & SOAP build hash table of locations of k-mers

The performance of Bowtie v0.9.6, SOAP v1.10, and Maq versions v0.6.6 and v0.7.1 on the server platform when aligning 2 M untrimmed reads from the 1,000 Genome project (National Center for Biotechnology Information Short Read Archive: SRR003084 for 36 base pairs [bp], SRR003092 for 50 bp, and SRR003196 for 76 bp). For each read length, the 2 M reads were randomly sampled from the FASTQ file downloaded from the Archive such that the average per-base error rate as measured by quality values was uniform across the three sets. All reads pass through Maq's "catfilter". Maq v0.7.1 was used for the 76-bp reads because v0.6.6 does not support reads longer than 63 bp. SOAP is excluded from the 76-bp experiment because it does not support reads longer than 60 bp. Other experimental parameters are identical to those of the experiments in Table 1. CPU, central processing unit.

Langmead et al. (2008)

# Burrows-Wheeler Transform

Text transform that is useful for compression & search.

banana

banana$
anana$b
nana$ba
ana$ban
na$bana
a$banan
$banana

→ sort →

$banana
a$banan
ana$ban
anana$b
banana$
nana$ba
na$bana

BWT(banana) =
annb$aa

Tends to put runs of the same character together.

Makes compression work well.

"bzip" is based on this.

# Another Example

appellee$

appellee$
ppellee$a
pellee$ap
ellee$app
llee$appe         **sort**
lee$appel   ⟶
ee$appell
e$appelle
$appellee

$appellee
appellee$
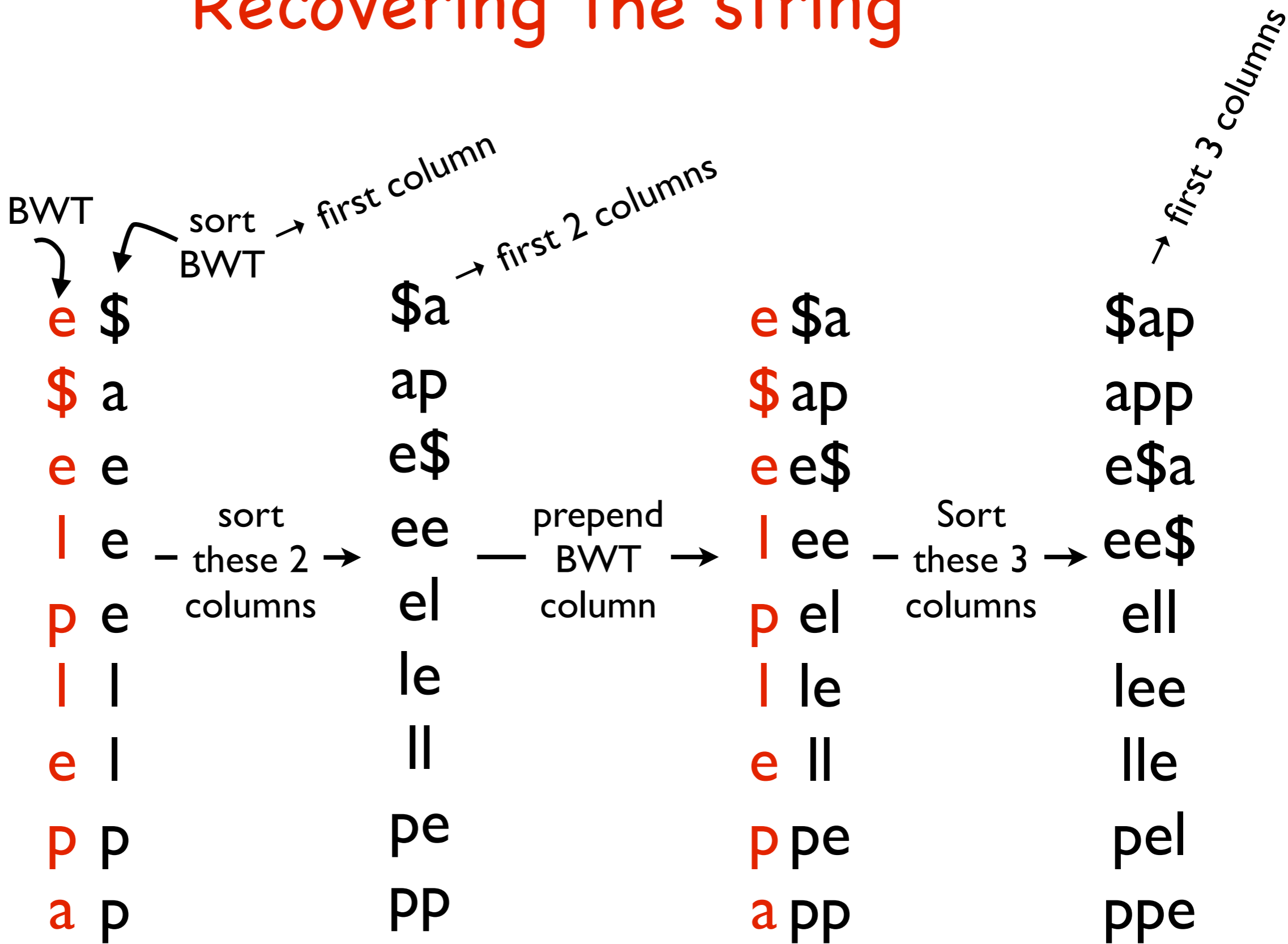e$appelle
ee$appell
ellee$app
lee$appel
llee$appe
pellee$ap
ppellee$a

BWT(appellee$) =
e$elplepa

Doesn't always improve
the compressibility...

# Recovering the string

BWT

sort BWT → first column

→ first 2 columns

→ first 3 columns

$appellee
appellee$
e$appelle
ee$appell
ellee$app
lee$appel
llee$appe
pellee$ap
ppellee$a

| e | $ |
|---|---|
| $ | a |
| e | e |
| l | e |
| P | e |
| l | l |
| e | l |
| P | p |
| a | p |

— sort these 2 columns →

$a
ap
e$
ee
el
le
ll
pe
pp

— prepend BWT column →

| e | $a |
|---|---|
| $ | ap |
| e | e$ |
| l | ee |
| P | el |
| l | le |
| e | ll |
| P | pe |
| a | pp |

— Sort these 3 columns →

$ap
app
e$a
ee$
ell
lee
lle
pel
ppe

# Inverse BWT

```
def inverseBWT(s):
    B = [s₁,s₂,s₃,...,sₙ]
    for i = 1..n:
        sort B
        prepend sᵢ to B[i]
    return row of B that ends with $
```

# Another BWT Example

| | | |
|---|---|---|
| dogwood$ | | $dogwoo**d** |
| ogwood$d | | d$dogwo**o** |
| gwood$do | | dogwood**$** |
| wood$dog | sort → | gwood$d**o** |
| ood$dogw | | od$dogw**o** |
| od$dogwo | | ogwood$**d** |
| d$dogwoo | | ood$dog**w** |
| $dogwood | | wood$do**g** |

last column →

BWT(dogwood$) =
do$oodwg

# do$oodwg   Another BWT Example

| Prepend | Sort | Prepend | Sort | Prepend | Sort | Prepend | Sort |
|---|---|---|---|---|---|---|---|
| d $ | $d | d $d | $do | d$do | $dog | d $dog | $dogw |
| o d | d$ | o d$ | d$d | od$d | d$do | o d$do | d$dog |
| $ d | do | $ do | dog | $ dog | dogw | $ dogw | dogwo |
| o g | gw | o gw | gwo | ogwo | gwoo | o gwoo | gwood |
| o o | od | o od | od$ | ood$ | od$d | o od$d | od$do |
| d o | og | d og | ogw | dogw | ogwo | d ogwo | ogwoo |
| w o | oo | w oo | ood | wood | ood$ | w ood$ | ood$d |
| g w | wo | g wo | woo | gwoo | wood | g wood | wood$ |

| Prepend | Sort | Prepend | Sort | Prepend | Sort |
|---|---|---|---|---|---|
| d $dogw | $dogwo | d $dogwo | $dogwoo | d $dogwoo | $dogwood |
| o d$dog | d$dogw | o d$dogw | d$dogwo | o d$dogwo | d$dogwoo |
| $ dogwo | dogwoo | $ dogwoo | dogwood | $ dogwood | dogwood$ |
| o gwood | gwood$ | o gwood$ | gwood$d | o gwood$d | gwood$do |
| o od$do | od$dog | o od$dog | od$dogw | o od$dogw | od$dogwo |
| d ogwoo | ogwood | d ogwood | ogwood$ | d ogwood$ | ogwood$d |
| w ood$d | ood$do | w ood$do | ood$dog | w ood$dog | ood$dogw |
| g wood$ | wood$d | g wood$d | wood$do | g wood$do | wood$dog |

# Searching with BWT: LF Mapping

LF Mapping

BWT(unabashable)

Σ

| | $ | a | b | e | h | l | n | s | u |
|---|---|---|---|---|---|---|---|---|---|
| $unabashable | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| abashable$un | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| able$unabash | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| ashable$unab | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| bashable$una | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| ble$unabasha | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| e$unabashabl | 0 | 2 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| hable$unabas | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| le$unabashab | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 0 |
| nabashable$u | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 0 |
| shable$unaba | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| unabashable$ | 0 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 1 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |

# of times letter appears before this position in the last column.

**LF Property:** The $i^{th}$ occurrence of a letter X in the last column corresponds to the $i^{th}$ occurrence of X in the first column.

# BWT Search

BWTSearch(aba)  Start from the **end** of the pattern

Step 1: Find the range of "a"s in the first column

Step 2: Look at the same range in the last column.

Step 3: "b" is the next pattern character. Set B = the LF mapping entry for b in the first row of the range.

Set E = the LF mapping entry for b in the last + 1 row of the range.

Step 4: Find the range for "b" in the first row, and use B and E to find the right subrange within the "b" range.

BWT(unabashable)

LF Mapping

Σ

| | $ | a | b | e | h | l | n | s | u |
|---|---|---|---|---|---|---|---|---|---|
| $unabashable | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| abashable$un | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| able$unabash | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| ashable$unab | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| bashable$una | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| ble$unabasha | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| e$unabashabl | 0 | 2 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| hable$unabas | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| le$unabashab | 0 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 0 |
| nabashable$u | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 0 |
| shable$unaba | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| unabashable$ | 0 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 1 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |

# BWT Searching Example 2

pattern = "bana"



**Top-left table** (label: a)

| | $ | a | b | n |
|---|---|---|---|---|
| $bananna | 0 | 0 | 0 | 0 |
| a$banann | 0 | 1 | 0 | 0 |
| ananna$b | 0 | 1 | 0 | 1 |
| anna$ban | 0 | 1 | 1 | 1 |
| bananna$ | 0 | 1 | 1 | 2 |
| na$banan | 1 | 1 | 1 | 2 |
| nanna$ba | 1 | 1 | 1 | 3 |
| nna$bana | 1 | 2 | 1 | 3 |
| | 1 | 3 | 1 | 3 |

**Top-middle table** (label: n)

| | $ | a | b | n |
|---|---|---|---|---|
| $bananna | 0 | 0 | 0 | 0 |
| a$banann | 0 | 1 | 0 | 0 |
| ananna$b | 0 | 1 | 0 | 1 |
| anna$ban | 0 | 1 | 1 | 1 |
| bananna$ | 0 | 1 | 1 | 2 |
| na$banan | 1 | 1 | 1 | 2 |
| nanna$ba | 1 | 1 | 1 | 3 |
| nna$bana | 1 | 2 | 1 | 3 |
| | 1 | 3 | 1 | 3 |

(B,E) = 0, 2

**Top-right table** (label: n)

| | $ | a | b | n |
|---|---|---|---|---|
| $bananna | 0 | 0 | 0 | 0 |
| a$banann | 0 | 1 | 0 | 0 |
| ananna$b | 0 | 1 | 0 | 1 |
| anna$ban | 0 | 1 | 1 | 1 |
| bananna$ | 0 | 1 | 1 | 2 |
| na$banan | 1 | 1 | 1 | 2 |
| nanna$ba | 1 | 1 | 1 | 3 |
| nna$bana | 1 | 2 | 1 | 3 |
| | 1 | 3 | 1 | 3 |

**Bottom-left table** (label: a)

| | $ | a | b | n |
|---|---|---|---|---|
| $bananna | 0 | 0 | 0 | 0 |
| a$banann | 0 | 1 | 0 | 0 |
| ananna$b | 0 | 1 | 0 | 1 |
| anna$ban | 0 | 1 | 1 | 1 |
| bananna$ | 0 | 1 | 1 | 2 |
| na$banan | 1 | 1 | 1 | 2 |
| nanna$ba | 1 | 1 | 1 | 3 |
| nna$bana | 1 | 2 | 1 | 3 |
| | 1 | 3 | 1 | 3 |

(B,E) = 1, 2

**Bottom-middle table** (label: a)

| | $ | a | b | n |
|---|---|---|---|---|
| $bananna | 0 | 0 | 0 | 0 |
| a$banann | 0 | 1 | 0 | 0 |
| ananna$b | 0 | 1 | 0 | 1 |
| anna$ban | 0 | 1 | 1 | 1 |
| bananna$ | 0 | 1 | 1 | 2 |
| na$banan | 1 | 1 | 1 | 2 |
| nanna$ba | 1 | 1 | 1 | 3 |
| nna$bana | 1 | 2 | 1 | 3 |
| | 1 | 3 | 1 | 3 |

(B,E) = 0, 1

**Bottom-right table** (label: b)

| | $ | a | b | n |
|---|---|---|---|---|
| $bananna | 0 | 0 | 0 | 0 |
| a$banann | 0 | 1 | 0 | 0 |
| ananna$b | 0 | 1 | 0 | 1 |
| anna$ban | 0 | 1 | 1 | 1 |
| bananna$ | 0 | 1 | 1 | 2 |
| na$banan | 1 | 1 | 1 | 2 |
| nanna$ba | 1 | 1 | 1 | 3 |
| nna$bana | 1 | 2 | 1 | 3 |
| | 1 | 3 | 1 | 3 |

# BWT Searching Notes

- Don't have to store the LF mapping. A more complex algorithm (later slides) lets you compute it in O(1) time in **compressed** data on the fly with some extra storage.

- To find the range in the first column corresponding to a character:

  - Pre-compute array C[c] = # of occurrences in the string of characters lexicographically < c.

  - Then start of the "a" range, for example, is: C["a"] + 1.

- Running time: O(|pattern|)

  - Finding the range in the first column takes O(1) time using the C array.

  - Updating the range takes O(1) time using the LF mapping.

# Relationship Between BWT and Suffix Arrays

s = appellee$
123456789

**BWT matrix**

$appellee
appellee$
e$appelle
ee$appell
ellee$app
lee$appel
llee$appe
pellee$ap
ppellee$a

**The suffixes are obtained by deleting everything after the $**

$
appellee$
e$
ee$
ellee$
lee$
llee$
pellee$
ppellee$

These are still in sorted order because "$" comes before everything else

**Suffix array (start position for the suffixes)**

9
1
8
7
4
6
5
3
2

− subtract 1 →

**Suffix position - 1 = the position of the last character of the BWT matrix**

s[9-1] = e
s[1-1] = $
s[8-1] = e
s[7-1] = l
s[4-1] = p
s[6-1] = l
s[5-1] = e
s[3-1] = p
s[2-1] = a

($ is a special case)

# Relationship Between BWT and Suffix Trees

- Remember: Suffix Array = suffix numbers obtained by traversing the leaf nodes of the (ordered) Suffix Tree from left to right.

- Suffix Tree $\Rightarrow$ Suffix Array $\Rightarrow$ BWT.

Ordered suffix tree for previous example:

$\Sigma = \{\$,e,l,p\}$
s = appellee$
   123456789

# Computing BWT in O(n) time

- Easy $O(n^2 \log n)$-time algorithm to compute the BWT (create and sort the BWT matrix explicitly).

- Several direct $O(n)$-time algorithms for BWT.
  These are space efficient.

- Also can use suffix arrays or trees:

  Compute the suffix array, use correspondence between suffix array and BWT to output the BWT.

  $O(n)$-time and $O(n)$-space, but the constants are large.

# Move-To-Front Coding

*To encode a letter, use its index in the current list, and then move it to the front of the list.*

Σ            do$oodwg

*List with all letters from the allowed alphabet* →

| Σ | code |
|---|------|
| $dgow | 1 |
| d$gow | 13 |
| od$gw | 132 |
| $odgw | 1322 |
| o$dgw | 13220 |
| o$dgw | 132202 |
| do$gw | 1322024 |
| wdo$g | 13220244 = MTF(do$oodwg) |

Benefits:
- Runs of the same letter will lead to runs of 0s.
- Common letters get small numbers, while rare letters get big numbers.

# Compressing BWT Strings

Lots of possible compression schemes will benefit from preprocessing with BWT (since it tends to group runs of the same letters together).

One good scheme proposed by Ferragina & Manzini:

replace runs of 0s
with the <u>count</u> of 0s
↓

PrefixCode(rle(MTF(BWT(S))))
↑

Huffman code that
uses more bits for
rare symbols

# Pseudocode for CountingOccurrences in BWT w/o stored LF mapping

```
function Count(S_bwt, P):

    c = P[p], i = p
    sp = C[c] + 1; ep = C[c+1]

    while (sp ≤ ep) and (i ≥ 2) do
      c = P[i-1]
      sp = C[c] + Occ(c, sp-1) + 1
      ep = C[c] + Occ(c, ep)
      i = i - 1

    if ep < sp then
      return "not found"
    else
      return ep - sp + 1
```
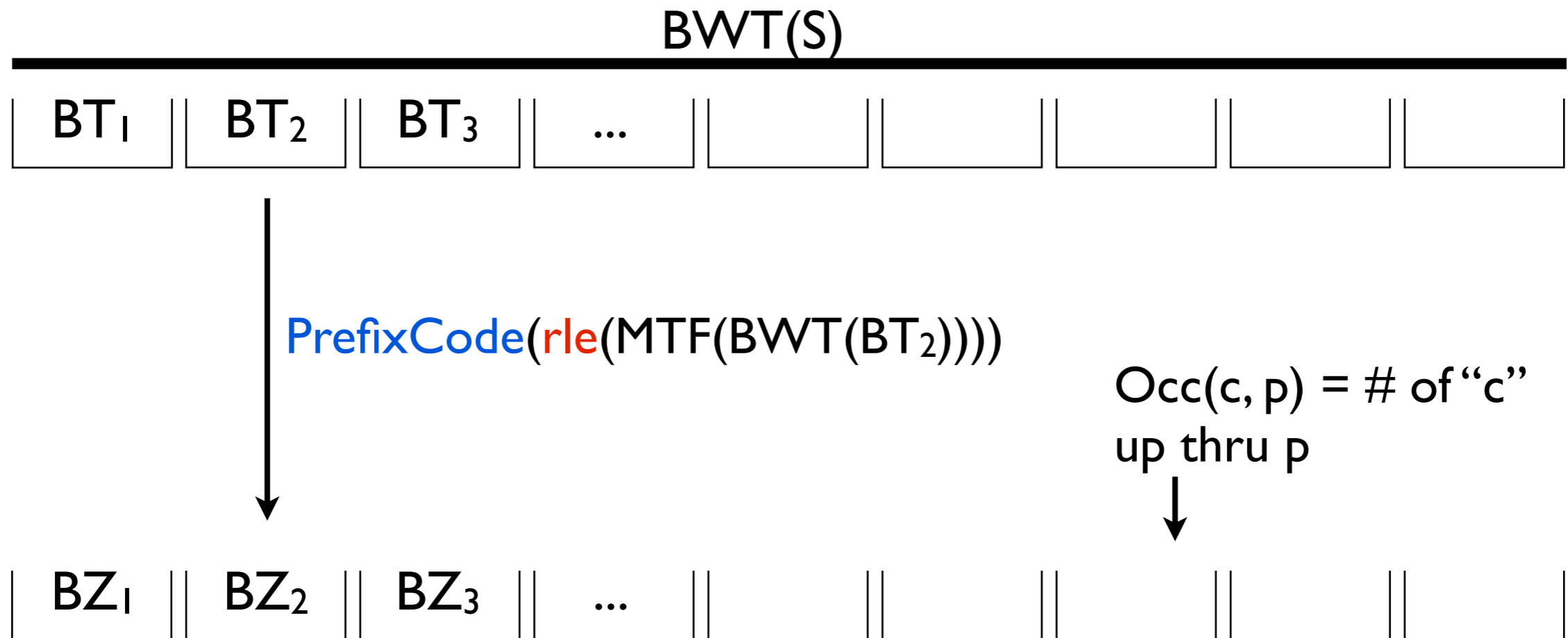
C[c] = index into first column where the "c"s begin.

$\mathbf{Occ}(c, p)$ = # of of $c$ in the first $p$ characters of BWT(S), aka the LF mapping.

# Computing Occ in Compressed String

Break BWT(S) into blocks of length L (we will decide on a value for L later):

BWT(S)

| $BT_1$ | $BT_2$ | $BT_3$ | ... | | | | | |

PrefixCode(rle(MTF(BWT($BT_2$))))

Occ(c, p) = # of "c" up thru p
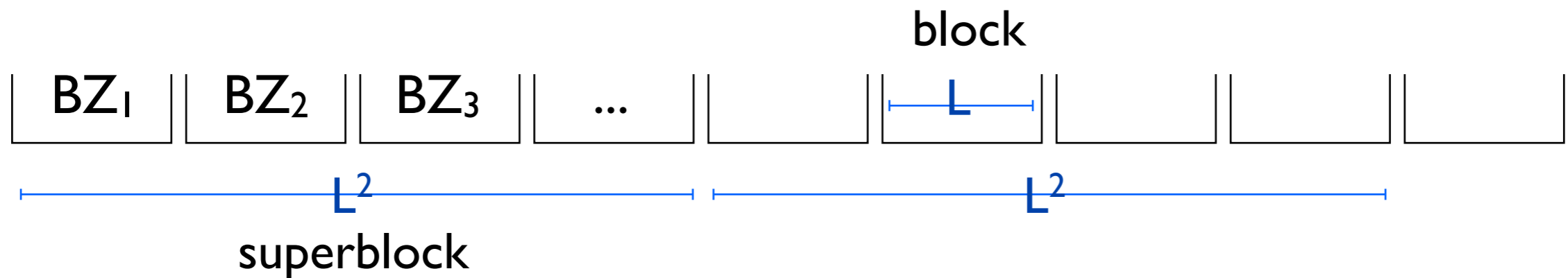
| $BZ_1$ | $BZ_2$ | $BZ_3$ | ... | | | | | |

Assumes every run of 0s is contained in a block [just for ease of explanation].

We will store some extra info for each block (and some groups of blocks) to compute Occ(c, p) quickly.

# Extra Info to Compute Occ

**block**: store $|\Sigma|$-long array giving # of occurrences of each character up thru and including this block *since the end of the last super block.*

block

| BZ₁ | BZ₂ | BZ₃ | ... | | $L$ | | | |

$L^2$        $L^2$

superblock

**superblock**: store $|\Sigma|$-long array giving # of occurrences of each character up thru *and including* this superblock
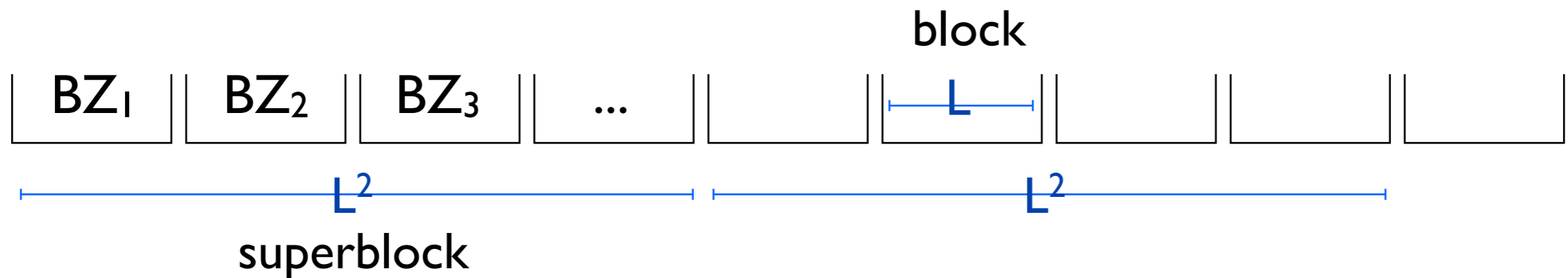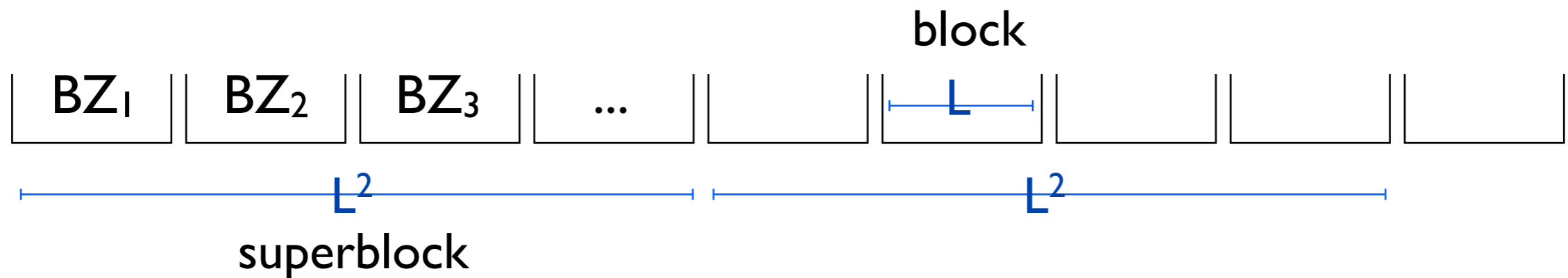
# Extra Info to Compute Occ

u = compressed length
Choose *L* = O(log *u*)

u/L blocks, each array is |Σ|log *L* long ⇒
$\frac{u}{L}\log L = \frac{u}{\log u}\log\log u$ total space.

**block**: store |Σ|-long array giving #
of occurrences of each character up
thru and including this block *since
the end of the last super block.*

block

| BZ₁ | BZ₂ | BZ₃ | ... | | | $L$ | | | |

$L^2$ ———————— $L^2$

superblock

**superblock**: store |Σ|-long
array giving # of occurrences
of each character up thru *and
including* this superblock

# Extra Info to Compute Occ

u = compressed length
Choose $L = O(\log u)$

u/L blocks, each array is $|\Sigma|\log L$ long $\Rightarrow$
$\frac{u}{L}\log L = \frac{u}{\log u}\log\log u$ total space.

**block**: store $|\Sigma|$-long array giving # of occurrences of each character up thru and including this block *since the end of the last super block.*



block

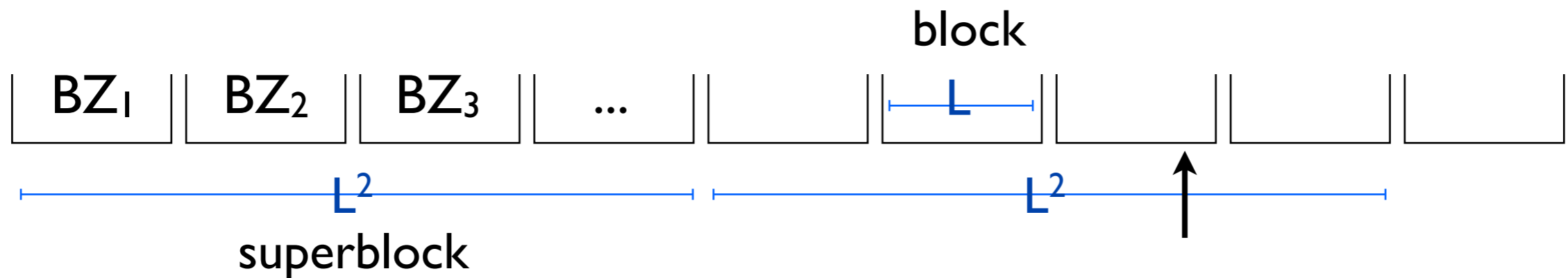| BZ₁ | BZ₂ | BZ₃ | ... | | | L | | | | |

$L^2$
superblock
$L^2$

**superblock**: store $|\Sigma|$-long array giving # of occurrences of each character up thru *and including* this superblock

$u/L^2$ superblocks, each array is $|\Sigma|\log u$ long $\Rightarrow \frac{u}{(\log u)^2}\log u = \frac{u}{\log u}$ total space.

# Extra Info to Compute Occ

$u$ = compressed length
Choose $L = O(\log u)$

block

| BZ$_1$ | BZ$_2$ | BZ$_3$ | ... | | | L | | | |

$L^2$
superblock

$L^2$

Occ(c, p) = # of "c" up thru p:

sum value at last superblock, value at end of previous block, but then need to handle *this block.*

Store an array: $M[c, k, BZ_i, MTF_i]$ = # of occurrences of c through the kth letter of a block of <u>type</u> (BZ$_i$, MTF$_i$).

Size: $O(|\Sigma|L2^L|\Sigma|) = O(L2^{L'}) = O(u^c\log u)$ for c < 1 (since the string is compressed)

# Recap

BWT useful for searching and compression.

BWT is *invertible*: given the BWT of a string, the string can be reconstructed!

BWT is computable in O(n) time.

Close relationships between Suffix Trees, Suffix Arrays, and BWT:

- Suffix array = order of the suffix numbers of the suffix tree, traversed left to right

- BWT = letters at positions given by the suffix array entries - 1

Even after compression, can search string quickly.