# Strings
## SFWRENG 2CO3: Data Structures and Algorithms

Jelle Hellings

Department of Computing and Software
McMaster University

McMaster
University

Winter 2024

## Strings over alphabets

An *alphabet* $\mathcal{A}$ is a finite set of distinct symbols.

A *string* over $\mathcal{A}$ is a sequence of symbols taken from $\mathcal{A}$.

# Strings over alphabets

An *alphabet* $\mathcal{A}$ is a finite set of distinct symbols.

A *string* over $\mathcal{A}$ is a sequence of symbols taken from $\mathcal{A}$.

### Examples

- ▶ A typical string over the roman alphabet {'a', ..., 'z', '␣'}:
  "hello", "hello␣world", "strings␣over␣alphabets", and "".

- ▶ A *bit string* is a sequence over {0, 1}:
  "0011", "1101100", "", and "0".

- ▶ A *DNA String* is a sequence over {$A$, $C$, $G$, $T$}:
  "", "AACATG", "AGT", and "AAACCCAAATTT".

- ▶ A *Unicode string* is a sequence over the unicode code points
  (149 186 symbols and counting).

- ▶ A *byte string* is a sequence over *bytes*.

# Operations on strings and alphabets

We assume the following basic operations:

- We can sequentially iterate over the symbols in a string.

- We can look up the $i$-th symbol in a string in $\Theta(1)$.
  (This can be hard in some practical settings: UTF-8 and UTF-16 strings do *not* support this).

- We assume that each alphabet $\mathcal{A}$ is an ordered list $L$ of symbols.

- For each $\sigma \in \mathcal{A}$, we can determine its position in $L$.

## BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

## BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):
1: *buckets* := $[0 \mid 0 \le i \le M-1]$.
2: **for all** $v \in L$ **do**
3:    *buckets*$[v]$ := *buckets*$[v] + 1$.
4: $k$ := 0.
5: **for all** $i := 0$ upto $M-1$ **do**
6:    **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:       $L[k]$ := $i$.
8:       $k$ := $k + 1$.

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):
1: $buckets := [0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:    $buckets[v] := buckets[v] + 1$. $\left.\begin{array}{l} \\ \\ \\ \end{array}\right\} \Theta(|L|)$
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:    **for all** $j := 0$ upto $buckets[i]$ **do**
7:       $L[k] := i$.
8:       $k := k + 1$.

## BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):

1: $buckets := [0 \mid 0 \le i \le M - 1]$.
2: **for all** $v \in L$ **do**
3:    $buckets[v] := buckets[v] + 1$. $\left.\rule{0pt}{2.5em}\right\} \Theta(|L|)$
4: $k := 0$.
5: **for all** $i := 0$ upto $M - 1$ **do**
6:    **for all** $j := 0$ upto $buckets[i]$ **do**
7:       $L[k] := i$. $\left.\rule{0pt}{3em}\right\} |L|$
8:       $k := k + 1$.

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):

1: *buckets* $:= [0 \mid 0 \le i \le M - 1]$. ⎱
2: **for all** $v \in L$ **do** ⎰ $\Theta(|L|)$
3:     *buckets*$[v] := $ *buckets*$[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M - 1$ **do** ⎱
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do** ⎰ $|L|$ ⎰ $\Theta(M + |L|)$
7:        $L[k] := i$.
8:        $k := k + 1$.

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):

1: $buckets := [0 \mid 0 \leq i \leq M - 1]$. $\left.\vphantom{\begin{array}{c}a\\b\end{array}}\right\} \Theta(M)$
2: **for all** $v \in L$ **do**
3:    $buckets[v] := buckets[v] + 1$. $\left.\vphantom{\begin{array}{c}a\\b\end{array}}\right\} \Theta(|L|)$
4: $k := 0$.
5: **for all** $i := 0$ upto $M - 1$ **do**
6:    **for all** $j := 0$ upto $buckets[i]$ **do**
7:       $L[k] := i$. $\left.\vphantom{\begin{array}{c}a\\b\\c\end{array}}\right\} |L| \left.\vphantom{\begin{array}{c}a\\b\\c\\d\\e\end{array}}\right\} \Theta(M + |L|)$
8:       $k := k + 1$.

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):

1: $buckets := [0 \mid 0 \leq i \leq M - 1]$. $\left.\rule{0pt}{1em}\right\} \Theta(M)$
2: **for all** $v \in L$ **do**
3: $\quad buckets[v] := buckets[v] + 1$. $\left.\rule{0pt}{2em}\right\} \Theta(|L|)$
4: $k := 0$.
5: **for all** $i := 0$ upto $M - 1$ **do**
6: $\quad$ **for all** $j := 0$ upto $buckets[i]$ **do**
7: $\quad\quad L[k] := i$.
8: $\quad\quad k := k + 1$.

$\left.\rule{0pt}{2em}\right\} |L|$ $\left.\rule{0pt}{3em}\right\} \Theta(M + |L|)$

$\left.\rule{0pt}{5em}\right\} \Theta(M + |L|)$

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*[$v$] := *buckets*[$v$] + 1.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*[$i$] **do**
7:        $L[k] := i$.
8:        $k := k + 1$.

$L$ = "GAGGATATGTAG".

| | |
|---|---|
| A: | 0 |
| C: | 0 |
| G: | 0 |
| T: | 0 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):
1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:    *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k$ := 0.
5: **for all** $i$ := 0 upto $M-1$ **do**
6:    **for all** $j$ := 0 upto *buckets*$[i]$ **do**
7:       $L[k]$ := $i$.
8:       $k$ := $k+1$.

$L$ = "GAGGATATGTAG".

| | |
|---|---|
| A: | 0 |
| C: | 0 |
| G: | 1 |
| T: | 0 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: $buckets := [0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:    $buckets[v] := buckets[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:    **for all** $j := 0$ upto $buckets[i]$ **do**
7:       $L[k] := i$.
8:       $k := k + 1$.

$L = $ "GAGGATATGTAG".

| | |
|---|---|
| A: | 1 |
| C: | 0 |
| G: | 1 |
| T: | 0 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \le i \le M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:         $L[k] := i$.
8:         $k := k + 1$.

$L$ = "GAGGATATGTAG".

| | |
|---|---|
| A: | 1 |
| C: | 0 |
| G: | 2 |
| T: | 0 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k$ := 0.
5: **for all** $i$ := 0 upto $M-1$ **do**
6:     **for all** $j$ := 0 upto *buckets*$[i]$ **do**
7:        $L[k]$ := $i$.
8:        $k$ := $k+1$.

$L$ = "GAGGATATGTAG".

|     |     |
| --- | --- |
| A:  | 1   |
| C:  | 0   |
| G:  | 3   |
| T:  | 0   |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):
1: *buckets* := $[0 \mid 0 \leq i \leq M - 1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k$ := 0.
5: **for all** $i$ := 0 upto $M - 1$ **do**
6:     **for all** $j$ := 0 upto *buckets*$[i]$ **do**
7:         $L[k]$ := $i$.
8:         $k$ := $k + 1$.

$L$ = "GAGGATATGTAG".

A:  | 2 |
C:  | 0 |
G:  | 3 |
T:  | 0 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:   *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k$ := 0.
5: **for all** $i$ := 0 upto $M-1$ **do**
6:   **for all** $j$ := 0 upto *buckets*$[i]$ **do**
7:     $L[k]$ := $i$.
8:     $k$ := $k + 1$.

$L$ = "GAGGATATGTAG".

|      |     |
|------|-----|
| A:   | 2   |
| C:   | 0   |
| G:   | 3   |
| T:   | 1   |

# BUCKETSORT: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BUCKETSORT($L$):

1: *buckets* $:= [0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v] :=$ *buckets*$[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:         $L[k] := i$.
8:         $k := k + 1$.

$L = $ "GAGGATATGTAG".

| | |
|---|---|
| A: | 3 |
| C: | 0 |
| G: | 3 |
| T: | 1 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:         $L[k] := i$.
8:         $k := k + 1$.

$L = $ "GAGGATATGTAG".

A: $\boxed{3}$
C: $\boxed{0}$
G: $\boxed{3}$
T: $\boxed{2}$

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):
1: *buckets* := $[0 \mid 0 \leq i \leq M - 1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k$ := 0.
5: **for all** $i$ := 0 upto $M - 1$ **do**
6:     **for all** $j$ := 0 upto *buckets*$[i]$ **do**
7:         $L[k]$ := $i$.
8:         $k$ := $k + 1$.

$L$ = "GAGGATATGTAG".

A: $\boxed{3}$
C: $\boxed{0}$
G: $\boxed{4}$
T: $\boxed{2}$

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \le i \le M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:         $L[k] := i$.
8:         $k := k+1$.

$L =$ "GAGGATATGTAG".

| | |
|---|---|
| A: | 3 |
| C: | 0 |
| G: | 4 |
| T: | 3 |

## BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: $buckets := [0 \mid 0 \le i \le M-1]$.
2: **for all** $v \in L$ **do**
3:    $buckets[v] := buckets[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:    **for all** $j := 0$ upto $buckets[i]$ **do**
7:       $L[k] := i$.
8:       $k := k + 1$.

$L$ = "GAGGATATGTAG".

| | |
|---|---|
| A: | 4 |
| C: | 0 |
| G: | 4 |
| T: | 3 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M - 1$.

**Algorithm** BucketSort($L$):

```
1:  buckets := [0 | 0 ≤ i ≤ M − 1].
2:  for all v ∈ L do
3:      buckets[v] := buckets[v] + 1.
4:  k := 0.
5:  for all i := 0 upto M − 1 do
6:      for all j := 0 upto buckets[i] do
7:          L[k] := i.
8:          k := k + 1.
```

$L$ = "GAGGATATGTAG".

|     |     |
| --- | --- |
| A:  | 4   |
| C:  | 0   |
| G:  | 5   |
| T:  | 3   |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:    *buckets*$[v]$ := *buckets*$[v] + 1$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:    **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:       $L[k] := i$.
8:       $k := k + 1$.

$L =$ "AAAAATATGTAG".

| | |
|---|---|
| A: | 4 |
| C: | 0 |
| G: | 5 |
| T: | 3 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := [0 | 0 ≤ i ≤ M − 1].
2: **for all** $v \in L$ **do**
3:  *buckets*[$v$] := *buckets*[$v$] + 1.
4: $k$ := 0.
5: **for all** $i$ := 0 upto $M − 1$ **do**
6:   **for all** $j$ := 0 upto *buckets*[$i$] **do**
7:    $L[k]$ := $i$.
8:    $k$ := $k + 1$.

$L$ = "AAAAATATGTAG".

A: | 4 |
C: | 0 |
G: | 5 |
T: | 3 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BucketSort($L$):

1: *buckets* := $[0 \mid 0 \le i \le M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:         $L[k] := i$.
8:         $k := k + 1$.

$L =$ "AAAAGGGGGTAG".

| | |
|---|---|
| A: | 4 |
| C: | 0 |
| G: | 5 |
| T: | 3 |

# BUCKETSORT: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** BUCKETSORT($L$):

1: *buckets* := $[0 \mid 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:     *buckets*$[v]$ := *buckets*$[v]$ + 1.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto *buckets*$[i]$ **do**
7:         $L[k] := i$.
8:         $k := k + 1$.

$L$ = "AAAAGGGGGTTT".

| | |
|---|---|
| A: | 4 |
| C: | 0 |
| G: | 5 |
| T: | 3 |

# BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** GBucketSort($L$, $r$):
1: *buckets* := $[[\ ]\ |\ 0 \leq i \leq M-1]$.
2: **for all** $v \in L$ **do**
3:    Append $v$ to *buckets*$[r(v)]$.
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:    **for all** $j := 0$ upto $|buckets[i]|$ **do**
7:       $L[k] := buckets[i][j]$.
8:       $k := k + 1$.

Generalization
Assume we have values that "represent" $0, \ldots, M-1$ via some function $r$.

## BucketSort: A special-purpose sort

*Assumption* We have $M$ distinct symbols with values in the range $0, \ldots, M-1$.

**Algorithm** GBucketSort($L$, $r$):
1: *buckets* := [ [ ] | $0 \leq i \leq M-1$ ].
2: **for all** $v \in L$ **do**
3:     Append $v$ to *buckets*[$r(v)$].
4: $k := 0$.
5: **for all** $i := 0$ upto $M-1$ **do**
6:     **for all** $j := 0$ upto $|$*buckets*[$i$]$|$ **do**
7:         $L[k] :=$ *buckets*[$i$][$j$].
8:         $k := k + 1$.

### Generalization
Assume we have values that "represent" $0, \ldots, M-1$ via some function $r$.

Notice that GBucketSort is *stable*.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

## RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:

# RADIXSORT: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RADIXSORT($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:     GBUCKETSORT($L, r_d$) with $r_d(S) = S[d]$.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:    Stable-sort $L$ on the $d$-th string symbols.
3:    GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.   $\Bigr\}\ \Theta(|L| + |\mathcal{A}|)$

# RADIXSORT: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RADIXSORT($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:    Stable-sort $L$ on the $d$-th string symbols. $\left.\rule{0pt}{2.2em}\right\}\Theta(|L| + |\mathcal{A}|)$
3:    GBUCKETSORT($L, r_d$) with $r_d(S) = S[d]$.

$\left.\rule{0pt}{3.2em}\right\}\Theta(k(|L| + |\mathcal{A}|))$

# RADIXSORT: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RADIXSORT($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:     GBUCKETSORT($L, r_d$) with $r_d(S) = S[d]$.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):

1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBucketSort($L, r_d$) with $r_d(S) = S[d]$.

$L = [$"AGCTCT",
"ATTAAC",
"GCGCGG",
"GGCGCG",
"TCTATG",
"TCACCG",
"AGCTGA",
"ATCTAA",
"GTCTGC",
"TGGACG"$]$

## RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

$L = [$"AGCTCT",     $L = [$"AGCTGA",
      "ATTAAC",           "ATCTAA",
      "GCGCGG",           "GTCTGC",
      "GGCGCG",           "ATTAAC",
      "TCTATG",  $\rightarrow$  "GCGCGG",
      "TCACCG",           "GGCGCG",
      "AGCTGA",           "TCTATG",
      "ATCTAA",           "TCACCG",
      "GTCTGC",           "TGGACG",
      "TGGACG"]         "AGCTCT"]

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto $0$ **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

$L = [$"AGCTGA",
"ATCTAA",
"GTCTGC",
"ATTAAC",
"GCGCGG",
"GGCGCG",
"TCTATG",
"TCACCG",
"TGGACG",
"AGCTCT"]

$\rightarrow$

$L = [$"ATCTAA",
"ATTAAC",
"GGCGCG",
"TCACCG",
"TGGACG",
"AGCTCT",
"AGCTGA",
"GTCTGC",
"GCGCGG",
"TCTATG"]

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto $0$ **do**
2:    Stable-sort $L$ on the $d$-th string symbols.
3:    GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

$L = [$"ATCTAA",
   "ATTAAC",
   "GGCGCG",
   "TCACCG",
   "TGGACG",
   "AGCTCT",
   "AGCTGA",
   "GTCTGC",
   "GCGCGG",
   "TCTATG"]

$\rightarrow$

$L = [$"ATTAAC",
   "TGGACG",
   "TCTATG",
   "TCACCG",
   "GCGCGG",
   "GGCGCG",
   "ATCTAA",
   "AGCTCT",
   "AGCTGA",
   "GTCTGC"]

## RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto $0$ **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:     GBucketSort($L, r_d$) with $r_d(S) = S[d]$.

$$L = [\text{"ATTAAC"},$$
$$\text{"TGGACG"},$$
$$\text{"TCTATG"},$$
$$\text{"TCACCG"},$$
$$\text{"GCGCGG"},$$
$$\text{"GGCGCG"},$$
$$\text{"ATCTAA"},$$
$$\text{"AGCTCT"},$$
$$\text{"AGCTGA"},$$
$$\text{"GTCTGC"}]$$

$\rightarrow$

$$L = [\text{"TCACCG"},$$
$$\text{"GGCGCG"},$$
$$\text{"ATCTAA"},$$
$$\text{"AGCTCT"},$$
$$\text{"AGCTGA"},$$
$$\text{"GTCTGC"},$$
$$\text{"TGGACG"},$$
$$\text{"GCGCGG"},$$
$$\text{"ATTAAC"},$$
$$\text{"TCTATG"}]$$

## RADIXSORT: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RADIXSORT($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBUCKETSORT($L$, $r_d$) with $r_d(S) = S[d]$.

$$L = [\text{"TCACCG"},$$
$$\text{"GGCGCG"},$$
$$\text{"ATCTAA"},$$
$$\text{"AGCTCT"},$$
$$\text{"AGCTGA"},$$
$$\text{"GTCTGC"},$$
$$\text{"TGGACG"},$$
$$\text{"GCGCGG"},$$
$$\text{"ATTAAC"},$$
$$\text{"TCTATG"}]$$

$\rightarrow$

$$L = [\text{"TCACCG"},$$
$$\text{"GCGCGG"},$$
$$\text{"TCTATG"},$$
$$\text{"GGCGCG"},$$
$$\text{"AGCTCT"},$$
$$\text{"AGCTGA"},$$
$$\text{"TGGACG"},$$
$$\text{"ATCTAA"},$$
$$\text{"GTCTGC"},$$
$$\text{"ATTAAC"}]$$

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

$$
L = [\text{"TCACCG"}, \\
\text{"GCGCGG"}, \\
\text{"TCTATG"}, \\
\text{"GGCGCG"}, \\
\text{"AGCTCT"}, \\
\text{"AGCTGA"}, \\
\text{"TGGACG"}, \\
\text{"ATCTAA"}, \\
\text{"GTCTGC"}, \\
\text{"ATTAAC"}]
$$

$\rightarrow$

$$
L = [\text{"AGCTCT"}, \\
\text{"AGCTGA"}, \\
\text{"ATCTAA"}, \\
\text{"ATTAAC"}, \\
\text{"GCGCGG"}, \\
\text{"GGCGCG"}, \\
\text{"GTCTGC"}, \\
\text{"TCACCG"}, \\
\text{"TCTATG"}, \\
\text{"TGGACG"}]
$$

## RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:     GBucketSort($L, r_d$) with $r_d(S) = S[d]$.

$L = [$"AGCTCT",
         "AGCTGA",
         "ATCTAA",
         "ATTAAC",
         "GCGCGG",
         "GGCGCG",
         "GTCTGC",
         "TCACCG",
         "TCTATG",
         "TGGACG"$]$.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:    Stable-sort $L$ on the $d$-th string symbols.
3:    GBucketSort($L, r_d$) with $r_d(S) = S[d]$.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBucketSort($L, r_d$) with $r_d(S) = S[d]$.

Correctness

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:     GBucketSort($L, r_d$) with $r_d(S) = S[d]$.

Correctness
*Invariant: In L, the suffix of the last $k - (d + 1)$ symbols is sorted.*

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:    Stable-sort $L$ on the $d$-th string symbols.
3:    GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

Correctness
*Invariant: In L, the suffix of the last $k - (d + 1)$ symbols is sorted.*

Generalization: strings with variable lengths up-to-$k$

# RADIXSORT: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RADIXSORT($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:    Stable-sort $L$ on the $d$-th string symbols.
3:    GBUCKETSORT($L, r_d$) with $r_d(S) = S[d]$.

Correctness
*Invariant: In L, the suffix of the last $k - (d + 1)$ symbols is sorted.*

Generalization: strings with variable lengths up-to-$k$
Let $S$ be a string of length $|S| < k$.
Interpret $S[|S|], \ldots, S[k - 1]$ as symbols that come before all other symbols.

The book calls this *least-significant-digit string sort*.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:     Stable-sort $L$ on the $d$-th string symbols.
3:     GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

Is RadixSort worth it?

# RADIXSORT: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RADIXSORT($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBUCKETSORT($L, r_d$) with $r_d(S) = S[d]$.

Is RADIXSORT worth it?
- ▶ Optimal sorts perform $\Theta(|L| \log(|L|))$ *comparisons*.
- ▶ Comparing two strings of length $k$ costs at-most $\Theta(k)$.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:    Stable-sort $L$ on the $d$-th string symbols.
3:    GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

Is RadixSort worth it?

▶ Optimal sorts perform $\Theta(|L| \log(|L|))$ *comparisons*.

▶ Comparing two strings of length $k$ costs at-most $\Theta(k)$.

▶ For $|L|$ *random strings*, comparisons are expected to cost $\Theta(\log_2(|L|))$.

# RadixSort: A special-purpose sort for strings

*Assumption.* We have strings of length $k$ over alphabet $\mathcal{A}$.

**Algorithm** RadixSort($L$):
1: **for** $d := k - 1$ downto 0 **do**
2:   Stable-sort $L$ on the $d$-th string symbols.
3:   GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$.

Is RadixSort worth it?
- ▶ Optimal sorts perform $\Theta(|L| \log(|L|))$ *comparisons*.
- ▶ Comparing two strings of length $k$ costs at-most $\Theta(k)$.
- ▶ For $|L|$ *random strings*, comparisons are expected to cost $\Theta(\log_2(|L|))$.

$\Theta(k(|L| + |\mathcal{A}|))$ versus $\Theta(k|L| \log(|L|))$ (or $\Theta(|L| \log^2(|L|))$ expected).

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

Consider GBucketSort($L$, $r_0$):

$$
\begin{array}{ll}
L = [\text{“ATTAAC”}, & L = [\text{“ATTAAC”}, \\
\quad\text{“GCGCGG”}, & \quad\text{“AGCTGA”}, \\
\quad\text{“GGCGCG”}, & \quad\text{“ATCTAA”}, \\
\quad\text{“TCTATG”}, & \quad\text{“GCGCGG”}, \\
\quad\text{“TCACCG”}, \;\rightarrow & \quad\text{“GGCGCG”}, \\
\quad\text{“AGCTGA”}, & \quad\text{“GTCTGC”}, \\
\quad\text{“ATCTAA”}, & \quad\text{“TCTATG”}, \\
\quad\text{“GTCTGC”}, & \quad\text{“TCACCG”}, \\
\quad\text{“TGGACG”}] & \quad\text{“TGGACG”}].
\end{array}
$$

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

Consider GBucketSort($L$, $r_0$):

$$
\begin{array}{lcl}
L = [\text{"ATTAAC"}, & & L = [\text{"ATTAAC"}, \\
\quad \text{"GCGCGG"}, & & \quad \text{"AGCTGA"}, \\
\quad \text{"GGCGCG"}, & & \quad \text{"ATCTAA"}, \\
\quad \text{"TCTATG"}, & & \quad \text{"GCGCGG"}, \\
\quad \text{"TCACCG"}, & \rightarrow & \quad \text{"GGCGCG"}, \\
\quad \text{"AGCTGA"}, & & \quad \text{"GTCTGC"}, \\
\quad \text{"ATCTAA"}, & & \quad \text{"TCTATG"}, \\
\quad \text{"GTCTGC"}, & & \quad \text{"TCACCG"}, \\
\quad \text{"TGGACG"}] & & \quad \text{"TGGACG"}].
\end{array}
$$

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

Consider GBucketSort($L$, $r_0$):

$$L = [\text{"ATTAAC"}, \qquad L = [\text{"ATTAAC"},$$

| | |
|---|---|
| "GCGCGG", | "AGCTGA", |
| "GGCGCG", | "ATCTAA", |
| "TCTATG", | "GCGCGG", |
| "TCACCG", $\rightarrow$ | "GGCGCG", |
| "AGCTGA", | "GTCTGC", |
| "ATCTAA", | "TCTATG", |
| "GTCTGC", | "TCACCG", |
| "TGGACG"] | "TGGACG"]. |

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

Consider GBucketSort($L$, $r_0$).

After *constructing* the buckets with respect to $r_0$,
we only need to resort the individual buckets.

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

Consider GBucketSort($L$, $r_0$).

After *constructing* the buckets with respect to $r_0$,
we only need to resort the individual buckets.

We will apply this idea *recursively*!

# Most-significant-digit string sort

RadixSort does not try to minimize the number of sorting rounds:
if $k$ is the length of the longest string in $L$, then RadixSort "reorders" the list $k$ times.

Consider GBucketSort($L$, $r_0$).

After *constructing* the buckets with respect to $r_0$,
we only need to resort the individual buckets.

We will apply this idea *recursively*!

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2:   GBucketSort($L$, $r_d$) with $r_d(S) = S[d]$,
     during which we further sort each
     individual bucket separately.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     $buckets := [[\ ] \mid 0 \le i \le |\mathcal{A}| - 1]$.
3:     **for all** $v \in L$ **do**
4:         Append $v$ to $buckets[v[d]]$.
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:         $k_{\text{start}} := k$.
8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:             $L[k] := buckets[i][j]$.
10:            $k := k + 1$.
11:        MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:   *buckets* := $[[\ ] \mid 0 \leq i \leq |\mathcal{A}| - 1]$.
3:   **for all** $v \in L$ **do**
4:     Append $v$ to *buckets*$[v[d]]$.
5:   $k := 0$.
6:   **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:     $k_{\text{start}} := k$.
8:     **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:       $L[k] := buckets[i][j]$.
10:       $k := k + 1$.
11:     MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"ATTAAC",
    "GCGCGG",
    "GGCGCG",
    "TCTATG",
    "TCACCG",
    "AGCTGA",
    "ATCTAA",
    "GTCTGC",
    "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:      $buckets := [[\ ] \mid 0 \leq i \leq |\mathcal{A}| - 1]$.
3:      **for all** $v \in L$ **do**
4:          Append $v$ to $buckets[v[d]]$.
5:      $k := 0$.
6:      **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:          $k_{\text{start}} := k$.
8:          **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:              $L[k] := buckets[i][j]$.
10:             $k := k + 1$.
11:        MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"ATTAAC",
      "AGCTGA",
      "ATCTAA",
      "TCTATG",
      "TCACCG",
      "AGCTGA",
      "ATCTAA",
      "GTCTGC",
      "TGGACG"].

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L, d$):

1:  **if** $d < k$ **and** $|L| > 1$ **then**
2:    $buckets := [\,[\,]\ |\ 0 \leq i \leq |\mathcal{A}| - 1\,]$.
3:    **for all** $v \in L$ **do**
4:      Append $v$ to $buckets[v[d]]$.
5:    $k := 0$.
6:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:      $k_{\text{start}} := k$.
8:      **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:        $L[k] := buckets[i][j]$.
10:        $k := k + 1$.
11:      MSD-Sort($L[k_{\text{start}} \ldots k], d + 1$).

$L = [$"ATTAAC",
"AGCTGA",
"ATCTAA",
"TCTATG",
"TCACCG",
"AGCTGA",
"ATCTAA",
"GTCTGC",
"TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     $buckets := [[\ ]\ |\ 0 \le i \le |\mathcal{A}| - 1]$.
3:     **for all** $v \in L$ **do**
4:         Append $v$ to $buckets[v[d]]$.
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:         $k_{\text{start}} := k$.
8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:             $L[k] := buckets[i][j]$.
10:            $k := k + 1$.
11:         MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
     "AGCTGA",
     "ATCTAA",
     "TCTATG",
     "TCACCG",
     "AGCTGA",
     "ATCTAA",
     "GTCTGC",
     "TGGACG"]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:    $buckets := [[\ ]\ |\ 0 \le i \le |\mathcal{A}| - 1]$.
3:    **for all** $v \in L$ **do**
4:       Append $v$ to $buckets[v[d]]$.
5:    $k := 0$.
6:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:       $k_{\text{start}} := k$.
8:       **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:          $L[k] := buckets[i][j]$.
10:          $k := k + 1$.
11:       MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
      "AGCTGA",
      "ATCTAA",
      "TCTATG",
      "TCACCG",
      "AGCTGA",
      "ATCTAA",
      "GTCTGC",
      "TGGACG"].

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     *buckets* := $[[\ ] \mid 0 \leq i \leq |\mathcal{A}| - 1]$.
3:     **for all** $v \in L$ **do**
4:        Append $v$ to *buckets*$[v[d]]$.
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:        $k_{\text{start}} := k$.
8:        **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:           $L[k] := buckets[i][j]$.
10:          $k := k + 1$.
11:      MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
      "AGCTGA",
      "ATCTAA",
      "TCTATG",
      "TCACCG",
      "AGCTGA",
      "ATCTAA",
      "GTCTGC",
      "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2: *buckets* := $[\,[\;] \mid 0 \le i \le |\mathcal{A}| - 1\,]$.
3: **for all** $v \in L$ **do**
4:  Append $v$ to *buckets*$[v[d]]$.
5: $k$ := 0.
6: **for all** $i$ := 0 upto $|\mathcal{A}| - 1$ **do**
7:  $k_{\text{start}}$ := $k$.
8:  **for all** $j$ := 0 upto $|\textit{buckets}[i]|$ **do**
9:   $L[k]$ := *buckets*$[i][j]$.
10:   $k$ := $k + 1$.
11:  MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
   "ATTAAC",
   "ATCTAA",
   "TCTATG",
   "TCACCG",
   "AGCTGA",
   "ATCTAA",
   "GTCTGC",
   "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:   $buckets := [[\ ]\ |\ 0 \le i \le |\mathcal{A}| - 1]$.
3:   **for all** $v \in L$ **do**
4:     Append $v$ to $buckets[v[d]]$.
5:   $k := 0$.
6:   **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:     $k_{\text{start}} := k$.
8:     **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:       $L[k] := buckets[i][j]$.
10:       $k := k + 1$.
11:     MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
    "ATTAAC",
    "ATCTAA",
    "TCTATG",
    "TCACCG",
    "AGCTGA",
    "ATCTAA",
    "GTCTGC",
    "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:    $buckets := [[\ ]\ |\ 0 \le i \le |\mathcal{A}| - 1]$.
3:    **for all** $v \in L$ **do**
4:       Append $v$ to $buckets[v[d]]$.
5:    $k := 0$.
6:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:       $k_{\text{start}} := k$.
8:       **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:          $L[k] := buckets[i][j]$.
10:         $k := k + 1$.
11:      MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
     "ATCTAA",
     "ATCTAA",
     "TCTATG",
     "TCACCG",
     "AGCTGA",
     "ATCTAA",
     "GTCTGC",
     "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     $buckets := [[\,] \mid 0 \leq i \leq |\mathcal{A}| - 1]$.
3:     **for all** $v \in L$ **do**
4:         Append $v$ to $buckets[v[d]]$.
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:         $k_{\text{start}} := k$.
8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:             $L[k] := buckets[i][j]$.
10:            $k := k + 1$.
11:         MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
      "ATCTAA",
      "ATCTAA",
      "TCTATG",
      "TCACCG",
      "AGCTGA",
      "ATCTAA",
      "GTCTGC",
      "TGGACG"].

# Most-significant-digit string sort

**Algorithm** MSD-SORT($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2:    *buckets* := $[[ \, ] \mid 0 \le i \le |\mathcal{A}| - 1]$.
3:    **for all** $v \in L$ **do**
4:       Append $v$ to *buckets*$[v[d]]$.
5:    $k := 0$.
6:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:       $k_{\text{start}} := k$.
8:       **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:          $L[k] := buckets[i][j]$.
10:          $k := k + 1$.
11:       MSD-SORT($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
    "ATCTAA",
    "ATCTAA",
    "TCTATG",
    "TCACCG",
    "AGCTGA",
    "ATCTAA",
    "GTCTGC",
    "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:   $buckets := [[\,] \mid 0 \le i \le |\mathcal{A}| - 1]$.
3:   **for all** $v \in L$ **do**
4:     Append $v$ to $buckets[v[d]]$.
5:   $k := 0$.
6:   **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:     $k_{\text{start}} := k$.
8:     **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:       $L[k] := buckets[i][j]$.
10:       $k := k + 1$.
11:     MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
    "ATCTAA",
    "ATTAAC",
    "TCTATG",
    "TCACCG",
    "AGCTGA",
    "ATCTAA",
    "GTCTGC",
    "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:    *buckets* := $[ [ \ ] \mid 0 \leq i \leq |\mathcal{A}| - 1 ]$.
3:    **for all** $v \in L$ **do**
4:      Append $v$ to *buckets*$[v[d]]$.
5:    $k := 0$.
6:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:      $k_{\text{start}} := k$.
8:      **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:        $L[k] := buckets[i][j]$.
10:       $k := k + 1$.
11:      MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$ "AGCTGA",
     "ATCTAA",
     "ATTAAC",
     "TCTATG",
     "TCACCG",
     "AGCTGA",
     "ATCTAA",
     "GTCTGC",
     "TGGACG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort(L, d):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     *buckets* := [ [ ] | $0 \leq i \leq |\mathcal{A}| - 1$].
3:     **for all** $v \in L$ **do**
4:         Append $v$ to *buckets*[$v[d]$].
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:         $k_{\text{start}} := k$.
8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:             $L[k] := buckets[i][j]$.
10:            $k := k + 1$.
11:         MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = $ ["AGCTGA",
     "ATCTAA",
     "ATTAAC",
     "TCTATG",
     "TCACCG",
     "AGCTGA",
     "ATCTAA",
     "GTCTGC",
     "TGGACG"].

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

  1: **if** $d < k$ **and** $|L| > 1$ **then**

  2:     $buckets := [[\ ] \mid 0 \le i \le |\mathcal{A}| - 1]$.

  3:     **for all** $v \in L$ **do**

  4:         Append $v$ to $buckets[v[d]]$.

  5:     $k := 0$.

  6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**

  7:         $k_{\text{start}} := k$.

  8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**

  9:             $L[k] := buckets[i][j]$.

10:             $k := k + 1$.

11:         MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",

      "ATCTAA",

      "ATTAAC",

      "GCGCGG",

      "GGCGCG",

      "GTCTGC",

      "ATCTAA",

      "GTCTGC",

      "TGGACG"$]$.

## Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2:   $buckets := [[\ ] \mid 0 \leq i \leq |\mathcal{A}| - 1]$.
3:   **for all** $v \in L$ **do**
4:     Append $v$ to $buckets[v[d]]$.
5:   $k := 0$.
6:   **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:     $k_{\text{start}} := k$.
8:     **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:       $L[k] := buckets[i][j]$.
10:       $k := k + 1$.
11:     MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
      "ATCTAA",
      "ATTAAC",
      "GCGCGG",
      "GGCGCG",
      "GTCTGC",
      "ATCTAA",
      "GTCTGC",
      "TGGACG"].

## Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):

1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     *buckets* := $[[\ ]\ |\ 0 \le i \le |\mathcal{A}| - 1]$.
3:     **for all** $v \in L$ **do**
4:         Append $v$ to *buckets*$[v[d]]$.
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:         $k_{\text{start}} := k$.
8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:             $L[k] := buckets[i][j]$.
10:            $k := k + 1$.
11:         MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
      "ATCTAA",
      "ATTAAC",
      "GCGCGG",
      "GGCGCG",
      "GTCTGC",
      "TGGACG",
      "TCTATG",
      "TCACCG"$]$.

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2:    $buckets := [[\ ]\ |\ 0 \leq i \leq |\mathcal{A}| - 1]$.
3:    **for all** $v \in L$ **do**
4:       Append $v$ to $buckets[v[d]]$.
5:    $k := 0$.
6:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:       $k_{\text{start}} := k$.
8:       **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:          $L[k] := buckets[i][j]$.
10:         $k := k + 1$.
11:      MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

$L = [$"AGCTGA",
    "ATCTAA",
    "ATTAAC",
    "GCGCGG",
    "GGCGCG",
    "GTCTGC",
    "TCACCG",
    "TCTATG",
    "TGGACG"].

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $d < k$ **and** $|L| > 1$ **then**
2:     $buckets := [\, [\,] \mid 0 \le i \le |\mathcal{A}| - 1\,]$.
3:     **for all** $v \in L$ **do**
4:         Append $v$ to $buckets[v[d]]$.
5:     $k := 0$.
6:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
7:         $k_{\text{start}} := k$.
8:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
9:             $L[k] := buckets[i][j]$.
10:           $k := k + 1$.
11:         MSD-Sort($L[k_{\text{start}} \ldots k)$, $d + 1$).

## Finetuning

We end up with many arrays *buckets* that each hold $|\mathcal{A}|$ lists!

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $|L| \leq |\mathcal{A}|$ **then**
2:    Use another algorithm to sort $L$.
3: **else if** $d < k$ **and** $|L| > 1$ **then**
4:    $buckets := [\,[\,]\,\mid 0 \leq i \leq |\mathcal{A}| - 1\,]$.
5:    **for all** $v \in L$ **do**
6:       Append $v$ to $buckets[v[d]]$.
7:    $k := 0$.
8:    **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
9:       $k_{\text{start}} := k$.
10:      **for all** $j := 0$ upto $|buckets[i]|$ **do**
11:         $L[k] := buckets[i][j]$.
12:         $k := k + 1$.
13:      MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

Finetuning
We end up with many arrays *buckets*
that each hold $|\mathcal{A}|$ lists!

# Most-significant-digit string sort

**Algorithm** MSD-Sort($L$, $d$):
1: **if** $|L| \leq |\mathcal{A}|$ **then**
2:     Use another algorithm to sort $L$.
3: **else if** $d < k$ **and** $|L| > 1$ **then**
4:     $buckets := [\,[\,]\mid 0 \leq i \leq |\mathcal{A}| - 1\,]$.
5:     **for all** $v \in L$ **do**
6:         Append $v$ to $buckets[v[d]]$.
7:     $k := 0$.
8:     **for all** $i := 0$ upto $|\mathcal{A}| - 1$ **do**
9:         $k_{\text{start}} := k$.
10:         **for all** $j := 0$ upto $|buckets[i]|$ **do**
11:             $L[k] := buckets[i][j]$.
12:             $k := k + 1$.
13:         MSD-Sort($L[k_{\text{start}} \ldots k]$, $d + 1$).

## Finetuning
We end up with many arrays *buckets*
that each hold $|\mathcal{A}|$ lists!

## Complexity
At-most $\Theta(k(|L| + |\mathcal{A}|))$.

# Sorting: Best practices

So which sort algorithm is the best?

# Sorting: Best practices

So which sort algorithm is the best?
Depends on the to-be-sorted input.

# Sorting: Best practices

So which sort algorithm is the best?
Depends on the to-be-sorted input.

Often, your standard sort algorithm will be sufficient.

# Tries: special-purpose sets and dictionaries

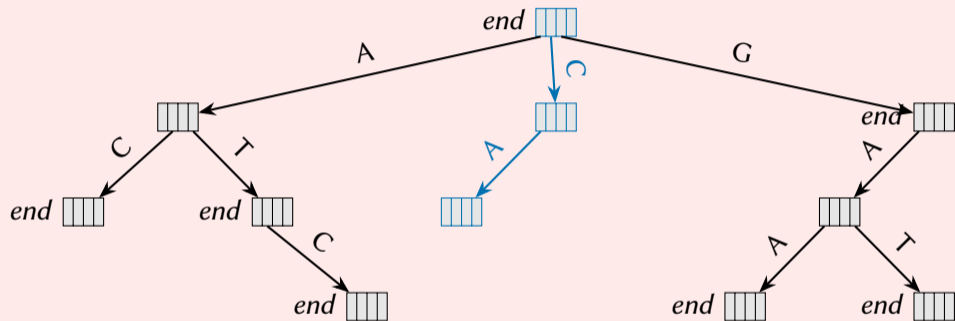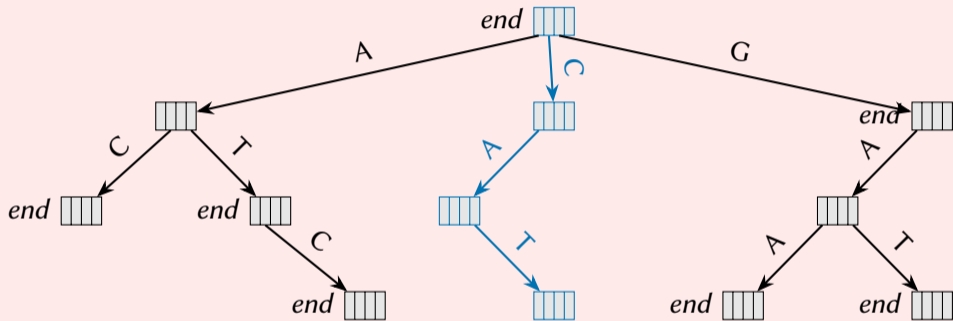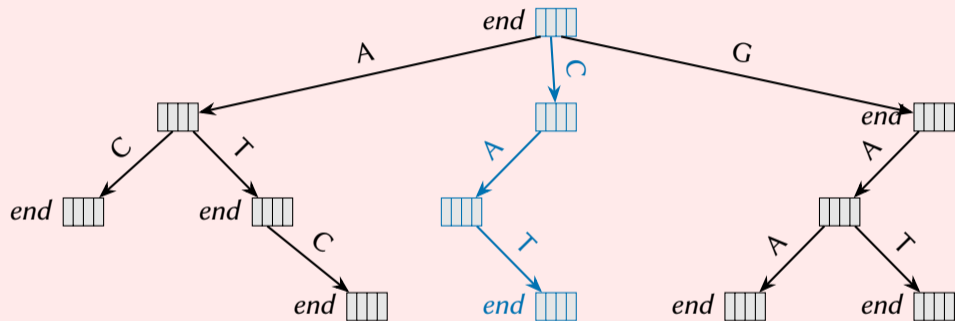*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.

A *Trie* is a set representation that can hold strings over $\mathcal{A}$ such that:
- strings of length $N$ can be *added* in $\Theta(N)$;
- strings of length $N$ can be *removed* in $\Theta(N)$;
- *checking* whether a string of length $N$ is in the set costs $\Theta(N)$;
- one can efficiently *iterate* over all strings in the set (in sorted order).

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.

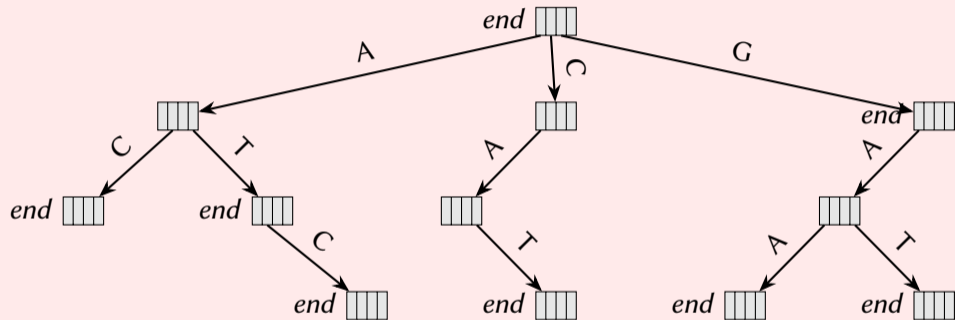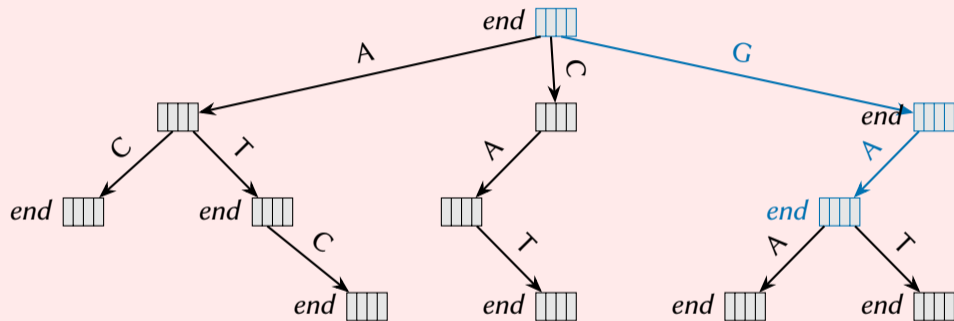A *Trie* is a set representation that can hold strings over $\mathcal{A}$ such that:
- strings of length $N$ can be *added* in $\Theta(N)$;
- strings of length $N$ can be *removed* in $\Theta(N)$;
- *checking* whether a string of length $N$ is in the set costs $\Theta(N)$;
- one can efficiently *iterate* over all strings in the set (in sorted order).

We have seen tries with $\mathcal{A} = \{0, 1\} \rightarrow$ BSSet in Example Assignment 3.

# Tries: special-purpose sets and dictionaries

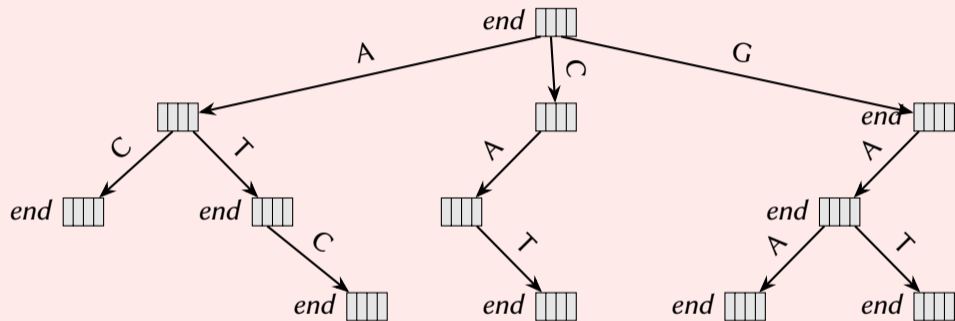*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.

Each node *n* in a *trie* $T$ over $\mathcal{A} = \{\sigma_1, \ldots, \sigma_M\}$ has
- a flag *n.end* that is true if the node *n* represents a string in $T$;
- at-most $M$ edges to children labeled $\sigma_1, \ldots, \sigma_M$.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.

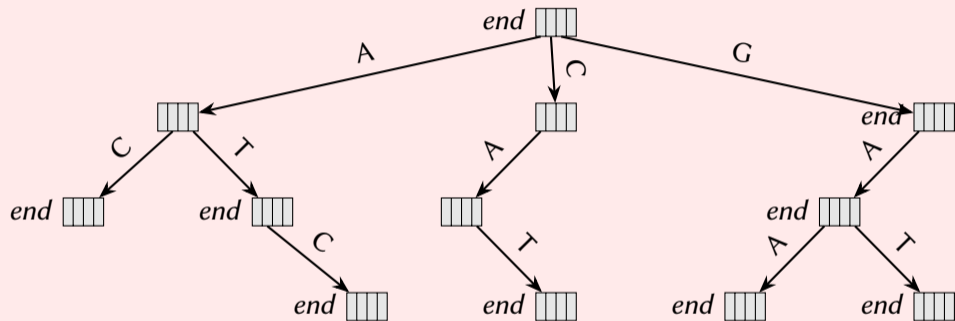Each node *n* in a *trie* $T$ over $\mathcal{A} = \{\sigma_1, \ldots, \sigma_M\}$ has
- a flag *n.end* that is true if the node *n* represents a string in $T$;
- at-most $M$ edges to children labeled $\sigma_1, \ldots, \sigma_M$.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.

# Tries: special-purpose sets and dictionaries

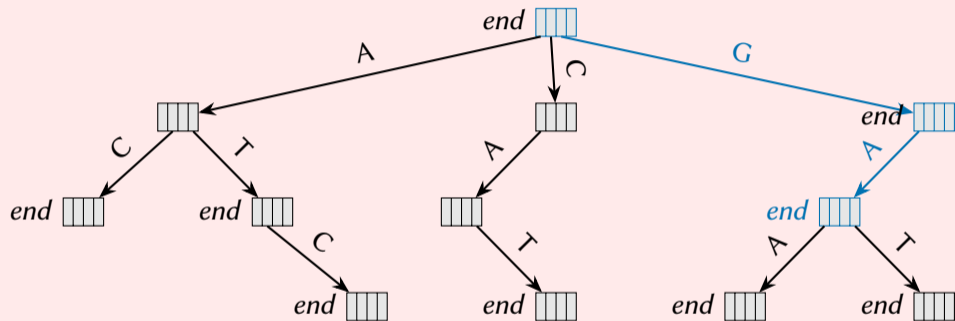*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



"", "AC"

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



"", "AC", "AT"

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



"", "AC", "AT", "ATC"
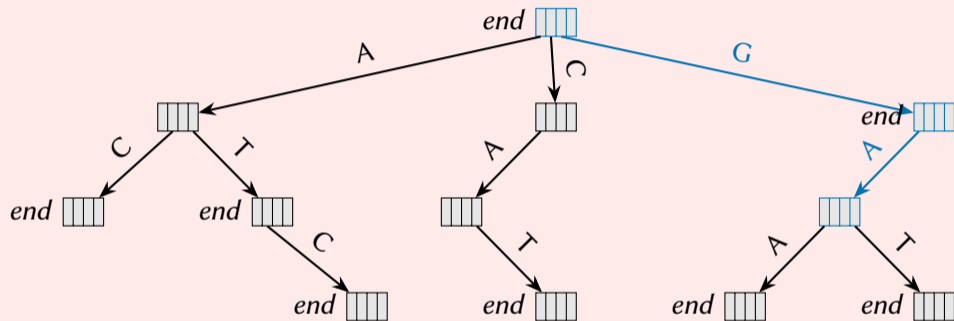
# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



"", "AC", "AT", "ATC", "G"

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



"", "AC", "AT", "ATC", "G", "GAA"

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



"", "AC", "AT", "ATC", "G", "GAA", "GAT".

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Adding a string
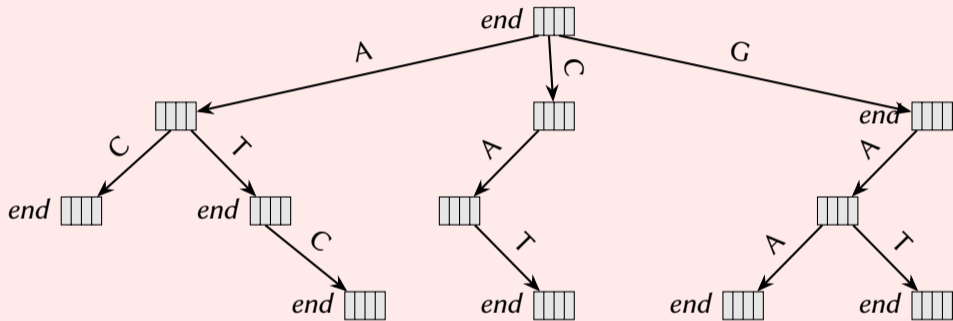
▶ Follow-or-make a path according to the string symbols.

▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



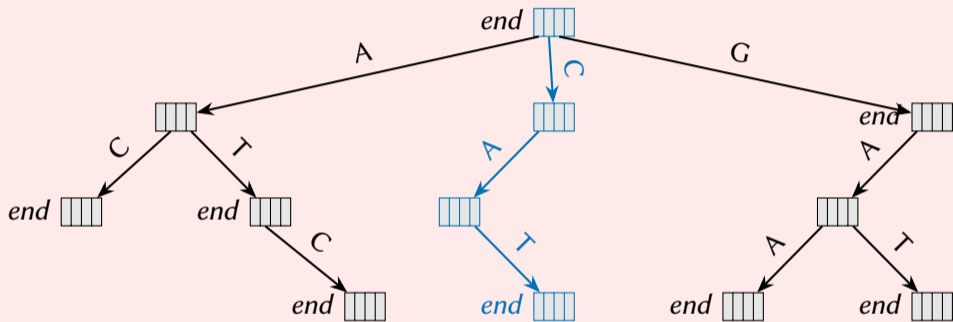## Adding a string "CAT"

- ▶ Follow-or-make a path according to the string symbols.
- ▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Adding a string "CAT"

- ▶ Follow-or-make a path according to the string symbols.
- ▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



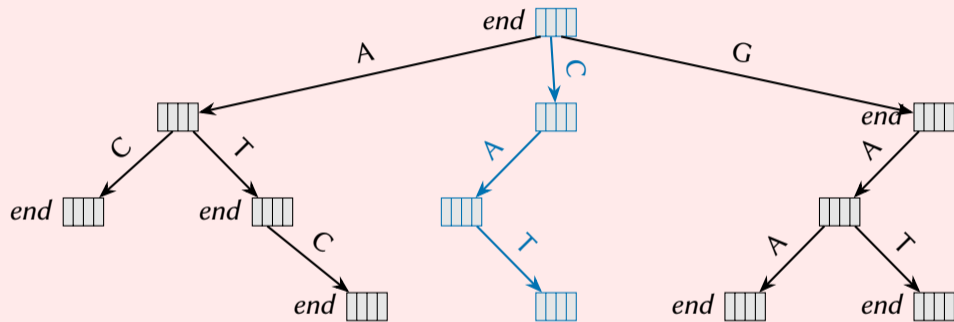## Adding a string "CAT"

▶ Follow-or-make a path according to the string symbols.

▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Adding a string "CAT"

- ▶ Follow-or-make a path according to the string symbols.
- ▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



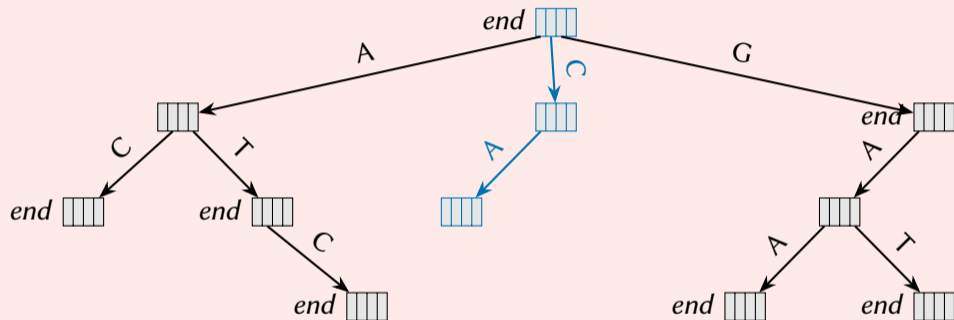Adding a string "CAT"

- ▶ Follow-or-make a path according to the string symbols.
- ▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



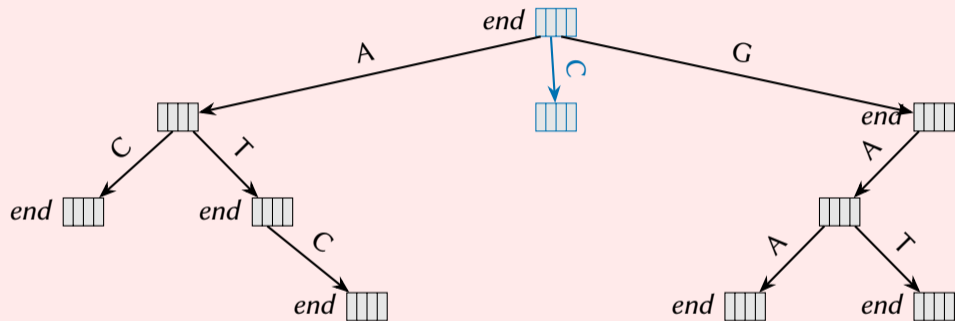## Adding a string "GA"

- ▶ Follow-or-make a path according to the string symbols.
- ▶ Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Adding a string "GA"

► Follow-or-make a path according to the string symbols.

► Set *n.end* on the last node *n* on this path.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



### Removing a string

- ▶ Follow a path according to the string symbols to node *n* and unset *n.end*.
- ▶ Remove *n* if *n* has no children.
- ▶ Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

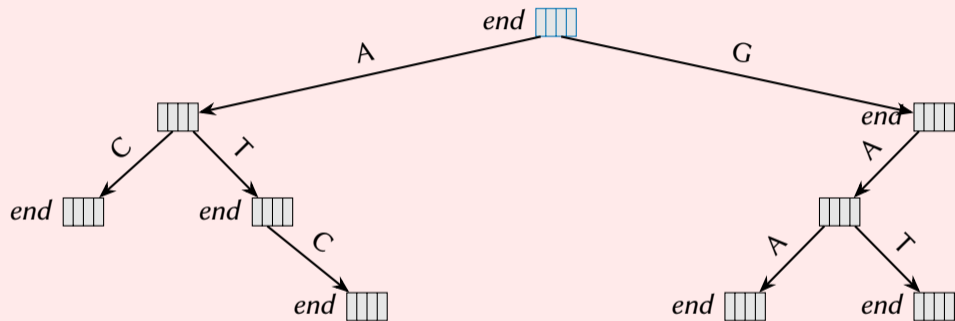*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "GA"

▶ Follow a path according to the string symbols to node *n* and unset *n.end*.

▶ Remove *n* if *n* has no children.

▶ Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

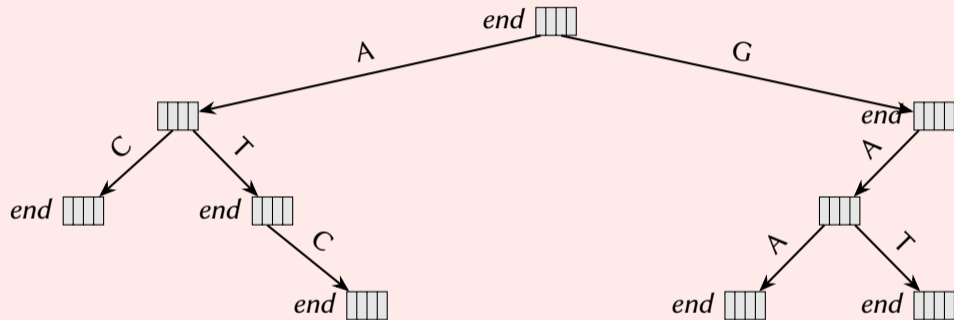*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "GA"

- Follow a path according to the string symbols to node *n* and unset *n.end*.
- Remove *n* if *n* has no children.
- Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

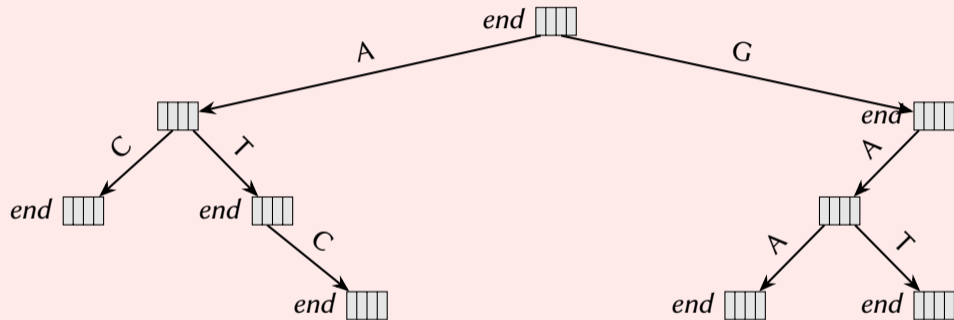*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "GA"

- ▶ Follow a path according to the string symbols to node *n* and unset *n.end*.
- ▶ Remove *n* if *n* has no children.
- ▶ Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

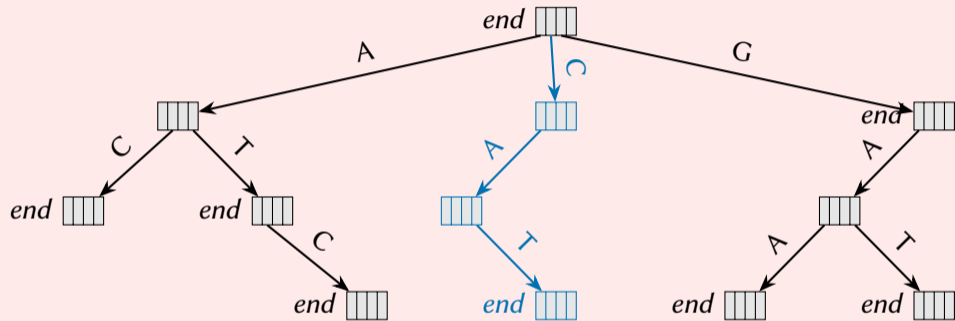*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "CAT"

- Follow a path according to the string symbols to node *n* and unset *n.end*.
- Remove *n* if *n* has no children.
- Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "CAT"
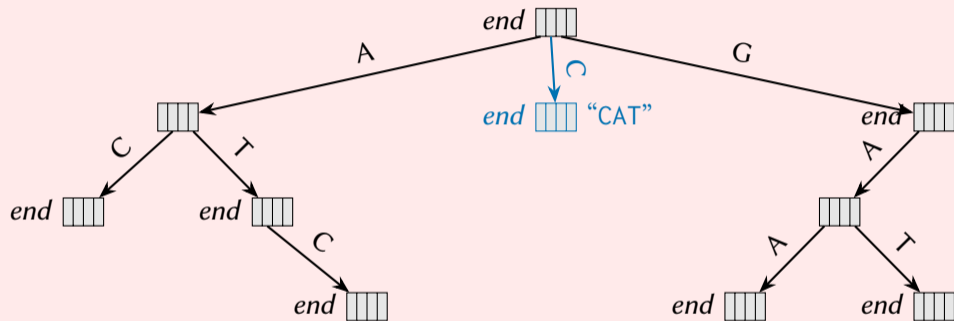
▶ Follow a path according to the string symbols to node *n* and unset *n.end*.

▶ Remove *n* if *n* has no children.

▶ Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "CAT"
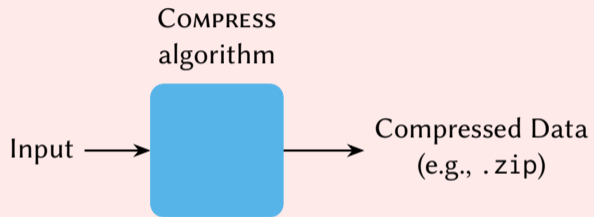
- Follow a path according to the string symbols to node *n* and unset *n.end*.
- Remove *n* if *n* has no children.
- Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "CAT"

- ▶ Follow a path according to the string symbols to node *n* and unset *n.end*.
- ▶ Remove *n* if *n* has no children.
- ▶ Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Removing a string "CAT"

- ▶ Follow a path according to the string symbols to node *n* and unset *n.end*.
- ▶ Remove *n* if *n* has no children.
- ▶ Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



### Removing a string "CAT"

- Follow a path according to the string symbols to node *n* and unset *n.end*.
- Remove *n* if *n* has no children.
- Recurse to the ancestors *m*: remove *m* if *m* has no children and *m.end* is unset.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Printing all strings in-order

Perform a pre-order traversal starting at the root. For each node *n*:

▶ print the path from root to node *n* if *n.end* is set;

▶ pre-order traverse all children in-order of alphabet symbols.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Printing all strings in-order with prefix $W$

▶ Follow a path according to the string symbols of $W$ to node $m$.

▶ Perform a pre-order traversal starting at the node $m$.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Finetuning

▶ To deal with big alphabets:
use a dictionary with $\mathcal{A}$-symbols as keys at each node to store all edges.

# Tries: special-purpose sets and dictionaries

*Assumption* We have an alphabet $\mathcal{A}$ with $M = |\mathcal{A}|$ symbols.



## Finetuning

- To deal with big alphabets:
  use a dictionary with $\mathcal{A}$-symbols as keys at each node to store all edges.

- To *compress* non-branching paths: nodes can represent strings of symbols.

# Data compression

Input

# Data compression



Compress algorithm

Input ⟶ Compressed Data (e.g., .zip)

# Data compression

# Data compression



*Lossless compression*: The **input** must be equivalent to the **output**!

# Limits of compression

### Theorem
*No algorithm A can compress every input I.*

# Limits of compression

### Theorem
*No algorithm A can compress every input I into compressed data $A(I)$ with $|A(I)| < |I|$.*

# Limits of compression

### Theorem

*No algorithm A can compress every input I into compressed data $A(I)$ with $|A(I)| < |I|$.*

### Proof

# Limits of compression

### Theorem

*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

### Proof

Intuituin 1  If A can compress every *input*, then $|A(A(I))| < |A(I)|$.

# Limits of compression

### Theorem

*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

### Proof

Intuituin 1  If A can compress every *input*, then $|A(A(I))| < |A(I)|$.
Repeaded application of A on the *input* will eventually lead to zero bits!

# Limits of compression

## Theorem

*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

## Proof

Intuituin 1    If A can compress every *input*, then $|A(A(I))| < |A(I)|$.

                   Repeaded application of A on the *input* will eventually lead to zero bits!

Intuition 2    Consider all possible *inputs* of $N$ bits: we have $2^N$ distinct values.

# Limits of compression

### Theorem

*No algorithm A can compress every input I into compressed data $A(I)$ with $|A(I)| < |I|$.*

### Proof

Intuituin 1     If A can compress every *input*, then $|A(A(I))| < |A(I)|$.

Repeaded application of A on the *input* will eventually lead to zero bits!

Intuition 2     Consider all possible *inputs* of $N$ bits: we have $2^N$ distinct values.

We need at-least $2^N$ distinct *outputs* of A on these inputs.

# Limits of compression

### Theorem
*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

### Proof

Intuituin 1   If A can compress every *input*, then $|A(A(I))| < |A(I)|$.
Repeaded application of A on the *input* will eventually lead to zero bits!

Intuition 2   Consider all possible *inputs* of $N$ bits: we have $2^N$ distinct values.
We need at-least $2^N$ distinct *outputs* of A on these inputs.

   ▶ We can compress at-most $2^M$ values to a size of $M < N$ bits.

# Limits of compression

## Theorem

*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

## Proof

Intuituin 1   If A can compress every *input*, then $|A(A(I))| < |A(I)|$.
Repeaded application of A on the *input* will eventually lead to zero bits!

Intuition 2   Consider all possible *inputs* of $N$ bits: we have $2^N$ distinct values.
We need at-least $2^N$ distinct *outputs* of A on these inputs.

- ▶ We can compress at-most $2^M$ values to a size of $M < N$ bits.
- ▶ We can compress *a small fraction* $\frac{2^M}{2^N} = 2^{M-N}$ of all inputs to $M$ bits.
  E.g., we can compress at-most $2^{32}$ out of $2^{64}$ values from 8 B to 4 B.

# Limits of compression

## Theorem

*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

## Proof

Intuituin 1   If A can compress every *input*, then $|A(A(I))| < |A(I)|$.
                 Repeaded application of A on the *input* will eventually lead to zero bits!

Intuition 2   Consider all possible *inputs* of $N$ bits: we have $2^N$ distinct values.
                We need at-least $2^N$ distinct *outputs* of A on these inputs.

- ▶ We can compress at-most $2^M$ values to a size of $M < N$ bits.
- ▶ We can compress *a small fraction* $\frac{2^M}{2^N} = 2^{M-N}$ of all inputs to $M$ bits.
  E.g., we can compress at-most $2^{32}$ out of $2^{64}$ values from 8 B to 4 B.
  Doing so will require more-than 8 B for some other inputs.

# Limits of compression

## Theorem
*No algorithm A can compress every input I into compressed data A(I) with $|A(I)| < |I|$.*

## Proof

Intuituin 1   If A can compress every *input*, then $|A(A(I))| < |A(I)|$.
Repeaded application of A on the *input* will eventually lead to zero bits!

Intuition 2   Consider all possible *inputs* of $N$ bits: we have $2^N$ distinct values.
We need at-least $2^N$ distinct *outputs* of A on these inputs.

- We can compress at-most $2^M$ values to a size of $M < N$ bits.
- We can compress *a small fraction* $\frac{2^M}{2^N} = 2^{M-N}$ of all inputs to $M$ bits.
  E.g., we can compress at-most $2^{32}$ out of $2^{64}$ values from 8 B to 4 B.
  Doing so will require more-than 8 B for some other inputs.

Conceptually: We need *structure* in the input to be able to reliably compress that input!

## A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\text{B} = 8N\text{bit}$.

How many bytes do we need to represent $\mathcal{S}$?

# A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\text{B} = 8N\text{bit}$.

## How many bytes do we need to represent $\mathcal{S}$?

▶ We have $|\mathcal{A}| = 4$ distinct values.

# A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\mathrm{B} = 8N\mathrm{bit}$.

## How many bytes do we need to represent $\mathcal{S}$?

▶ We have $|\mathcal{A}| = 4$ distinct values.

▶ We can represent 4 distinct values with 2 bit: "00", "01", "10", and "11".

# A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\text{B} = 8N\text{bit}$.

## How many bytes do we need to represent $\mathcal{S}$?

▶ We have $|\mathcal{A}| = 4$ distinct values.

▶ We can represent 4 distinct values with 2 bit: "00", "01", "10", and "11".

▶ A single byte holds 8 bit.

# A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\mathrm{B} = 8N\mathrm{bit}$.

## How many bytes do we need to represent $\mathcal{S}$?

- We have $|\mathcal{A}| = 4$ distinct values.
- We can represent 4 distinct values with 2 bit: "00", "01", "10", and "11".
- A single byte holds 8 bit.

Hence, we need at-most $\frac{8N}{4} = 2N\mathrm{bit}$.
A compression ratio of 0.25.

# A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\mathrm{B} = 8N\,\mathrm{bit}$.

## How many bytes do we need to represent $\mathcal{S}$?

- We have $|\mathcal{A}| = 4$ distinct values.
- We can represent 4 distinct values with 2 bit: "00", "01", "10", and "11".
- A single byte holds 8 bit.

Hence, we need at-most $\frac{8N}{4} = 2N\,\mathrm{bit}$.
A compression ratio of 0.25.

## From *bits* to *bytes*

We can store *four DNA characters* in one byte. Can we store $\mathcal{S}$ in $\frac{2N}{8}\mathrm{B}$ using the above?

# A simple structure: Small alphabets

Consider DNA strings over the alphabet $\mathcal{A} = \{A, C, T, G\}$.

An usual DNA string $\mathcal{S}$ represented by $N$ characters takes up $N\text{B} = 8N\text{bit}$.

## How many bytes do we need to represent $\mathcal{S}$?

- ▶ We have $|\mathcal{A}| = 4$ distinct values.
- ▶ We can represent 4 distinct values with 2 bit: "00", "01", "10", and "11".
- ▶ A single byte holds 8 bit.

Hence, we need at-most $\frac{8N}{4} = 2N\text{bit}$.
A compression ratio of 0.25.

## From *bits* to *bytes*

We can store *four DNA characters* in one byte. Can we store $\mathcal{S}$ in $\frac{2N}{8}\text{B}$ using the above?

*No!* Where in the last byte would our string end?

E.g., "ACTGA" takes 10 bit (1.25 B).

# A common structure: Repetition

Consider the following string of bits:

0000000000000001111111000000001111111111000000000011

# A common structure: Repetition

Consider the following string of bits:

000000000000000111111100000000011111111110000000000011

15 zeros   7 ones   8 zeros   11 ones   10 zeros   2 ones

# A common structure: Repetition

Consider the following string of bits:

$$\underbrace{000000000000000}_{\text{15 zeros}}\underbrace{1111111}_{\text{7 ones}}\underbrace{00000000}_{\text{8 zeros}}\underbrace{11111111111}_{\text{11 ones}}\underbrace{0000000000}_{\text{10 zeros}}\underbrace{11}_{\text{2 ones}}$$

| Number | (in 4-bit binary) |
|--------|-------------------|
| 15     | 1111              |
| 7      | 0111              |
| 8      | 1000              |
| 11     | 1011              |
| 10     | 1010              |
| 2      | 0010              |

# A common structure: Repetition

Consider the following string of bits:

$$\underbrace{1111}_{0s}\underbrace{0111}_{1s}\underbrace{1000}_{0s}\underbrace{101}_{1s}\underbrace{110}_{0s}\underbrace{100010}_{1s}$$

| Number | (in 4-bit binary) |
|--------|-------------------|
| 15 | 1111 |
| 7 | 0111 |
| 8 | 1000 |
| 11 | 1011 |
| 10 | 1010 |
| 2 | 0010 |

From $15 + 7 + 8 + 11 + 10 + 2 = 53$ bit to $6 \cdot 4 = 24$ bit.

# A common structure: Repetition

Consider the following string of bits:

00000000000000000111111100000000111111111100000000001111

| | 17 zeros | 7 ones | 8 zeros | 11 ones | 10 zeros | 2 ones |

| Number | (in 4-bit binary) |
|--------|-------------------|
| 15 | 1111 |
| 7 | 0111 |
| 8 | 1000 |
| 11 | 1011 |
| 10 | 1010 |
| 2 | 0010 |

# A common structure: Repetition

Consider the following string of bits:

```
1111000000100111100010111010001 0
```

0s  1s  0s  1s  0s  1s  0s  1s

| Number | (in 4-bit binary) |
|--------|-------------------|
| 15     | 1111              |
| 7      | 0111              |
| 8      | 1000              |
| 11     | 1011              |
| 10     | 1010              |
| 2      | 0010              |

From $17 + 7 + 8 + 11 + 10 + 2 = 55$ bit to $8 \cdot 4 = 32$ bit.

# A common structure: Repetition

Consider the following string of bits:

1111000000100111100010111010001

0s   1s   0s   1s   0s   1s   0s   1s

| Number | (in 4-bit binary) |
|--------|-------------------|
| 15     | 1111              |
| 7      | 0111              |
| 8      | 1000              |
| 11     | 1011              |
| 10     | 1010              |
| 2      | 0010              |

Run-length encoding: *simple idea* with good results on *bitmaps*.

# Another common structure: Using symbol frequencies

Consider *simple* text written in the English language.

## Another common structure: Using symbol frequencies

Consider *simple* text written in the English language.
The text uses the following 66 symbols "frequently":

- ▶ Digits 0123456789: 10 symbols.
- ▶ Lower-case letters "a"–"z": 26 symbols.
- ▶ upper-case letters "A"–"Z": 26 symbols.
- ▶ Punctuation " ", ".", ",", "!": 4 symbols.

# Another common structure: Using symbol frequencies

Consider *simple* text written in the English language.
The text uses the following 66 symbols "frequently":

- ▶ Digits 0123456789: 10 symbols.
- ▶ Lower-case letters "a"–"z": 26 symbols.
- ▶ upper-case letters "A"–"Z": 26 symbols.
- ▶ Punctuation " ", ".", ",", "!": 4 symbols.

Stored normally, *each* symbol occupies $1\,B = 8\,bit$.

# Another common structure: Using symbol frequencies

Consider *simple* text written in the English language.
The text uses the following 66 symbols "frequently":

- ▶ Digits 0123456789: 10 symbols.
- ▶ Lower-case letters "a"–"z": 26 symbols.
- ▶ upper-case letters "A"–"Z": 26 symbols.
- ▶ Punctuation " ", ".", ",", "!": 4 symbols.

Stored normally, *each* symbol occupies 1 B = 8 bit.

Even in these "frequent" symbols, some are much rarer than others: "x" versus "e".

# Another common structure: Using symbol frequencies

Consider *simple* text written in the English language.
The text uses the following 66 symbols "frequently":

- ▶ Digits 0123456789: 10 symbols.
- ▶ Lower-case letters "a"–"z": 26 symbols.
- ▶ upper-case letters "A"–"Z": 26 symbols.
- ▶ Punctuation " ", ".", ",", "!": 4 symbols.

Stored normally, *each* symbol occupies 1 B = 8 bit.

Even in these "frequent" symbols, some are much rarer than others: "x" versus "e".

*Idea.* Use fewer bits for frequent characters, more for rare characters.

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

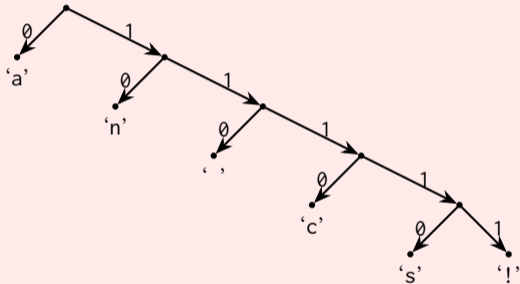| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     |             |
| 'n'    | 5     |             |
| ' '    | 4     |             |
| 'c'    | 3     |             |
| 's'    | 1     |             |
| '!'    | 1     |             |

The string has *6 distinct symbols*: at-least 3 bits if *all the same length*.

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

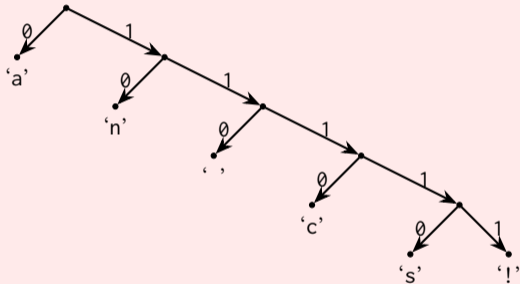| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a' | 6 | 000 |
| 'n' | 5 | 001 |
| ' ' | 4 | 010 |
| 'c' | 3 | 011 |
| 's' | 1 | 100 |
| '!' | 1 | 101 |

The string has *6 distinct symbols*: at-least 3 bits if *all the same length*.

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

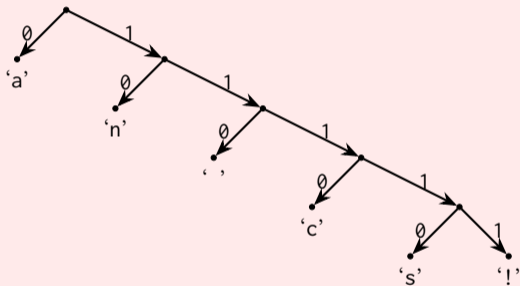| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a' | 6 | 0 |
| 'n' | 5 | 1 |
| ' ' | 4 | 00 |
| 'c' | 3 | 01 |
| 's' | 1 | 10 |
| '!' | 1 | 11 |

*Attempt 1.* The most-frequent symbols get the shortest bit patterns.

```
a n n a   c a n   s c a n   a   c a n !
0 1 1 0 00 01 0 1 00 10 01 0 1 00 0 00 01 0 1 11
```

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a' | 6 | 0 |
| 'n' | 5 | 1 |
| ' ' | 4 | 00 |
| 'c' | 3 | 01 |
| 's' | 1 | 10 |
| '!' | 1 | 11 |

*Attempt 1.* The most-frequent symbols get the shortest bit patterns.

```
a n n a   c a n   s c a n   a   c a n !
0 1 1 0 00 01 0 1 00 10 01 0 1 00 0 00 01 0 1 11
     01100001010010010100000010111
```

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 1           |
| ' '    | 4     | 00          |
| 'c'    | 3     | 01          |
| 's'    | 1     | 10          |
| '!'    | 1     | 11          |

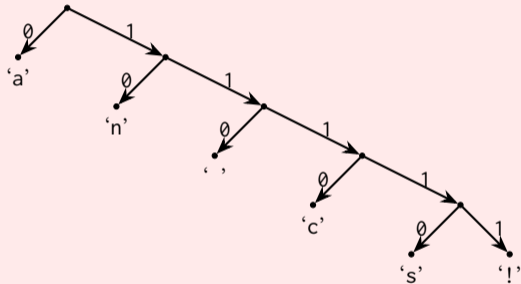*Attempt 1.* The most-frequent symbols get the shortest bit patterns.

01100001010010010100000010111 ← 29 bit *instead of at-least* 60 bit.

## Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 1           |
| ' '    | 4     | 00          |
| 'c'    | 3     | 01          |
| 's'    | 1     | 10          |
| '!'    | 1     | 11          |

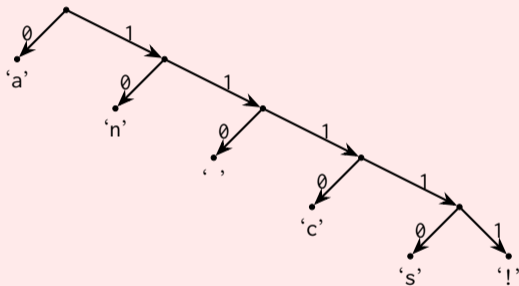*Attempt 1.* The most-frequent symbols get the shortest bit patterns.

01100001010010010100000010111 ← 29 bit *instead of at-least* 60 bit.
Y
a?

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 1           |
| ' '    | 4     | 00          |
| 'c'    | 3     | 01          |
| 's'    | 1     | 10          |
| '!'    | 1     | 11          |

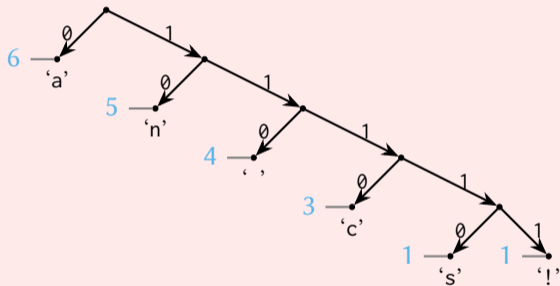*Attempt 1.* The most-frequent symbols get the shortest bit patterns.

01100001010010010100000010111 ← 29 bit *instead of at-least* 60 bit.
⌣
c?

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a' | 6 | 0 |
| 'n' | 5 | 1 |
| ' ' | 4 | 00 |
| 'c' | 3 | 01 |
| 's' | 1 | 10 |
| '!' | 1 | 11 |



*Attempt 1.* The most-frequent symbols get the shortest bit patterns.

$$011000010100100101000000010111 \leftarrow 29\,\text{bit } \textit{instead of at-least } 60\,\text{bit.}$$

*Issue.* The bit pattern of one symbol (e.g., a) is a *prefix* of other symbols!

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

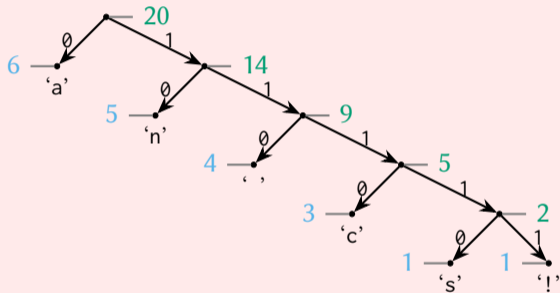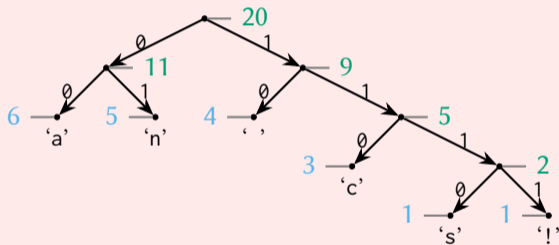| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a' | 6 | 0 |
| 'n' | 5 | 10 |
| ' ' | 4 | 110 |
| 'c' | 3 | 1110 |
| 's' | 1 | 11110 |
| '!' | 1 | 11111 |



*Attempt 2.* The most-frequent symbols get the shortest *prefix-free* bit patterns.

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 10          |
| ' '    | 4     | 110         |
| 'c'    | 3     | 1110        |
| 's'    | 1     | 11110       |
| '!'    | 1     | 11111       |



*Attempt 2.* The most-frequent symbols get the shortest *prefix-free* bit patterns.

```
a  n  n  a     c    a  n     s     c    a  n     a     c    a  n  !
0 10 10 0 110 1110 0 10 110 11110 1110 0 10 110 0 110 1110 0 10 11111
```

## Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 10          |
| ' '    | 4     | 110         |
| 'c'    | 3     | 1110        |
| 's'    | 1     | 11110       |
| '!'    | 1     | 11111       |



*Attempt 2.* The most-frequent symbols get the shortest *prefix-free* bit patterns.

```
a  n  n  a     c   a  n     s     c   a  n     a     c   a  n    !
0  10 10 0  110 1110 0  10 110 11110 1110 0  10 110 0  110 1110 0  10 11111
   0101001101110010110111101110010110011011100101111  ← 50 bit.
```

13/20

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 10          |
| ' '    | 4     | 110         |
| 'c'    | 3     | 1110        |
| 's'    | 1     | 11110       |
| '!'    | 1     | 11111       |



*Attempt 2.* The most-frequent symbols get the shortest *prefix-free* bit patterns.

```
0101001101110010110111101110010110011011100101111 ← 50 bit.
Y
a
```

## Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a' | 6 | 0 |
| 'n' | 5 | 10 |
| ' ' | 4 | 110 |
| 'c' | 3 | 1110 |
| 's' | 1 | 11110 |
| '!' | 1 | 11111 |



*Attempt 2.* The most-frequent symbols get the shortest *prefix-free* bit patterns.

```
0101001101110010110111101110010110011011100101111 ← 50 bit.
 a n n a ' '  c a n ' '  s  c a n ' ' a ' '  c a n  '!
```

## Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 10          |
| ' '    | 4     | 110         |
| 'c'    | 3     | 1110        |
| 's'    | 1     | 11110       |
| '!'    | 1     | 11111       |



*Attempt 2.* The most-frequent symbols get the shortest *prefix-free* bit patterns.

0101001101110010110111101110010110011011100101111 ← 50 bit.

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 10          |
| ' '    | 4     | 110         |
| 'c'    | 3     | 1110        |
| 's'    | 1     | 11110       |
| '!'    | 1     | 11111       |



*Questions.* How to construct the bit patterns and are these patterns optimal?

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 0           |
| 'n'    | 5     | 10          |
| ' '    | 4     | 110         |
| 'c'    | 3     | 1110        |
| 's'    | 1     | 11110       |
| '!'    | 1     | 11111       |



*Questions.* How to construct the bit patterns and are these patterns optimal?

$$6 \cdot 1 + 5 \cdot 2 + 4 \cdot 3 + 3 \cdot 4 + 1 \cdot 5 + 1 \cdot 5 = 50.$$

# Another common structure: Using symbol frequencies

Consider the string "anna can scan a can!".

| Symbol | Count | Bit pattern |
|--------|-------|-------------|
| 'a'    | 6     | 00          |
| 'n'    | 5     | 01          |
| ' '    | 4     | 10          |
| 'c'    | 3     | 110         |
| 's'    | 1     | 1110        |
| '!'    | 1     | 1111        |



*Questions.* How to construct the bit patterns and are these patterns optimal?

$$6 \cdot 2 + 5 \cdot 2 + 4 \cdot 2 + 3 \cdot 3 + 1 \cdot 4 + 1 \cdot 4 = 47.$$

# Huffman coding

### Problem
Given an alphabet $\mathcal{A}$ and symbol-frequencies $f : \mathcal{A} \rightarrow [0, 1]$,
Produce *prefix-free bit patterns* for all symbols in $\mathcal{A}$ such that these patterns are optimal.

# Huffman coding

### Problem
Given an alphabet $\mathcal{A}$ and symbol-frequencies $f : \mathcal{A} \rightarrow [0, 1]$,
Produce *prefix-free bit patterns* for all symbols in $\mathcal{A}$ such that these patterns are optimal.

*Optimal* No other bit patterns will compress strings $\mathcal{S}$ over $\mathcal{A}$ more
       *Assuming* symbols counts in $\mathcal{S}$ agree with $f$.

# Huffman coding

## Problem

Given an alphabet $\mathcal{A}$ and symbol-frequencies $f : \mathcal{A} \to [0, 1]$,
Produce *prefix-free bit patterns* for all symbols in $\mathcal{A}$ such that these patterns are optimal.

*Optimal* No other bit patterns will compress strings $\mathcal{S}$ over $\mathcal{A}$ more
   *Assuming* symbols counts in $\mathcal{S}$ agree with $f$.

## Algorithm HuffmanPFTrie($f$):

1: $Q :=$ a min-priority queue.
2: **for all** $\sigma \in \mathcal{A}$ **do**
3:    Make a leaf-node $n$ labeled $\sigma$.
4:    Add $(n, f(\sigma))$ to $Q$ with priority $f(\sigma)$.
5: **while** $|Q| \geq 2$ **do**
6:    $(n_0, p_0) :=$ DelMin$(Q)$, $(n_1, p_1) :=$ DelMin$(Q)$.
7:    Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:    Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9: **return** $n$ with $(n, p) :=$ DelMin$(Q)$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |

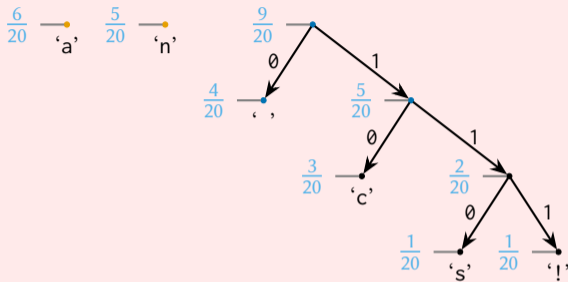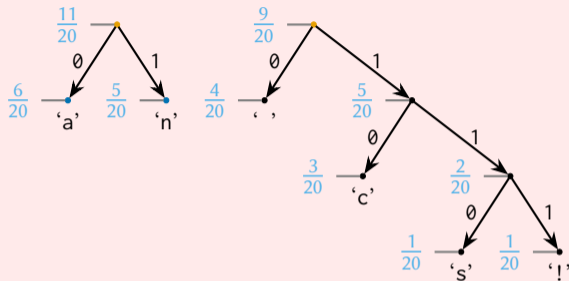$\frac{6}{20}$ •'a'  $\frac{5}{20}$ •'n'  $\frac{4}{20}$ •' ',  $\frac{3}{20}$ •'c'  $\frac{1}{20}$ •'s'  $\frac{1}{20}$ •'!'

**Algorithm** HUFFMANPFTRIE($f$):

1: $Q$ := a min-priority queue.
2: **for all** $\sigma \in \mathcal{A}$ **do**
3:    Make a leaf-node $n$ labeled $\sigma$.
4:    Add $(n, f(\sigma))$ to $Q$ with priority $f(\sigma)$.
5: ...

# Huffman coding

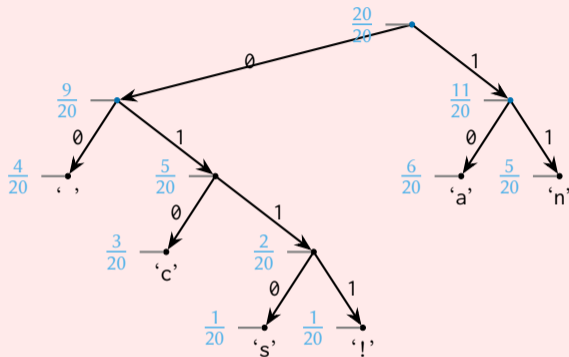| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |

$\frac{6}{20} \longrightarrow \bullet$ 'a'   $\frac{5}{20} \longrightarrow \bullet$ 'n'   $\frac{4}{20} \longrightarrow \bullet$ ','   $\frac{3}{20} \longrightarrow \bullet$ 'c'   $\frac{1}{20} \longrightarrow \bullet$ 's'   $\frac{1}{20} \longrightarrow \bullet$ '!'

**Algorithm** HUFFMANPFTRIE($f$):

4: ...
5: **while** $|Q| \geq 2$ **do**
6:   $(n_0, p_0) := \text{DELMIN}(Q), (n_1, p_1) := \text{DELMIN}(Q)$.
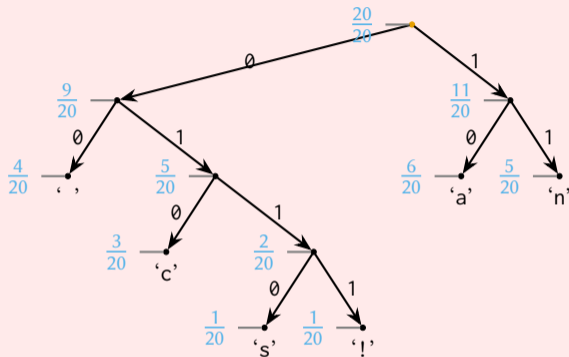7:   Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:   Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9: **return** $n$ with $(n, p) := \text{DELMIN}(Q)$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |

$\frac{6}{20}$ →  •  'a'   $\frac{5}{20}$ →  •  'n'   $\frac{4}{20}$ →  •  ' '   $\frac{3}{20}$ →  •  'c'   $\frac{1}{20}$ →  •  's'   $\frac{1}{20}$ →  •  '!'

**Algorithm** HUFFMANPFTRIE($f$):

4: …
5: **while** $|Q| \geq 2$ **do**
6:     $(n_0, p_0) := \text{DELMIN}(Q)$, $(n_1, p_1) := \text{DELMIN}(Q)$.
7:     Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:     Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
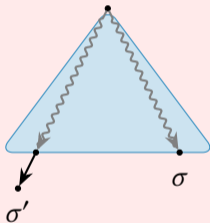9: **return** $n$ with $(n, p) := \text{DELMIN}(Q)$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |



**Algorithm** HUFFMANPFTRIE($f$):

4:  …
5:  **while** $|Q| \geq 2$ **do**
6:      $(n_0, p_0) := $ DELMIN$(Q)$, $(n_1, p_1) := $ DELMIN$(Q)$.
7:      Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:      Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9:  **return** $n$ with $(n, p) := $ DELMIN$(Q)$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |



**Algorithm** HUFFMANPFTRIE($f$):

4: …
5: **while** $|Q| \geq 2$ **do**
6:    $(n_0, p_0) :=$ DELMIN$(Q)$, $(n_1, p_1) :=$ DELMIN$(Q)$.
7:    Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:    Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9: **return** $n$ with $(n, p) :=$ DELMIN$(Q)$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |



**Algorithm** HUFFMANPFTRIE($f$):

4:  …
5:  **while** $|Q| \geq 2$ **do**
6:     $(n_0, p_0) := $ DELMIN($Q$), $(n_1, p_1) := $ DELMIN($Q$).
7:     Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:     Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9:  **return** $n$ with $(n, p) := $ DELMIN($Q$).

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |



**Algorithm** HUFFMANPFTRIE($f$):

4: …
5: **while** $|Q| \geq 2$ **do**
6:    $(n_0, p_0) :=$ DELMIN($Q$), $(n_1, p_1) :=$ DELMIN($Q$).
7:    Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:    Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9: **return** $n$ with $(n, p) :=$ DELMIN($Q$).

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |



**Algorithm** HUFFMANPFTRIE($f$):

4: ...
5: **while** $|Q| \geq 2$ **do**
6: $\quad (n_0, p_0) := \text{DELMIN}(Q), (n_1, p_1) := \text{DELMIN}(Q)$.
7: $\quad$ Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8: $\quad$ Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9: **return** $n$ with $(n, p) := \text{DELMIN}(Q)$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a' | 6 | $\frac{6}{20}$ |
| 'n' | 5 | $\frac{5}{20}$ |
| ' ' | 4 | $\frac{4}{20}$ |
| 'c' | 3 | $\frac{3}{20}$ |
| 's' | 1 | $\frac{1}{20}$ |
| '!' | 1 | $\frac{1}{20}$ |



**Algorithm** HUFFMANPFTRIE($f$):

4:  …
5:  **while** $|Q| \geq 2$ **do**
6:      $(n_0, p_0) := \text{DELMIN}(Q), (n_1, p_1) := \text{DELMIN}(Q)$.
7:      Create a node $n$ with children $n_0$ labeled 0, $n_1$ labeled 1.
8:      Add node $(n, p_0 + p_1)$ to $Q$ with priority $p_0 + p_1$.
9:  **return** $n$ with $(n, p) := \text{DELMIN}(Q)$.

# Huffman coding

### Property 1

Let $\sigma \in \mathcal{A}$ be the symbol with lowest frequency $f$.

Any optimal prefix-free trie for $\mathcal{A}, f$ can be changed such that the path to $\sigma$ is longest.

# Huffman coding

### Property 1
Let $\sigma \in \mathcal{A}$ be the symbol with lowest frequency $f$.
Any optimal prefix-free trie for $\mathcal{A}, f$ can be changed such that the path to $\sigma$ is longest.

# Huffman coding

### Property 2

Let $\sigma_1, \sigma_2 \in \mathcal{A}$ be the symbols with lowest frequency $f$.
Any optimal prefix-free trie for $\mathcal{A}, f$ can be changed such that
symbols $\sigma_1, \sigma_2$ are children of the same node.

# Huffman coding

### Property 2

Let $\sigma_1, \sigma_2 \in \mathcal{A}$ be the symbols with lowest frequency $f$.
Any optimal prefix-free trie for $\mathcal{A}, f$ can be changed such that
symbols $\sigma_1, \sigma_2$ are children of the same node.

# Huffman coding

### Property 3

Let $\sigma_1, \sigma_2 \in \mathcal{A}$ be symbols represented by children $n_0, n_1$ of node $n$ in trie $T$.
Let $T'$ be the prefix-free trie for $\mathcal{A}', f'$ with

- $\mathcal{A}' = \mathcal{A} \setminus \{\sigma_2\}$;
- $f' = \{\sigma \mapsto f(\sigma) \mid \sigma \in \mathcal{A} \setminus \{\sigma_1, \sigma_2\}\} \cup \{\sigma_1 \mapsto f(\sigma_1) + f(\sigma_2)\}$; and
- leafs $n_0, n_1$ removed from $n$ and $n$ made to represent $\sigma_1$.

The trie $T$ is optimal for $\mathcal{A}, f$ if and only if $T'$ is optimal for $\mathcal{A}', f'$.

# Huffman coding

### Property 3

Let $\sigma_1, \sigma_2 \in \mathcal{A}$ be symbols represented by children $n_0, n_1$ of node $n$ in trie $T$.
Let $T'$ be the prefix-free trie for $\mathcal{A}', f'$ with

- $\mathcal{A}' = \mathcal{A} \setminus \{\sigma_2\}$;
- $f' = \{\sigma \mapsto f(\sigma) \mid \sigma \in \mathcal{A} \setminus \{\sigma_1, \sigma_2\}\} \cup \{\sigma_1 \mapsto f(\sigma_1) + f(\sigma_2)\}$; and
- leafs $n_0, n_1$ removed from $n$ and $n$ made to represent $\sigma_1$.

The trie $T$ is optimal for $\mathcal{A}, f$ if and only if $T'$ is optimal for $\mathcal{A}', f'$.

# Huffman coding

| Symbol | Count | Frequency |
|--------|-------|-----------|
| 'a'    | 6     | $\frac{6}{20}$ |
| 'n'    | 5     | $\frac{5}{20}$ |
| ' '    | 4     | $\frac{4}{20}$ |
| 'c'    | 3     | $\frac{3}{20}$ |
| 's'    | 1     | $\frac{1}{20}$ |
| '!'    | 1     | $\frac{1}{20}$ |



### Theorem
*The HuffmanPFTrie algorithm builds an optimal prefix-free code.*

### Proof (sketch)
HuffmanPFTrie follows Property 1–3.

# Beyond Huffman: Frequent strings

Huffman looks at frequent symbols from an alphabet.

- ▶ We can generalize these ideas to *frequent* sequences of symbols.
- ▶ Tries can be used to efficiently manage *frequency* data for substrings in an input.
- ▶ Challenge: which substrings to consider?
  E.g., fixed length, maximum length, . . . .

# Beyond Huffman: Frequent strings

Huffman looks at frequent symbols from an alphabet.

- ▶ We can generalize these ideas to *frequent* sequences of symbols.
- ▶ Tries can be used to efficiently manage *frequency* data for substrings in an input.
- ▶ Challenge: which substrings to consider?
  E.g., fixed length, maximum length, ....

*Many* variations of this idea used in practice, e.g., .zip, .gif, ....

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $S$ (the haystack) and $P$ (the needle or pattern),
return the first position in $S$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($S$, $P$):
1: **for** $i := 0$ upto $|S| - 1$ **do**
2:     **if** pattern $P$ starts at $S[i]$ **then**
3:        **return** $i$.

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BasicStringSearch($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:    **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:       **return** $i$.

Searching $P =$ "*example*"

```
" a n   e x a m p l e   o f   w o r d s "
   e x a m p l e
```

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}, P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:    **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:       **return** $i$.

Searching $P =$ "*example*"
" a n   e x a m p l e   o f   w o r d s "
    e x a m p l e

# Basic substring search: Searching a needle in a haystack

### Problem

Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):

1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:     **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:         **return** $i$.

Searching $P = $ "*example*"

```
" a n   e x a m p l e   o f   w o r d s "
      e x a m p l e
```

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:     **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:         **return** $i$.

Searching $P =$ "*example*"
```
            " a n   e x a m p l e   o f   w o r d s "
                    e x a m p l e
```

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:   **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:     **return** $i$.

Searching $P =$ "*example*"
```
" a n   e x a m p l e   o f   w o r d s "
      e x a m p l e
```

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:   **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:     **return** $i$.

Searching $P =$ "*great*"

```
" a n   e x a m p l e   o f   w o r d s "
                        g r e a t
```

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:   **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:     **return** $i$.

Searching $P =$ "*great*"
                  " a n   e x a m p l e   o f   w o r d s "
                                        g r e a t

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:    **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:       **return** $i$.

Searching $P =$ "*great*"

```
" a n   e x a m p l e   o f   w o r d s "
                          g r e a t
```

# Basic substring search: Searching a needle in a haystack

## Problem

Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:     **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:         **return** $i$.

Searching $P =$ "*great*"

                "  a n    e x a m p l e    o f    w o r d s "

                                         g r e a t

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BasicStringSearch($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:    **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:       **return** $i$.

Searching $P =$ "*great*"
```
" a n   e x a m p l e   o f   w o r d s "
                              g r e a t
```

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - 1$ **do**
2:     **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:         **return** $i$.

Searching $P =$ "*great*"
        "  a n     e x a m p l e     o f     w o r d s "
                                        g r e a t

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BasicStringSearch($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - |P|$ **do**
2:    **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:       **return** $i$.

Searching $P =$ "*great*"
                    " a n   e x a m p l e   o f   w o r d s "
                                          g r e a t

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - |P|$ **do**
2:    **if** pattern $P$ starts at $\mathcal{S}[i]$ **then**
3:       **return** $i$.

### Complexity

# Basic substring search: Searching a needle in a haystack

## Problem

Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - |P|$ **do**
2:    **if** MATCHSTRING($\mathcal{S}$, $P$, $i$) **then**
3:       **return** $i$.

**Algorithm** MATCHSTRING($\mathcal{S}$, $P$, $i$):
4: **for** $j := 0$ upto $|P| - 1$ **do**
5:    **if** $\mathcal{S}[i + j] \neq P[j]$ **then**
6:       **return** false.
7: **return** true.

## Complexity

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):
1: **for** $i := 0$ upto $|\mathcal{S}| - |P|$ **do**
2:    **if** MATCHSTRING($\mathcal{S}$, $P$, $i$) **then**
3:      **return** $i$.

### Complexity

**Algorithm** MATCHSTRING($\mathcal{S}$, $P$, $i$):
4: **for** $j := 0$ upto $|P| - 1$ **do**
5:    **if** $\mathcal{S}[i + j] \neq P[j]$ **then**
6:      **return** false.
7: **return** true.

$\Theta(|P|)$

# Basic substring search: Searching a needle in a haystack

### Problem

Given strings $\mathcal{S}$ (the haystack) and $P$ (the needle or pattern),
return the first position in $\mathcal{S}$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($\mathcal{S}$, $P$):

1: **for** $i := 0$ upto $|\mathcal{S}| - |P|$ **do**
2:     **if** MATCHSTRING($\mathcal{S}$, $P$, $i$) **then** $\left.\rule{0pt}{3em}\right\}\Theta((|\mathcal{S}| - |P|)|P|)$
3:         **return** $i$.

### Complexity

# Basic substring search: Searching a needle in a haystack

### Problem
Given strings $S$ (the haystack) and $P$ (the needle or pattern),
return the first position in $S$ at which $P$ occurs (if any).

**Algorithm** BASICSTRINGSEARCH($S$, $P$):
1: **for** $i := 0$ upto $|S| - |P|$ **do**
2:    **if** MATCHSTRING($S$, $P$, $i$) **then** $\left.\right\}$ $\Theta((|S| - |P|)|P|) = \Theta(|S||P|)$
3:       **return** $i$.

### Complexity

# Substring search: Room for improvement

Searching $P = $ "*string*"

"  a   s t r o n g     s t r i n g  "

# Substring search: Room for improvement

Searching $P$ = "*string*"
```
" a   s t r o n g    s t r i n g "
  s t r i n g
```

# Substring search: Room for improvement

Searching $P$ = "*string*"

```
" a   s t r o n g     s t r i n g "
    s t r i n g
```

# Substring search: Room for improvement

Searching $P$ = "*string*"

```
" a   s t r o n g   s t r i n g "
      s t r i n g
```

# Substring search: Room for improvement

Searching $P = $ "*string*"
```
" a   s t r o n g     s t r i n g "
        s t r i n g
```

# Substring search: Room for improvement

Searching $P = $ "*string*"

```
" a   s t r o n g     s t r i n g "
      s t r i n g
```

# Substring search: Room for improvement

Searching $P$ = "*string*"

```
" a   s t r o n g    s t r i n g "
      s t r i n g
```

# Substring search: Room for improvement

Searching $P$ = "*string*"

```
" a   s t r o n g   s t r i n g "
            s t r i n g
```

# Substring search: Room for improvement

Searching $P =$ "*string*"

```
" a   s t r o n g   s t r i n g "
            s t r i n g
```

# Substring search: Room for improvement

Searching $P$ = "*string*"

```
          " a   s t r o n g    s t r i n g "
                              s t r i n g
```

# Substring search: Room for improvement

Searching $P$ = "*string*"

```
" a   s t r o n g   s t r i n g "
                    s t r i n g
```

# Substring search: Room for improvement

Searching $P =$ "*string*"

```
" a   s t r o n g    s t r i n g "
               s t r i n g
```

# Substring search: Room for improvement

Searching $P = $ "*string*"

```
" a   s t r o n g   s t r i n g "
                    s t r i n g
```

# Substring search: Room for improvement

Searching $P =$ "*ACACGT*"

" A C A C A C A C G T "

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

```
" A C A C A C A C G T "
  A C A C G T
```

# Substring search: Room for improvement

Searching $P = $ "*ACACGT*"

```
" A C A C A C A C G T "
  A C A C G T
```

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

<div align="center">

" A C A C A C A C G T "

A C A C G T

</div>

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

```
    " A C A C A C A C G T "
      A C A C G T
```

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

                                               " A C A C A C A C G T "

                                                A C A C G T

# Substring search: Room for improvement

Searching $P = $ "*ACACGT*"

"  A C A C A C A C G T  "
   A C A C G T

# Substring search: Room for improvement

Searching $P = $ "*ACACGT*"

```
" A C A C A C A C G T "
      A C A C G T
```

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

```
" A C A C A C A C G T "
      A C A C G T
```

Searching $P$ = "*ACACGT*"

```
" A C A C A C A C G T "
      A C A C G T
```

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

```
" A C A C A C A C G T "
        A C A C G T
```

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

```
" A C A C A C A C G T "
        A C A C G T
```

# Substring search: Room for improvement

Searching $P =$ "*ACACGT*"

```
" A C A C A C A C G T "
        A C A C G T
```

# Substring search: Room for improvement

Searching $P$ = "*ACACGT*"

```
" A C A C A C A C G T "
        A C A C G T
```

Searching $P$ = "*ACACGT*"

" A C A C A C A C G T "
            A C A C G T

# Encode search patterns as an automaton

### Finite automata

A finite automaton is a *graph* with

- a single *initial* node;
- zero-or-more *final* nodes;
- edges labeled with *symbols*.

# Encode search patterns as an automaton

## Finite automata

A finite automaton is a *graph* with

- a single *initial* node;
- zero-or-more *final* nodes;
- edges labeled with *symbols*.



Typically, we refer to nodes as *states* and edges as *transitions*.

# Encode search patterns as an automaton

### Finite automata
A finite automaton is a *graph* with

- a single *initial* node;
- zero-or-more *final* nodes;
- edges labeled with *symbols*.



### Running a finite automaton
We can use a string as *input* to the automaton to decide which path to follow.

# Encode search patterns as an automaton

## Finite automata

A finite automaton is a *graph* with

- a single *initial* node;
- zero-or-more *final* nodes;
- edges labeled with *symbols*.



## Running a finite automaton

We can use a string as *input* to the automaton to decide which path to follow.

"ACACACACGT".

# Encode search patterns as an automaton

### Finite automata
A finite automaton is a *graph* with

- a single *initial* node;
- zero-or-more *final* nodes;
- edges labeled with *symbols*.



### Running a finite automaton
We can use a string as *input* to the automaton to decide which path to follow.
For efficiency: we want a *deterministic* automaton: an automaton without choices!

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton
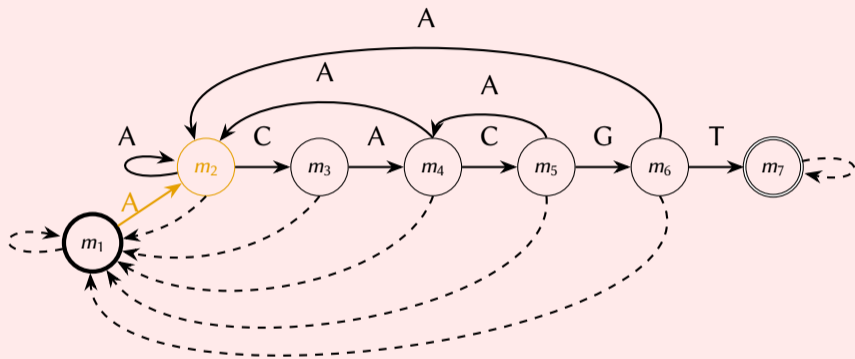
# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton

# Encode search patterns as an automaton
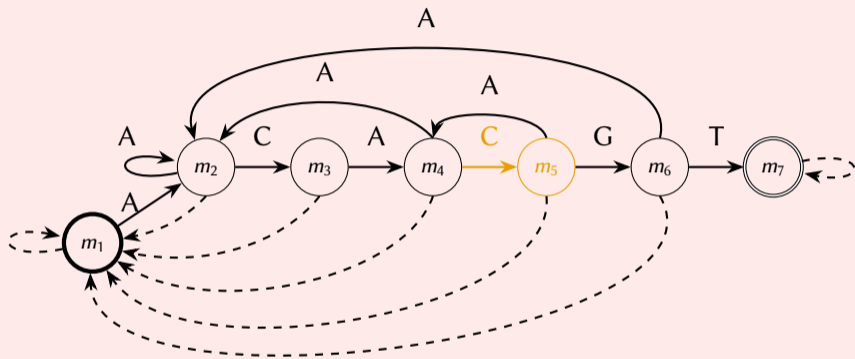
# Encode search patterns as an automaton

# Encode search patterns as an automaton



Searching $P$ = "*ACACGT*"

"  A  C  A  C  A  C  A  C  G  T  "

# Encode search patterns as an automaton



Searching $P$ = "ACACGT"
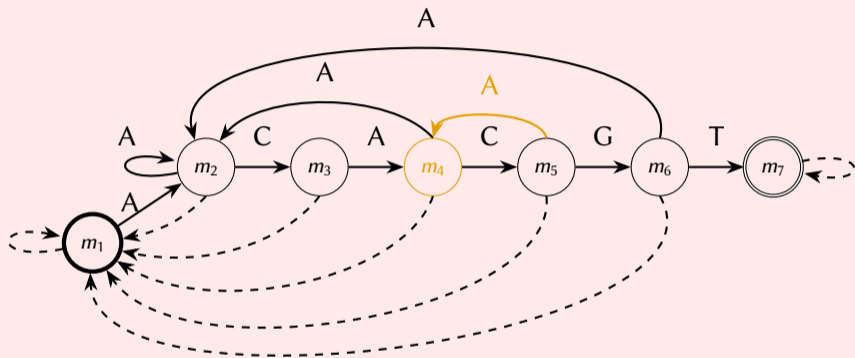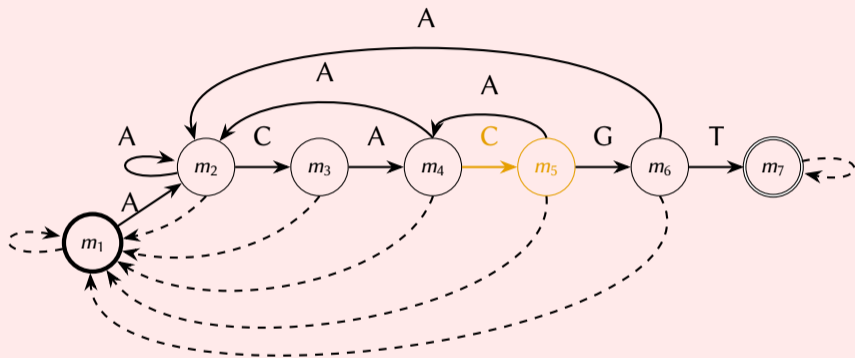
<p style="text-align:center">" A C A C A C G T "</p>

# Encode search patterns as an automaton



Searching $P$ = "$ACACGT$"

"A C A C A C A C G T"

# Encode search patterns as an automaton



Searching $P = $ "*ACACGT*"

"  A  C  A  C  A  C  A  C  G  T  "

# Encode search patterns as an automaton



Searching $P$ = "*ACACGT*"

"  A  C  A  C  A  C  A  C  G  T  "

# Encode search patterns as an automaton



Searching $P$ = "*ACACGT*"

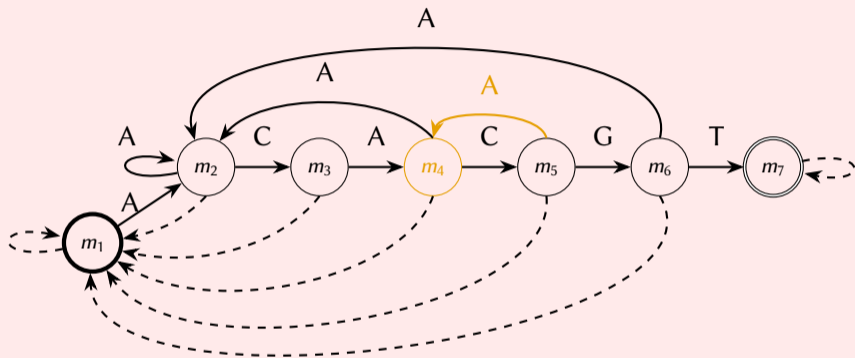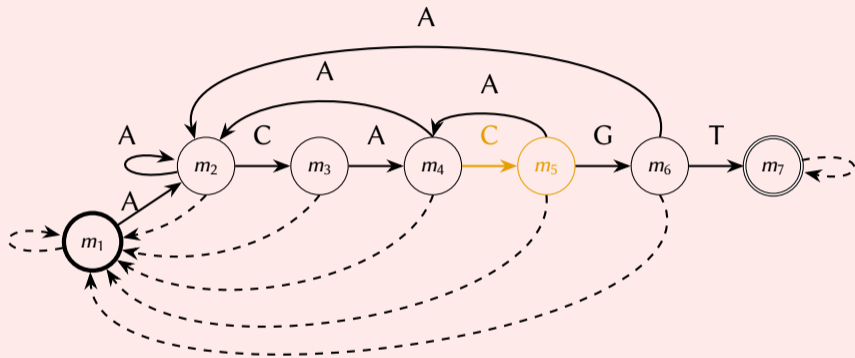                    " A C A C A C A C G T "

# Encode search patterns as an automaton



Searching $P$ = "$ACACGT$"

"  A  C  A  C  A  C  A  C  G  T  "

# Encode search patterns as an automaton



Searching $P$ = "*ACACGT*"

"  A  C  A  C  A  C  A  C  G  T  "

# Encode search patterns as an automaton



Searching $P$ = "$ACACGT$"
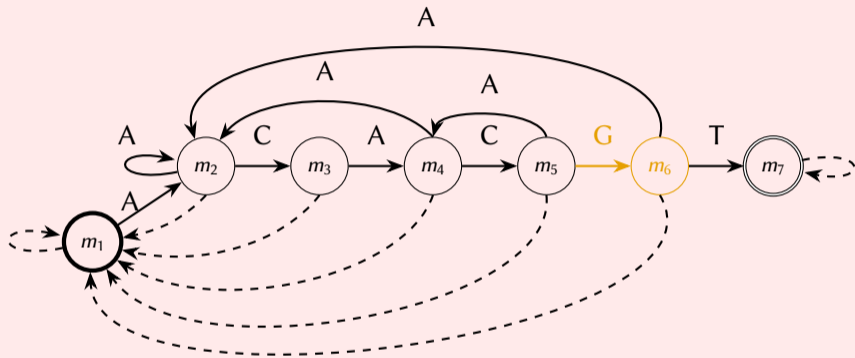
"  A  C  A  C  A  C  A  C  G  T  "

# Encode search patterns as an automaton



Searching $P$ = "*ACACGT*"

"  A  C  A  C  A  C  A  C  G  T  "
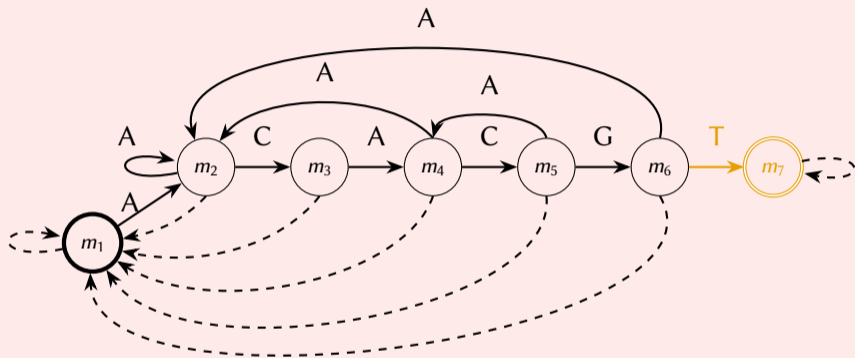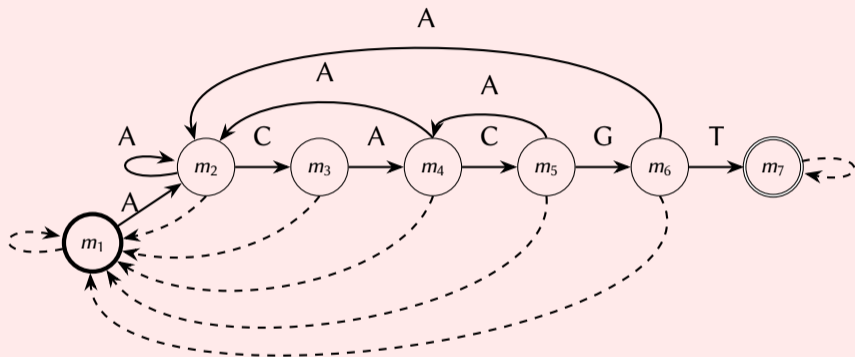
# Encode search patterns as an automaton



Searching $P$ = "*ACACGT*"

" A C A C A C A C G T "

# Encode search patterns as an automaton



Complexity of running a *deterministic* finite automaton with input $\mathcal{S}$

- Always in exactly *one* state.
- We perform at-most $\Theta(|\mathcal{S}|)$ *state transitions* in the automaton.
- Need efficient representation of the *transitions* (per state): e.g., hash table.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

- ▶ $\emptyset$ describes $\emptyset$.
- ▶ $\sigma \in \mathcal{A}$ describes $\{\sigma\}$.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

- ▶ $\emptyset$ describes $\emptyset$.
- ▶ $\sigma \in \mathcal{A}$ describes $\{\sigma\}$.

Now let $e, e_1, e_2$ be regular expressions describing sets $R, R_1, R_2$.

- ▶ $(e)$ desribes $R$.
- ▶ $e_1 e_2$ describes $\{\textsc{Concatenate}(\mathcal{S}_1, \mathcal{S}_2) \mid \mathcal{S}_1 \in R_1 \land \mathcal{S}_2 \in R_2\}$.
- ▶ $e_1 \mid e_2$ describes $R_1 \cup R_2$.
- ▶ $e^*$ describes any sequence of strings in $R$.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

- ▶ $\emptyset$ describes $\emptyset$.
- ▶ $\sigma \in \mathcal{A}$ describes $\{\sigma\}$.

Now let $e, e_1, e_2$ be regular expressions describing sets $R, R_1, R_2$.

- ▶ $(e)$ desribes $R$.
- ▶ $e_1 e_2$ describes $\{\text{Concatenate}(\mathcal{S}_1, \mathcal{S}_2) \mid \mathcal{S}_1 \in R_1 \land \mathcal{S}_2 \in R_2\}$.
- ▶ $e_1 \mid e_2$ describes $R_1 \cup R_2$.
- ▶ $e^*$ describes any sequence of strings in $R$.

Examples

$$\text{moose} \mid \text{mouse}$$

$$\text{sub}^*\text{section}$$

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

- ▶ $\emptyset$ describes $\emptyset$.
- ▶ $\sigma \in \mathcal{A}$ describes $\{\sigma\}$.

Now let $e, e_1, e_2$ be regular expressions describing sets $R, R_1, R_2$.

- ▶ $(e)$ desribes $R$.
- ▶ $e_1 e_2$ describes $\{\textsc{Concatenate}(\mathcal{S}_1, \mathcal{S}_2) \mid \mathcal{S}_1 \in R_1 \wedge \mathcal{S}_2 \in R_2\}$.
- ▶ $e_1 \mid e_2$ describes $R_1 \cup R_2$.
- ▶ $e^*$ describes any sequence of strings in $R$.

Examples

moose | mouse

sub$^*$section

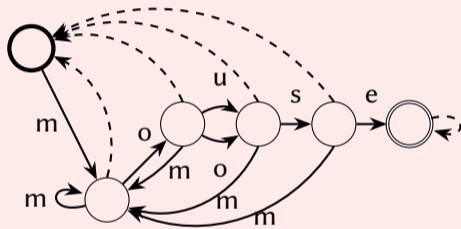Claim: Every regular expression is equivalent to a deterministic finite automaton
See SFWRENG 2FA3: Discrete Mathematics with Applications II.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

Claim: Every regular expression is equivalent to a deterministic finite automaton
See SFWRENG 2FA3: Discrete Mathematics with Applications II.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

Claim: Every regular expression is equivalent to a deterministic finite automaton
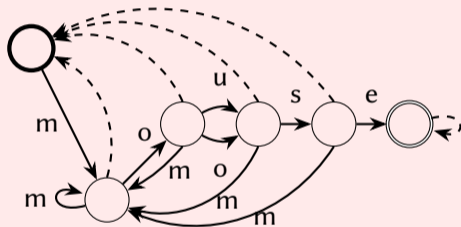See SFWRENG 2FA3: Discrete Mathematics with Applications II.



Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Claim: Every regular expression is equivalent to a deterministic finite automaton
See SFWRENG 2FA3: Discrete Mathematics with Applications II.



## Deterministic finite automata can grow *very* large
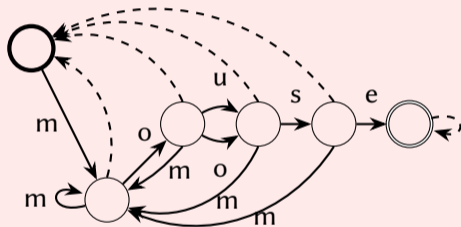Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
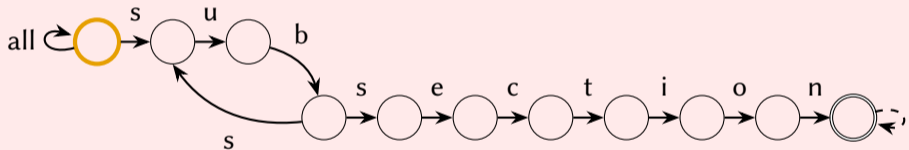Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"
    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
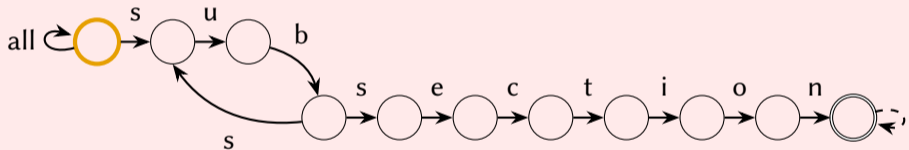
# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub$^*$section"

"a   s u b s t r i n g   i n   a   s u b s e c t i o n"

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
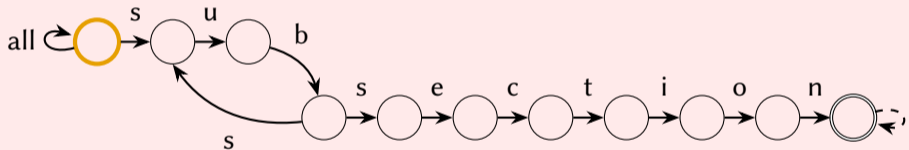(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub*section"
    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
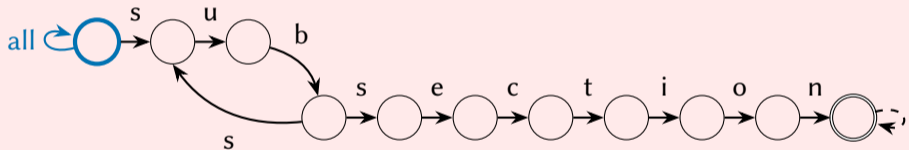(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
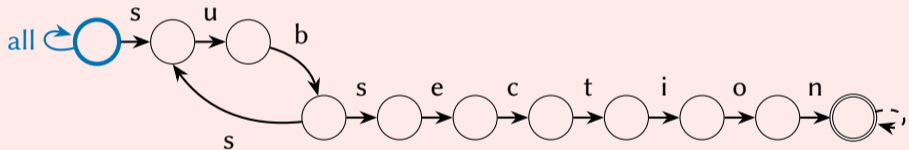(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
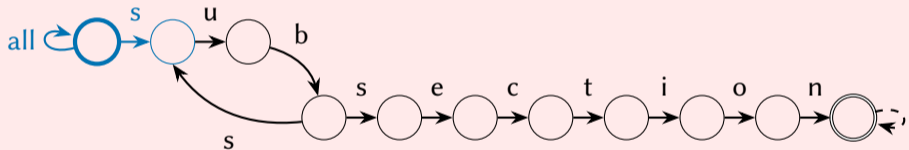(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
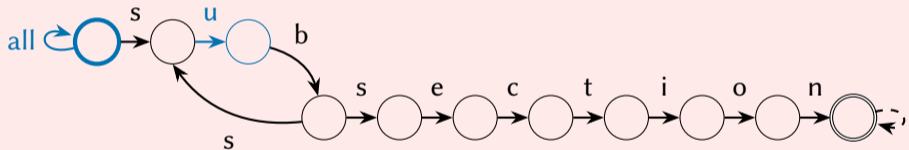(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub$^*$section"

     " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
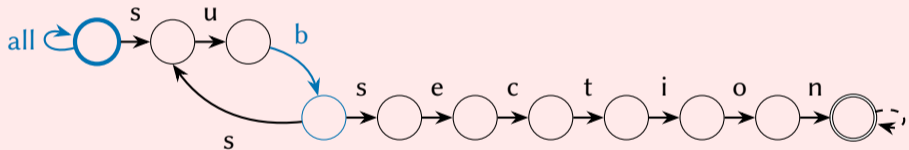(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub*section"

```
" a   s u b s t r i n g   i n   a   s u b s e c t i o n "
```

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
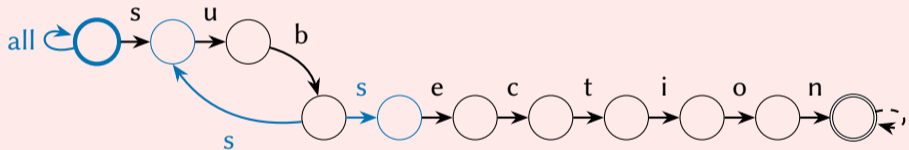(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
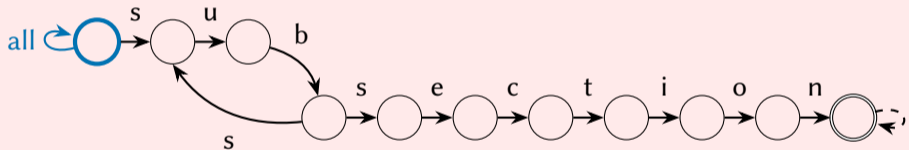(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
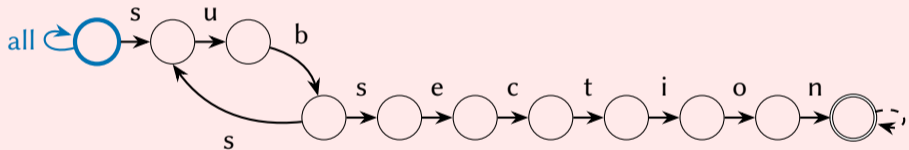(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea*. Keep track of the *set of states* we can be in while walking to the string.
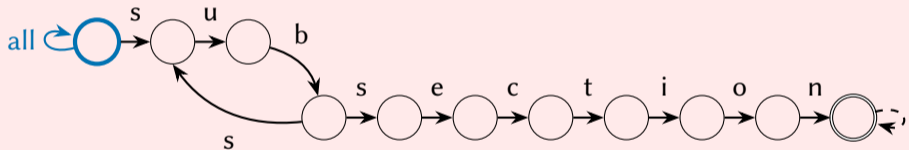(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
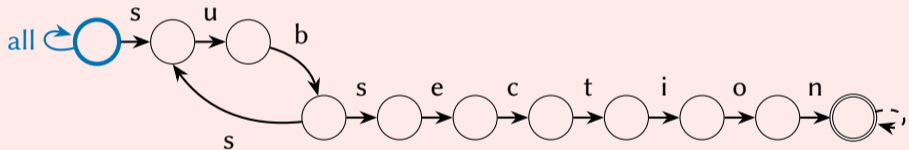(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub*section"

"  a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
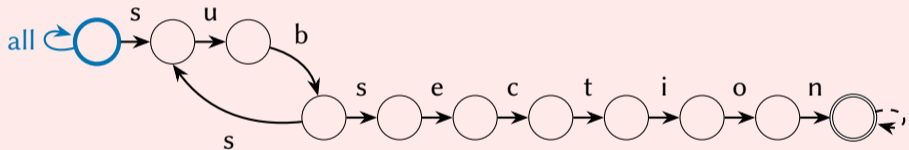(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

```
" a   s u b s t r i n g   i n   a   s u b s e c t i o n "
```

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
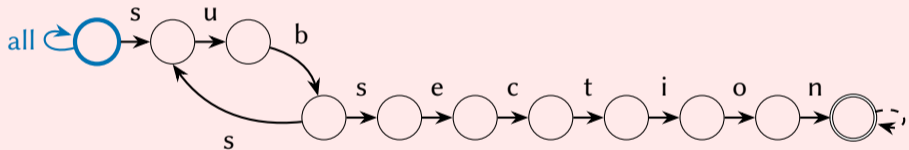(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
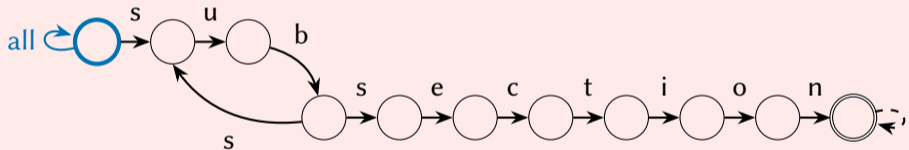(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub*section"

```
" a   s u b s t r i n g   i n   a   s u b s e c t i o n "
```

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
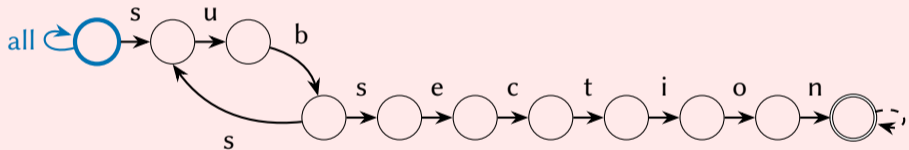(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
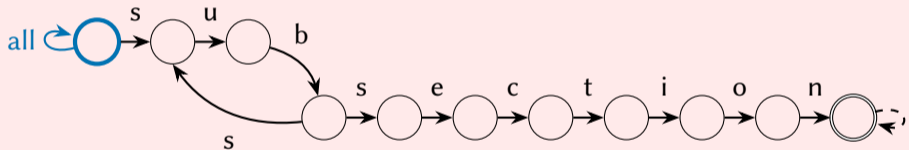(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
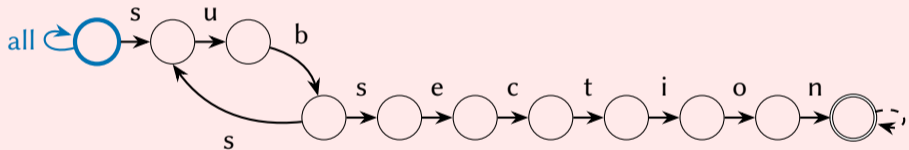(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub*section"

    " a    s u b s t r i n g    i n    a    s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
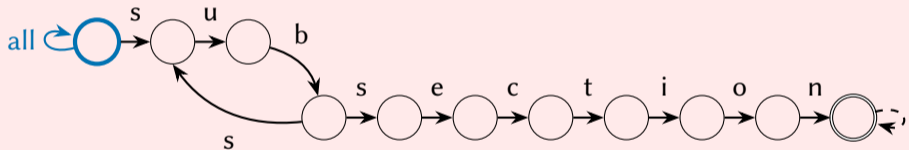(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub*section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
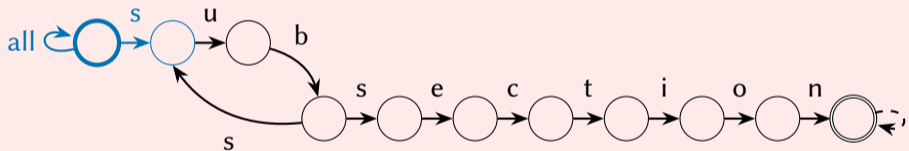(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
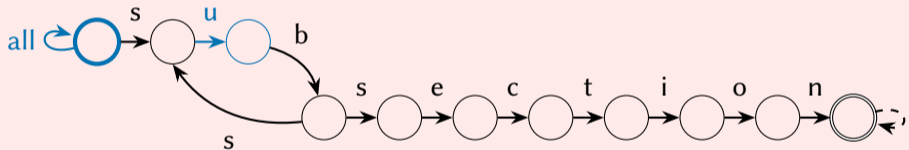(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

"  a   s u b s t r i n g   i n   a   s u b s e c t i o n  "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
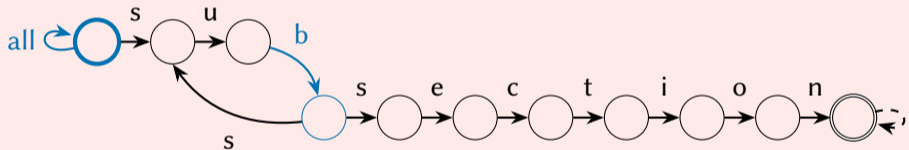(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
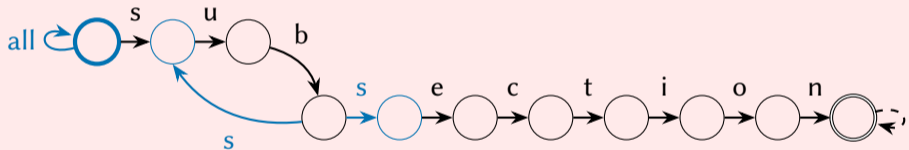(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
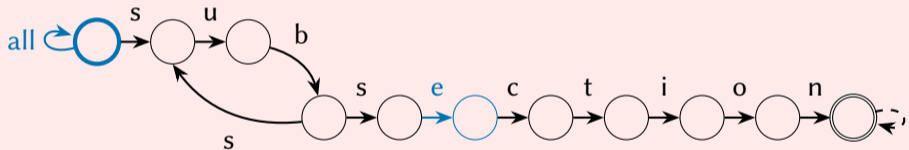(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P =$ "sub*section"
    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P = $ "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea*. Keep track of the *set of states* we can be in while walking to the string.
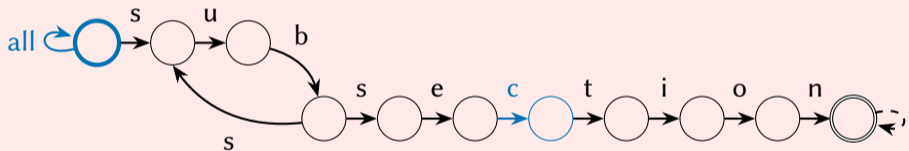(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large

Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.



Searching $P$ = "sub$^*$section"

    " a   s u b s t r i n g   i n   a   s u b s e c t i o n "

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
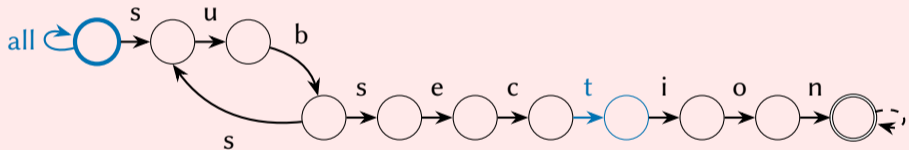(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Deterministic finite automata can grow *very* large
Alternatively, you can run a *nondeterminstic* finite automaton: an automaton with choices!
Lower costs to construct the automaton, higher costs to run them.

Searching $P$ = "sub*section"

```
" a   s u b s t r i n g   i n   a   s u b s e c t i o n "
```

*Idea.* Keep track of the *set of states* we can be in while walking to the string.
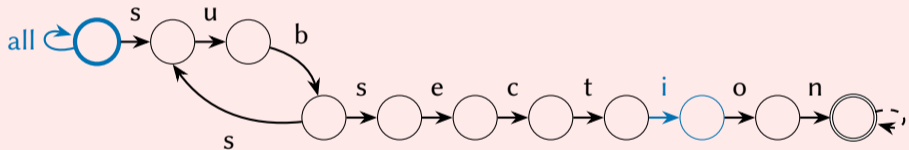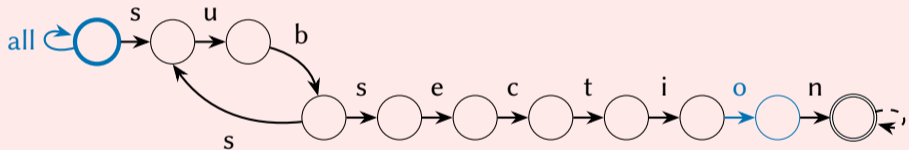(A *powerset construction* guided by the string one is searching in).

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Final remarks on regular expressions

- Running an $N$-state deterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(|\mathcal{S}|)$.
- Running an $N$-state nondeterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(N|\mathcal{S}|)$.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Final remarks on regular expressions

- ▶ Running an $N$-state deterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(|\mathcal{S}|)$.
- ▶ Running an $N$-state nondeterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(N|\mathcal{S}|)$.
- ▶ For single words, a deterministic finite automaton can be easily constructed:
  see the Knuth-Morris-Pratt algorithm in the book.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Final remarks on regular expressions

▶ Running an $N$-state deterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(|\mathcal{S}|)$.

▶ Running an $N$-state nondeterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(N|\mathcal{S}|)$.

▶ For single words, a deterministic finite automaton can be easily constructed:
  see the Knuth-Morris-Pratt algorithm in the book.

▶ Automata-based searching in string can be very fast:
  core component in lexers, parsers, and compilers.

# More general search patterns: Regular expressions

A *regular expression* describes a set of *strings*.

## Final remarks on regular expressions

- ▶ Running an $N$-state deterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(|\mathcal{S}|)$.

- ▶ Running an $N$-state nondeterministic finite automaton on $\mathcal{S}$ can be done in $\Theta(N|\mathcal{S}|)$.

- ▶ For single words, a deterministic finite automaton can be easily constructed: see the Knuth-Morris-Pratt algorithm in the book.

- ▶ Automata-based searching in string can be very fast: core component in lexers, parsers, and compilers.

- ▶ Many *RegExp* libraries are *regular expression like*: they support non-regular features. Consequently, many such libraries use *shamefully slow* backtracking algorithms: Worst-case exponential complexity, even when searching for *simple patterns*.

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P$ = "*great*"

```
" a n   e x a m p l e   o f   w o r d s "
  g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P$ = "*great*"

```
" a n   e x a m p l e   o f   w o r d s "
  g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?
Assume we checked up-till position $i$.

*Observation.* If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$.*

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P = $ "*great*"

```
" a n    e x a m p l e    o f    w o r d s "
             g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?

Assume we checked up-till position $i$.

*Observation.* If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$.*

▶ We see a symbol that is in $P$: *P could have started somewhere in $S[i \ldots i + |P|)$.*

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P =$ "*great*"

```
" a n     e x a m p l e     o f     w o r d s "
                    g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?
Assume we checked up-till position $i$.

*Observation.* If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$.*

▶ We see a symbol that is in $P$: *P could have started somewhere in $S[i \ldots i + |P|)$.*
   Based on the symbol, decide where to start looking.

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P =$ "*great*"

```
" a n   e x a m p l e   o f   w o r d s "
                  g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?
Assume we checked up-till position $i$.

*Observation.* If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

- ▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$.*
- ▶ We see a symbol that is in $P$: *P could have started somewhere in $S[i \ldots i + |P|)$.*
  Based on the symbol, decide where to start looking.

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P$ = "*great*"
```
          " a n    e x a m p l e    o f    w o r d s "
                                      g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?
Assume we checked up-till position $i$.

*Observation.* If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$.*

▶ We see a symbol that is in $P$: *P could have started somewhere in $S[i \ldots i + |P|)$.*
Based on the symbol, decide where to start looking.

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P = $ "*great*"

```
          " a n     e x a m p l e     o f     w o r d s "
                                        g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?
Assume we checked up-till position $i$.

*Observation.* If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$.*

▶ We see a symbol that is in $P$: *P could have started somewhere in $S[i \ldots i + |P|)$.*
  Based on the symbol, decide where to start looking.

▶ If the symbols match (never in this example): we might have found pattern $P$.
  Inspect from the end-of-$P$ to the begin to check.

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P =$ "*great*"

```
" a n   e x a m p l e   o f   w o r d s "
                                  g r e a t
```

How can I skip checking most letters in $\mathcal{S}$?
Assume we checked up-till position $i$.

*Observation*. If we compare the last symbol from our pattern with $S[i + |P| - 1]$, then

- ▶ We see a symbol that is not even in $P$: *P cannot occur in $S[i \ldots i + |P|)$*.
- ▶ We see a symbol that is in $P$: *P could have started somewhere in $S[i \ldots i + |P|)$.*
  Based on the symbol, decide where to start looking.
- ▶ If the symbols match (never in this example): we might have found pattern $P$.
  Inspect from the end-of-$P$ to the begin to check.
  *When it does not matches: another opportunity to jump!*

# Substring search in less-than $\Theta(|\mathcal{S}|)$ steps

Searching $P = $ "*great*"

```
" a n    e x a m p l e    o f    w o r d s "
                                  g r e a t
```

## How can I skip checking most letters in $\mathcal{S}$?

With some preprocessing on $P$ one can precompute how to jump around optimally:
the Boyer-Moore algorithm.