

Basic notions of algorithms

Mihaela Zavolan

Programming for Life Sciences 2024

Why do we need to talk about algorithms?

Algorithm = specification ('recipe') for solving a certain type of problem

- What's the need?
 - "Not reinvent the wheel" – once we have an optimal solution to a problem we can keep applying it in whatever context
 - Solutions to many and many different types of problem can be specified in terms of combinations of steps that come from a limited repertoire – e.g. parse a sentence into its parts vs. find secondary structures in RNAs vs. find binding sites for transcription factors in a genome. If we have a solution to one problem and we specify it in abstract terms, then we can apply it to other problems, e.g. by changing the type of variables/entities that we operate with.

What properties should an algorithm have?

- Unambiguous – should be interpretable by a machine in only one way
- General – applicable to a class of problems, should handle all possible parameter values
- Efficient – execution time should be as short as possible
- Memory efficient – should not require more memory than available

Key algorithm analyses:

- Correctness - Invariants
- Resource bounds – Scaling of memory and computing time

Estimating the resource requirements of a program

When we write a program to solve a certain problem we expect to it to return a solution on any specific set of inputs. But will it?

Analysis of algorithms is a branch of computer science that concerns itself with determining the resources (time and memory) that are needed by a specific algorithm to arrive at an answer.

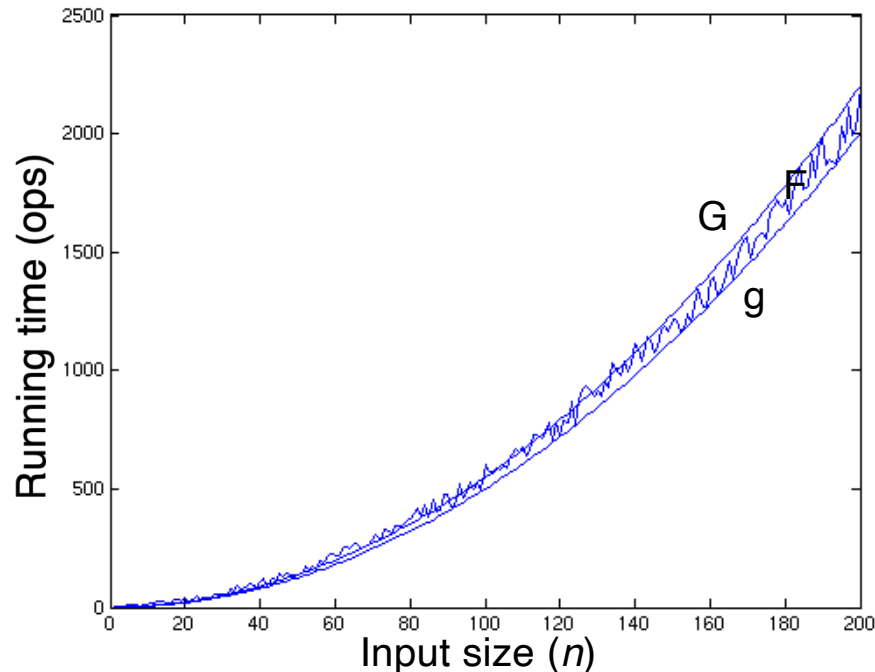
Important are not only (or less so) the absolute estimates of time and memory that the program needs when running on a specific input, but rather how these quantities scale with the size of the input.

Let's say we want to find out whether a specific sequence has been deposited in the NCBI sequence database. How long do we expect this to take?

Similarly, let's say that we want to find out whether there is any peptide annotated in a given genome that has a predicted mass equal to a mass that we measured experimentally. How long would this take?

Time complexity of a program: Big-O/little-o notation

Assume that the running time of a program defined as the number of 'atomic' operations (additions/subtractions/comparisons/etc) varies with the size of the input as shown in the figure, and let this empirical function be $F(n)$.



We say that

$$\begin{array}{ll} G(n) \text{ is } O(F(n)) & \text{if } F(n) \leq G(n) \\ g(n) \text{ is } o(F(n)) & \text{if } F(n) \geq g(n) \end{array}$$

for all n sufficiently large.

What we are interested in is the functions G and g that satisfy the properties above, and are as close as possible to F .

These will have the form

$$a * F + b$$

with a and b constants.

Time complexity of a program: Big-O/little-o notation

Let's take an example: a program that reads in a sequence of numbers and then computes their mean and variance. In Python, the program would look like this:

```
import re

p = re.compile('[^\d\.\-\+]', re.IGNORECASE)
inputSeq = []

input_var = 'yes'
while(input_var.lower() != 'quit' ):
    input_var = input("Enter a number or 'quit' to stop: ")
    if(input_var.lower() != 'quit' and not p.search(input_var)):
        print("Adding " + str(input_var) + " to the list")
        inputSeq.append(float(input_var))

mean = 0.0
var = 0.0
num_elem = len(inputSeq)

for i in range(num_elem):
    mean += inputSeq[i]
    var += (inputSeq[i] * inputSeq[i])

mean /= num_elem
var /= num_elem

print("Mean is " + str(mean) + " and variance is " + str(var - mean * mean))
```

How does the run time of this program scales with the size of the input?

Ubiquitous computational tasks: sort and search

Let's take another toy example: we get as input a list of numbers and we want to output it sorted.

A simple way to do it:

define the start index of the unsorted part of the list to be 0 (first element)

iterate until start index $< L$

find the smallest element in the list that starts at start index start and ends at L

swap this element with the element at start index

12	4	5	18	9	7	11
0						L-1

min = 12 min = 4 min = 4 min = 4 min = 4 min = 4 min = 4

start = 0

12	4	5	18	9	7	11
----	---	---	----	---	---	----

min = 12 min = 5 min = 5 min = 5 min = 5 min = 5

start = 1

4	12	5	18	9	7	11
---	----	---	----	---	---	----

min = 12 min = 12 min = 9 min = 7 min = 7

start = 2

4	5	12	18	9	7	11
---	---	----	----	---	---	----

How does the running time scale with the length of the list?

Ubiquitous computational tasks: sort and search

Iteration 1: L-1 comparisons

Iteration 2: L-2 comparisons

...

Iteration L-1: 1 comparison

$$(L - 1) + (L - 2) + \dots + 1 = \frac{L(L - 1)}{2} = O(L^2)$$

min = 12 min = 4 min = 4 min = 4 min = 4 min = 4 min = 4

start = 0

12	4	5	18	9	7	11
----	---	---	----	---	---	----

min = 12 min = 5 min = 5 min = 5 min = 5 min = 5

start = 1

4	12	5	18	9	7	11
---	----	---	----	---	---	----

min = 12 min = 12 min = 9 min = 7 min = 7

start = 2

4	5	12	18	9	7	11
---	---	----	----	---	---	----

How does the running time scale with the length of the list?

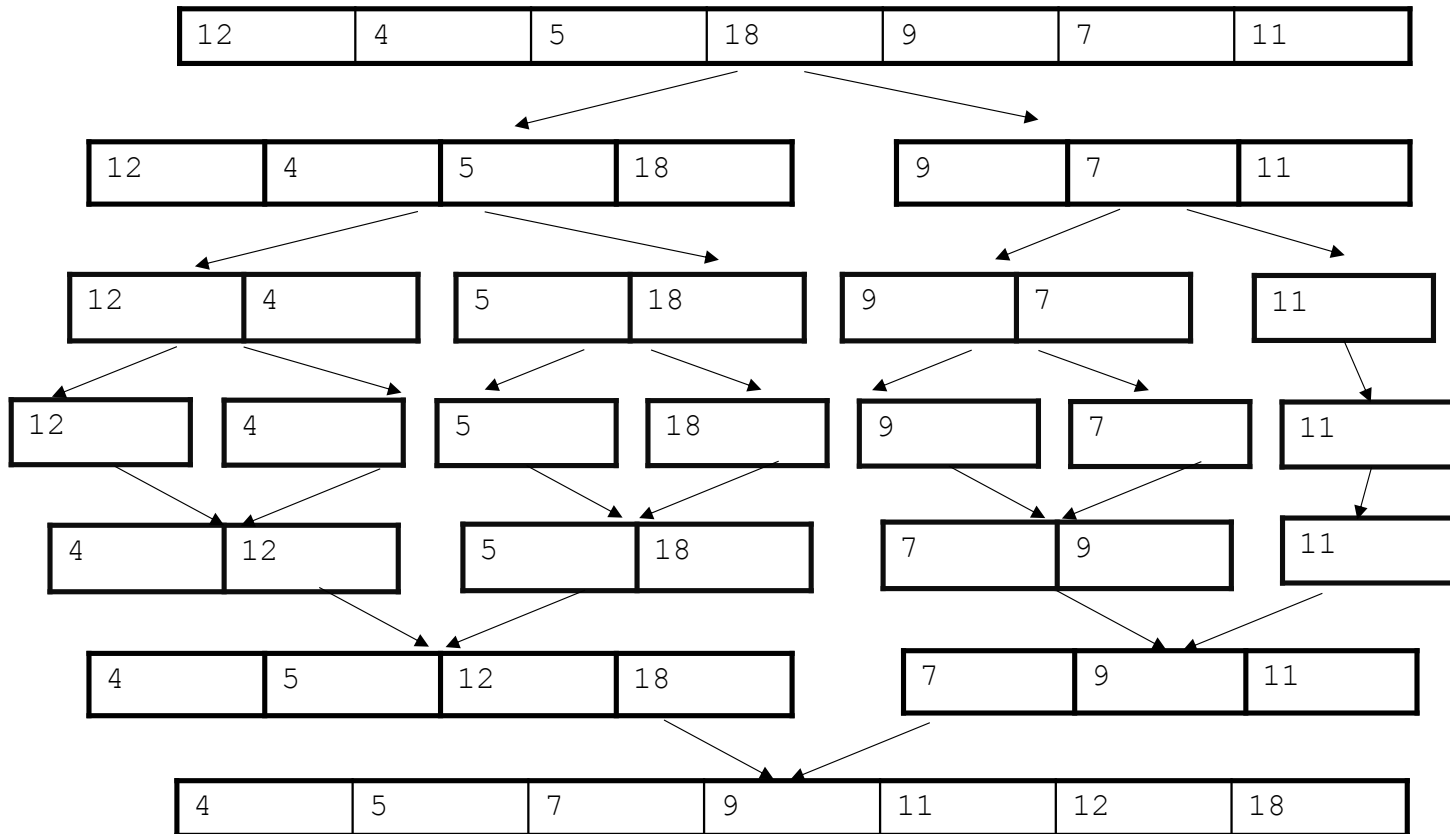
Ubiquitous computational tasks: sort and search

How would you prove the correctness of this algorithm?

	min = 12 min = 4 min = 4 min = 4 min = 4 min = 4 min = 4						
start = 0	12	4	5	18	9	7	11
	min = 12 min = 5 min = 5 min = 5 min = 5 min = 5 min = 5						
start = 1	4	12	5	18	9	7	11
	min = 12 min = 12 min = 9 min = 7 min = 7						
start = 2	4	5	12	18	9	7	11

Can we do better than quadratic time in sorting an array?

Merge sort



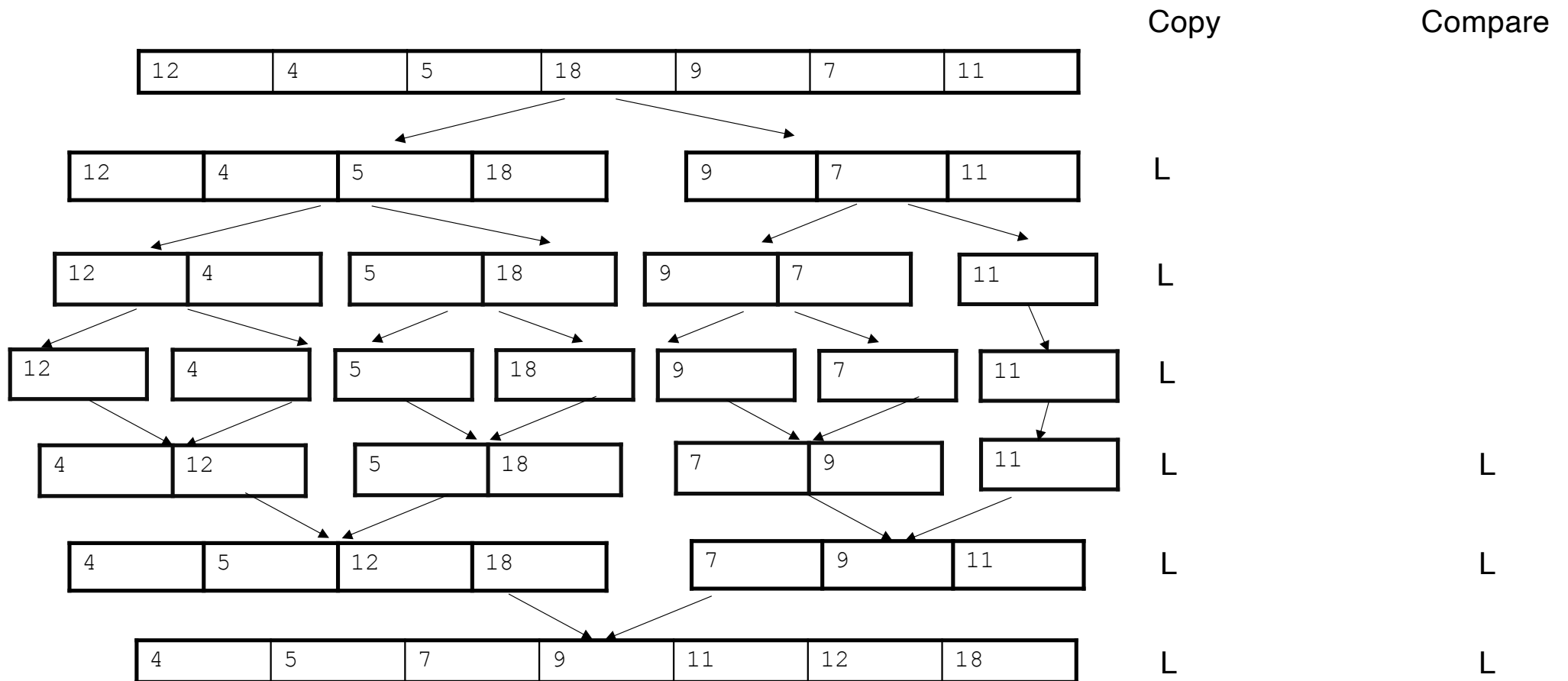
Merge sort

```
def mergeSort(A):
    if len(A) < 2:
        return
    left, right = A[ : len(A) // 2], A[len(A) // 2 : ]
    mergeSort(left)
    mergeSort(right)
    merge(left, right, A)

def merge(left, right, A):
    i, j, k = 0, 0, 0
    while i < len(left) and j < len(right):
        if left[i] <= right[j]:
            A[k], i, k = left[i], i + 1, k + 1
        else:
            A[k], j, k = right[j], j + 1, k + 1
    while i < len(left):
        A[k], i, k = left[i], i + 1, k + 1
    while j < len(right):
        A[k], j, k = right[j], j + 1, k + 1
```

Merge sort

How many operations?



Roughly $c \log n$ steps per recursion level. How many recursion levels?

Merge sort

How many operations?

More generally, we want to estimate the time complexity of an algorithm that, at each step, divides the problem into 2, solves the subproblems and then it takes linear time to merge the solutions of the two subproblems.

Without worrying about integer divisions, this means: $T(N) = 2 T(N/2) + N$.

Here we have to use a trick: ignoring that we deal with discrete steps, and writing N as 2^n we have

$$T(2^n) = 2 * T(2^{n-1}) + 2^n$$

Dividing both sides by 2^n , we obtain

$$T(2^n)/2^n = T(2^{n-1})/2^{n-1} + 1 \text{ and telescoping further we get}$$

$$\begin{aligned} T(2^n)/2^n &= T(2^{n-1})/2^{n-1} + 1 = (2 * T(2^{n-2}) + 2^{n-1}) / 2^{n-1} + 1 \\ &= T(2^{n-2})/2^{n-2} + 1 + 1 \end{aligned}$$

...

$$= 1 + 1 + \dots + 1 \text{ (n times)}$$

Thus $T(2^n)/2^n = n$, and if we make the reverse substitution, $2^n = N$, we get $T(N) = N * \log_2(N)$.

This is the complexity of the merge sort algorithm.

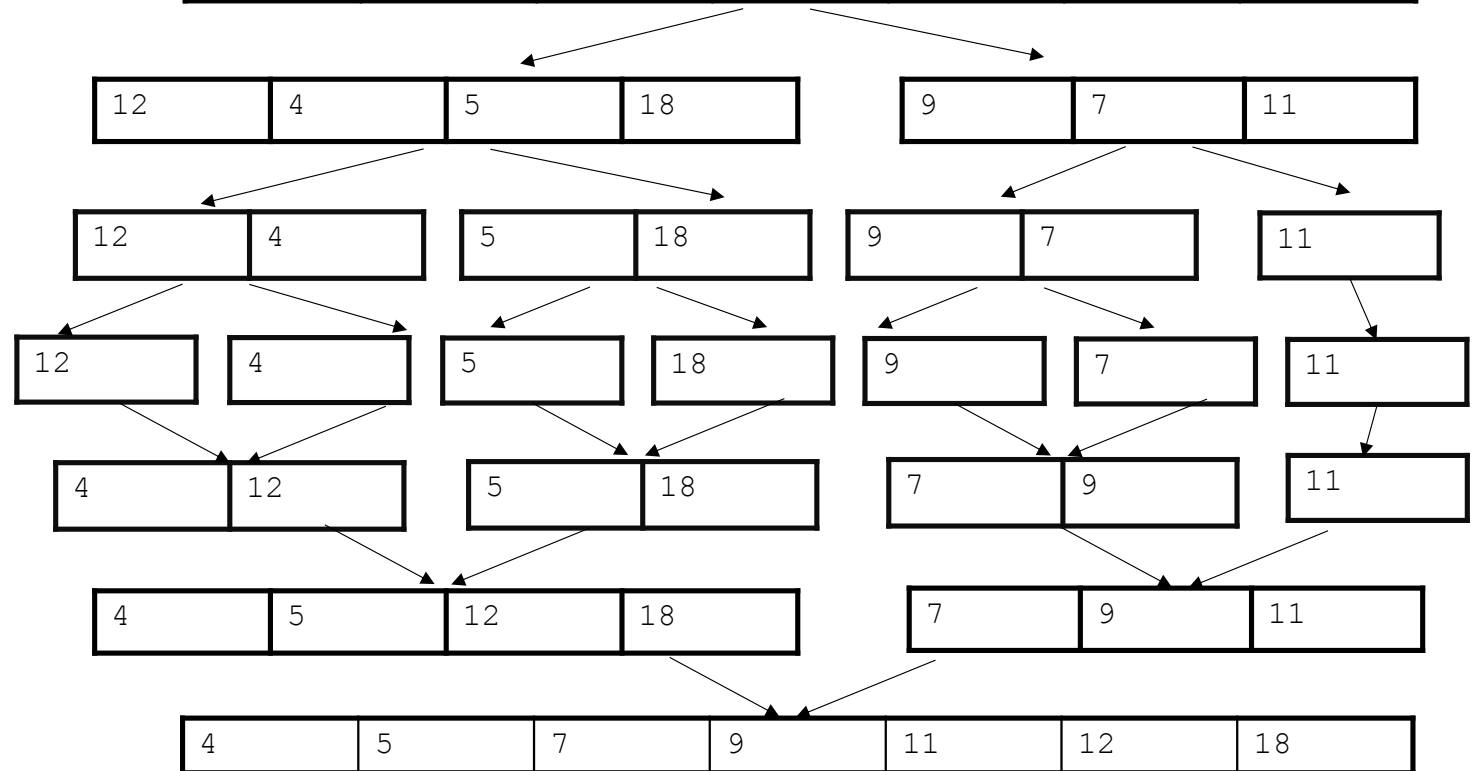
Trading space for time

12	4	5	18	9	7	11
----	---	---	----	---	---	----

Bubble sort – in place

12	4	5	18	9	7	11
----	---	---	----	---	---	----


Merge sort – multiple copies of the list



Database searching

NCBI's 'Nucleotide' repository of sequence information (<https://www.ncbi.nlm.nih.gov/genbank/statistics/>)

Release	Date	Bases (GenBank)	Sequences	Bases (WGS)	Sequences
262	Aug 2024	3'675'462'701'077	251'998'350	29'643'594'176'326	3'569'715'357

 NCBI [Resources](#) [How To](#) [Sign in to NCBI](#)

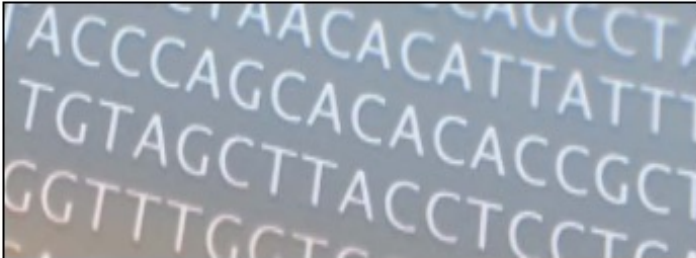
Nucleotide

Nucleotide

Advanced

Search

Help



Nucleotide

The Nucleotide database is a collection of sequences from several sources, including GenBank, RefSeq, TPA and PDB. Genome, gene and transcript sequence data provide the foundation for biomedical research and discovery.

Using Nucleotide

- [Quick Start Guide](#)
- [FAQ](#)
- [Help](#)
- [GenBank FTP](#)
- [RefSeq FTP](#)

Nucleotide Tools

- [Submit to GenBank](#)
- [LinkOut](#)
- [E-Utilities](#)
- [BLAST](#)
- [Batch Entrez](#)

Other Resources

- [GenBank Home](#)
- [RefSeq Home](#)
- [Gene Home](#)
- [SRA Home](#)
- [INSDC](#)

How can we search this large database efficiently?

Upper bound: linear traversal of the sequences in the DB, at each position looking for a perfect match

DB organization



Time to find the perfect match: logarithmic in the size of the database

Lower bound: constant time identification of sequence-positions where a perfect match is to be found

Binary search

Let's say that we have a list of objects that is sorted by some measure.

We are given a new object and we want to find out whether it is already present in the list.

Binary search

A concrete example: we have a list of ions that can be obtained in a mass spectrometry experiment, sorted by their mass/charge ratio.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

We are given a new m/z (say 31) and we want to find out whether it (and its associated ion) is already in the list.

The way we will do it is by repeatedly halving the interval in which the number can possibly be.

Binary search

```
def binarySearch (num, numList):  
    begin = 0  
    end = len(numList)-1  
    while(end > begin):  
        mid = int((begin+end)/2)  
        if(num == numList[mid]):  
            return True  
        elif(num > numList[mid]):  
            begin = mid + 1  
        elif (num < numList[mid]):  
            end = mid  
    if(end == begin):  
        if(num == numList[end]):  
            return True  
    return False
```

Binary search

Let's take a simple example: we have a list of ions that can be obtained in a mass spectrometry experiment, sorted by their mass/charge ratio.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

We are given a new m/z (say 31) and we want to find out whether it (and its associated ion) is already in the list.

Binary search

Let's take a simple example: we have a list of ions that can be obtained in a mass spectrometry experiment, sorted by their mass/charge ratio.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

We are given a new m/z (say 31) and we want to find out whether it (and its associated ion) is already in the list.

Binary search

Let's take a simple example: we have a list of ions that can be obtained in a mass spectrometry experiment, sorted by their mass/charge ratio.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

We are given a new m/z (say 31) and we want to find out whether it (and its associated ion) is already in the list.

Binary search

Let's take a simple example: we have a list of ions that can be obtained in a mass spectrometry experiment, sorted by their mass/charge ratio.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

We are given a new m/z (say 31) and we want to find out whether it (and its associated ion) is already in the list.

Binary search

Let's take a simple example: we have a list of ions that can be obtained in a mass spectrometry experiment, sorted by their mass/charge ratio.

<i>i</i>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

We are given a new m/z (say 38) and we want to find out whether it (and its associated ion) is already in the list.

Binary search

How does the running time of binary search scale with the size of the list?

Analogous to merge sort, at every iteration we reduce the size of the interval in which we have to look.

$$T(N) = T\left(\frac{N}{2}\right) + c = T\left(\frac{N}{4}\right) + c + c = \dots = T(1) + c \dots + c$$

And with the same substitution trick

$$T(2^n) = T(2^{n-1}) + c = T(2^{n-2}) + 2c = \dots = nc$$

$$T(N) = c \log_2(N)$$

Trees

Tree data structures support searches in $O(\log(n))$ time, n being the number of items in the tree.

They exploit some underlying order in the data.

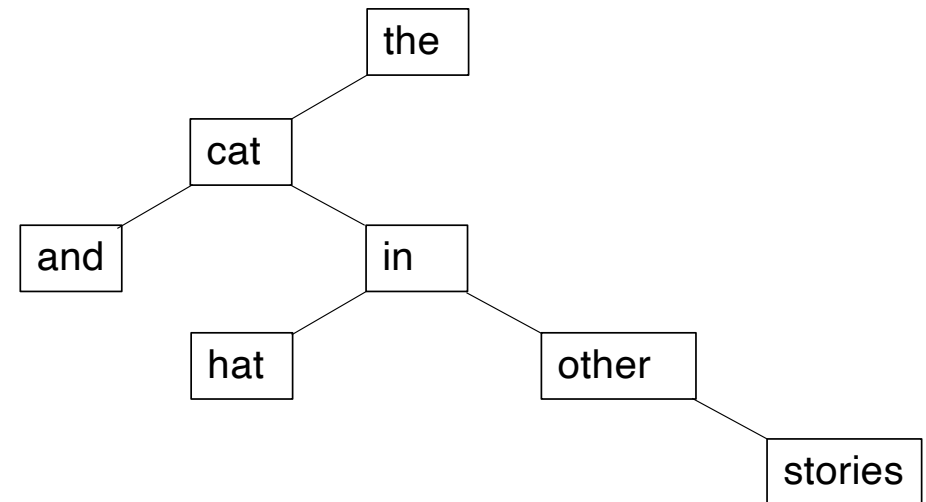
For example, let's say that we have some message, composed of words (could be the genome and the genes in it) and we want to find out if arbitrary input words occur in that message.

We organize the message in a binary search tree, which satisfies the condition that all nodes in the left subtree rooted at a specific node have 'smaller' key values than the key of the root, and all nodes in the right subtree have higher key values.

Note: here we compare key values lexicographically.

Difference with respect to array?

the cat in the hat and other stories



Importance of tree structure

Given the list

`L = ['the', 'cat', 'in', 'the', 'hat', 'and', 'other', 'stories']`

Construct the binary tree of the words that occur in the list, inserting all elements of the list into the binary tree, one at a time, from left-most to right-most.

What does the tree that you get look like?

How many steps do you need, on average, to find one of these words in the constructed binary tree?

Now do the same for the list

`L' = ['and', 'cat', 'hat', 'in', 'other', 'stories', 'the']`

What do you conclude?

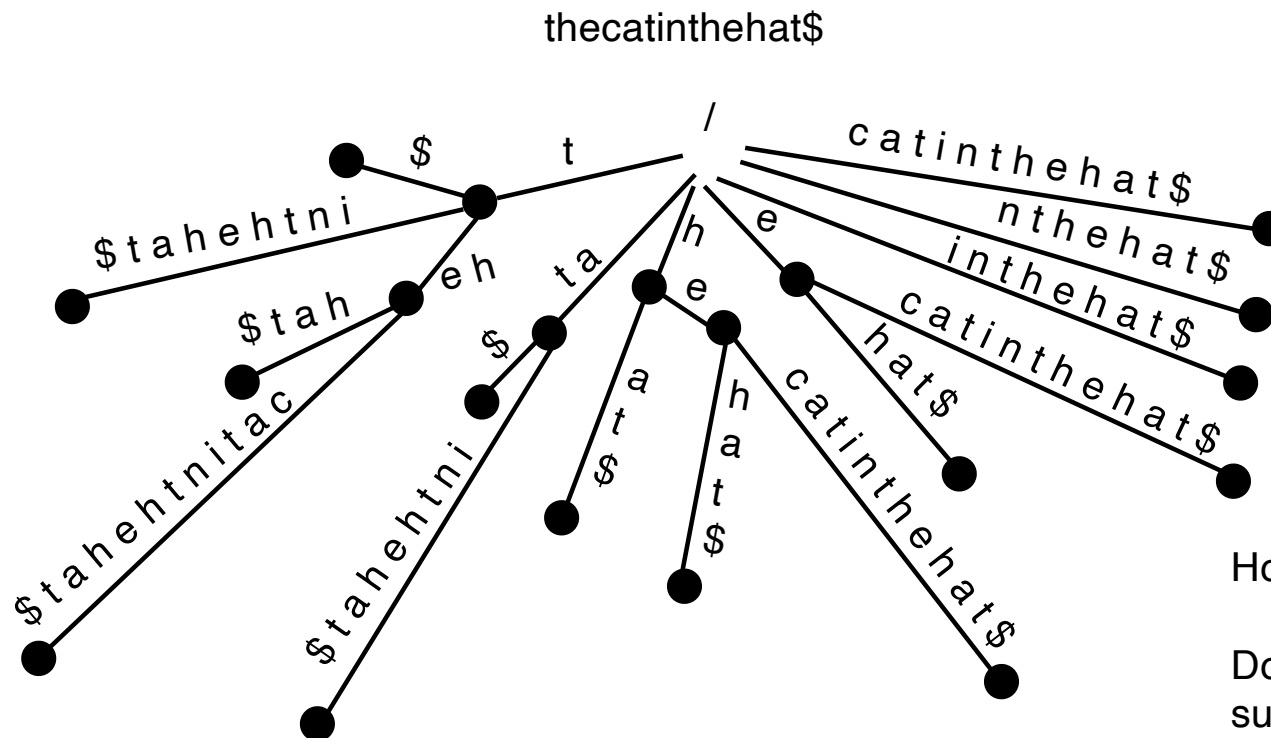
Estimating time/memory requirements

1. Looping through a list – searching, copying, passing list as argument to function
2. Nested for loops – traversing all elements of a matrix, searching multiple elements in a database
3. Order of magnitude calculation – determine the number of digits of a number, search in a n -ary tree

Suffix trees

For large sequence data such as genomes, a special kind of tree structure is frequently used, the suffix tree.

This exploits more of the structure in the underlying data (keeping track of the sequence in which letters occur).



How would you prove this?

Do you see any catch in having such an efficient search?

Search time with a suffix tree is linear in the query length.

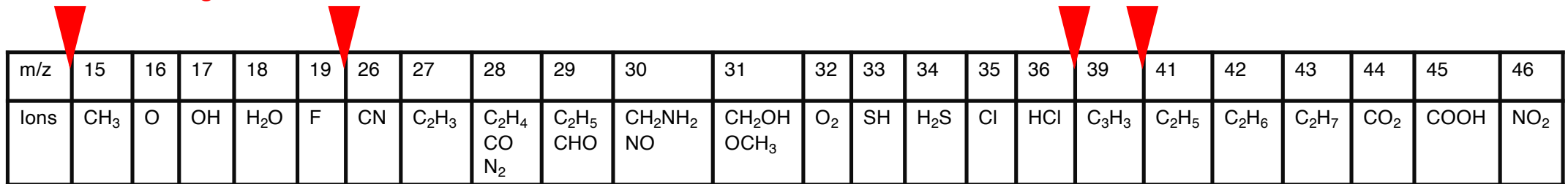
Hashing

Supports search operations in constant time, $O(1)$.

Idea: store elements in a table, at indices that can be computed directly from the elements in constant time.

An example we have seen before:

Missing indices



The diagram shows a hash table with 22 slots, indexed from 15 to 46. Red triangles point to indices 15, 19, 39, and 41, which are labeled as 'Missing indices' in red text. The table contains the following data:

m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

Hashing

Supports search operations in constant time, $O(1)$.

Idea: store elements in a table, at indices that can be computed directly from the elements in constant time.

An example we have seen before:

m/z	15	16	17	18	19	26	27	28	29	30	31	32	33	34	35	36	39	41	42	43	44	45	46
Ions	CH ₃	O	OH	H ₂ O	F	CN	C ₂ H ₃	C ₂ H ₄ CO N ₂	C ₂ H ₅ CHO	CH ₂ NH ₂ NO	CH ₂ OH OCH ₃	O ₂	SH	H ₂ S	Cl	HCl	C ₃ H ₃	C ₂ H ₅	C ₂ H ₆	C ₂ H ₇	CO ₂	COOH	NO ₂

Difference from the sorted list: we could not drop indices. The hash table would have 46 entries, only 23 populated with values, the others empty (or 23 populated with the value True and the rest with the value False).

Main challenges with hash tables: ensure that they are
densely populated
without clashes (multiple elements hashing to the same index).

Applications of hashing

Dictionaries: indexing by strings

```
#string argument has to be a sequence of bytes
#b'stringValue' #default representation utf-8
def hash(string):
    x = string[0] << 7
    for chr in string[1:]:
        x = ((1000003 * x) ^ chr) & (1<<32)
    return x
```

Applications of hashing

Data transfer, computer security etc.

SHA – Secure Hash Algorithm

Name	Input (message) size (bits)	Output size (bits)
MD5	Unlimited	128
SHA-1	$2^{64}-1$	160
SHA-2	$2^{64}-1/2^{128}-1$	224-512 (depending on variant)
SHA-3	Unlimited	224-arbitrary (depending on variant)

Standard data structures and their methods

1. Queue
2. Stack
3. Tree
4. Heap
5. Graphs

- What operations (methods) are defined on these data structures?
- How do they scale with the size of the data structure?
- What are shallow and deep copies?
- What are stable vs unstable algorithms?

Algorithms in bioinformatics

- **Global Alignment :**
 - Needleman-Wunsch
- **Local Alignment :**
 - Smith-Waterman
 - BLAST
- **Multiple Sequence Alignment:**
 - ProbCons
- **Motif finding**
 - GibbsSampler
 - MEME
- **Assembly :**
 - Debrujin Graph Assembly
 - Euler Assembly
- **Read Mapping :**
 - Suffix Trees
 - Burrows Wheeler Alignment

Improving efficiency – Burroughs-Wheeler transform

Suffix trees take a lot of space to save time. We can construct a much more space-efficient data structure.

```
thecatinthehat$
012345678911111
      01234
```

Make all circular permutations

Sort lexicographically

Lexicographically-
sorted suffices

BWT

0	thecatinthehat\$	14	\$thecatinthehat	14	\$ thecatinthehat	14	\$thecatinthehat
14	\$thecatinthehat	12	at\$thecatintheh	12	at \$ thecatintheh	12	at\$thecatinthe h
13	t\$thecatintheha	4	atinthehat\$thec	4	atinthehat \$ thec	4	atinthehat\$thec c
12	at\$thecatintheh	3	catinthehat\$the	3	catinthehat \$ the	3	catinthehat\$the e
11	hat\$thecatinthe	2	ecatinthehat\$th	2	ecatinthehat \$ th	2	ecatinthehat\$th h
10	ehat\$thecatint	10	ehat\$thecatint	10	ehat \$ thecatint	10	ehat\$thecatint h
9	hehat\$thecatint	11	hat\$thecatint	11	hat \$ thecatint	11	hat\$thecatint e
8	thehat\$thecatin	1	hecatinthehat\$t	1	hecatinthehat \$ t	1	hecatinthehat\$t t
7	nthehat\$thecati	9	hehat\$thecatint	9	hehat \$ thecatint	9	hehat\$thecatint t
6	inthehat\$thecat	6	inthehat\$thecat	6	inthehat \$ thecat	6	inthehat\$thecat t
5	tinthehat\$theca	7	nthehat\$thecati	7	nthehat \$ thecati	7	nthehat\$thecati i
4	atinthehat\$thec	13	t\$thecatintheha	13	t \$ thecatintheha	13	t\$thecatintheha a
3	catinthehat\$the	0	thecatinthehat\$	0	thecatinthehat \$	0	thecatinthehat\$ \$
2	ecatinthehat\$th	8	thehat\$thecatin	8	thehat \$ thecatin	8	thehat\$thecatin n
1	hecatinthehat\$t	5	tinthehat\$theca	5	tinthehat \$ theca	5	tinthehat\$theca a

Improving efficiency – Burroughs-Wheeler transform

		$t_1 h_1 e_1 c_1 a_1 t_2 i_1 n_1 t_3 h_2 e_2 h_3 a_2 t_4 \$$ 0 1 2 3 4 5 6 7 8 9 1 1 1 1 1	$t_2 h_2 e_1 c_1 a_2 t_4 i_1 n_1 t_3 h_3 e_2 h_1 a_1 t_1 \$$ 0 1 2 3 4 5 6 7 8 9 1 1 1 1 1
BWT	F	L	
14 \$thecatinthehat t	\$thecatinthehat t	14 \$ thecatintheha t₄	14 \$ thecatintheha t₁
12 at\$thecatinthe h	at\$thecatinthe h	12 a₂ t\$thecatinthe h₃	12 a₁ t\$thecatinthe h₁
4 atinthehat\$the c	atinthehat\$the c	4 a₁ tinthehat\$the c₁	4 a₂ tinthehat\$the c₁
3 catinthehat\$the e	catinthehat\$the e	3 c₁ atinthehat\$the e₁	3 c₁ atinthehat\$the e₁
2 ecatinthehat\$th h	ecatinthehat\$th h	2 e₁ catinthehat\$th h₁	2 e₁ catinthehat\$th h₂
10 ehat\$thecatinth h	ehat\$thecatinth h	10 e₂ hat\$thecatinth h₂	10 e₂ hat\$thecatinth h₃
11 hat\$thecatinthe e	hat\$thecatinthe e	11 h₃ at\$thecatinthe e₂	11 h₁ at\$thecatinthe e₂
1 hecatinthehat\$t t	hecatinthehat\$t t	1 h₁ ecatinthehat\$t t₁	1 h₂ ecatinthehat\$t t₂
9 hehat\$thecatint t	hehat\$thecatint t	9 h₂ ehat\$thecatint t₃	9 h₃ ehat\$thecatint t₃
6 inthehat\$thecat t	inthehat\$thecat t	6 i₁ nthehat\$thecat t₂	6 i₁ nthehat\$thecat t₄
7 nthehat\$thecati i	nthehat\$thecati i	7 n₁ thehat\$thecati i₁	7 n₁ thehat\$thecati i₁
13 t\$thecatintheha a	t\$thecatintheha a	13 t₄ \$thecatintheha a₂	13 t₁ \$thecatintheha a₁
0 thecatinthehat\$ \$	thecatinthehat\$ \$	0 t₁ hecatinthehat\$	0 t₂ hecatinthehat\$
8 thehat\$thecatin n	thecat\$thecatin n	8 t₃ hehat\$thecatin n₁	8 t₃ hehat\$thecatin n₁
5 tinthehat\$theca a	tinthehat\$theca a	5 t₂ inthehat\$theca a₂	5 t₄ inthehat\$theca a₂
Can be stored efficiently: \$1a2c1e2h3i1n1t4	sorted	preceding char	Order of character occurrence is preserved in L column
			Relabel by rank in F

Which row begins with the second e-starting suffix?

0 \$ thecatinthehat**t**₁
1 **a**₁t\$thecatinthe**h**₁
2 **a**₂tinthehat\$the**c**₁
3 **c**₁atinthehat\$the**e**₁
4 **e**₁catinthehat\$th**h**₂
5 **e**₂hat\$thecatinth**h**₃
6 **h**₁at\$thecatinthe**e**₂
7 **h**₂ecatinthehat\$**t**₂
8 **h**₃ehat\$thecatin**t**₃
9 **i**₁inthehat\$thecat**t**₄
10 **n**₁thehat\$thecati**i**₁
11 **t**₁\$thecatintheha**a**₁
12 **t**₂hecatinthehat\$**s**₂
13 **t**₃hehat\$thecatin**n**₁
14 **t**₄inthehat\$theca**a**₂

If we know the letter count in the text

a : 2

c : 1

e : 2

h : 3

i : 1

n : 1

t : 4

The index of the row where e_2 should be is $2(a) + 1(c) + 2(e) = 5$

Can we reconstruct the text if we only have the F and L columns?

We'll do it from the end of the text towards the beginning, knowing that

- the character in L **precedes** the one in F and
- the rank-order of characters in first and last column is **preserved**.

F	L
\$	thecatinthehat t ₁
a ₁	t\$thecatinthe h ₁
a ₂	tinthehat\$the c ₁
c ₁	atinthehat\$the e ₁
e ₁	catinthehat\$th h ₂
e ₂	hat\$thecatint h ₃
h ₁	at\$thecatint e ₂
h ₂	ecatinthehat\$t t ₂
h ₃	ehat\$thecatint t ₃
i ₁	nthehat\$thecat t ₄
n ₁	thehat\$thecati i ₁
t ₁	\$thecatinthe a ₁
t ₂	hecatinthehat\$ \$ ₂
t ₃	hehat\$thecatin n ₁
t ₄	inthehat\$theca a ₂

F	L
\$	t ₁
a ₁	h ₁
a ₂	c ₁
c ₁	e ₁
e ₁	h ₂
e ₂	h ₃
h ₁	e ₂
h ₂	t ₂
h ₃	t ₃
i ₁	t ₄
n ₁	i ₁
t ₁	a ₁
t ₂	\$ ₂
t ₃	n ₁
t ₄	a ₂

\$	
t ₁ \$	preceding?
t ₁ \$	row of first t?
a ₁ t ₁ \$	preceding?
a ₁ t ₁ \$	row of first a?
h ₁ a ₁ t ₁ \$	preceding?
h ₁ a ₁ t ₁ \$	row of first h?
e ₂ h ₁ a ₁ t ₁ \$	preceding?
e ₂ h ₁ a ₁ t ₁ \$	row of second e?

•
 •
 •
 t₂h₂e₁c₁a₂t₄i₁n₁t₃h₃e₂h₁a₁t₁\$

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthehat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁	h ₂
e ₂ hat\$thecatinth h ₃		e ₂	h ₃
h ₁ at\$thecatinthe e ₂		h ₁	e ₂
h ₂ ecatinthehat\$t t ₂		h ₂	t ₂
h ₃ ehat\$thecatint t ₃		h ₃	t ₃
i ₁ nthehat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecati i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatinthehat\$ \$ ₂		t ₂	\$
t ₃ hehat\$thecatin n ₁		t ₃	n ₁
t ₄ inthehat\$theca a ₂		t ₄	a ₂

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

Which of these **e**'s is preceded by an **h**?

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthehat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁	h ₂
e ₂ hat\$thecatinth h ₃		e ₂	h ₃
h ₁ at\$thecatinthe e ₂		h ₁	e ₂
h ₂ ecatinthehat\$ t ₂		h ₂	t ₂
h ₃ ehat\$thecatint t ₃		h ₃	t ₃
i ₁ nthehat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecati i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatinthehat\$ s ₂		t ₂	s
t ₃ hehat\$thecatin n ₁		t ₃	n ₁
t ₄ inthehat\$theca a ₂		t ₄	a ₂

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

Which of these **e**'s is preceded by an **h**?

Locate **h₂** and **h₃** in F

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthehat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁ →	h ₂
e ₂ hat\$thecatinth h ₃		e ₂ →	h ₃
h ₁ at\$thecatinthe e ₂		h ₁	e ₂
h ₂ ecatintthehat\$ t ₂		h ₂	t ₂
h ₃ ehat\$thecatint t ₃		h ₃	t ₃
i ₁ ntthehat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecati i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatintthehat\$ \$ ₂		t ₂	\$
t ₃ hehat\$thecatin n ₁		t ₃	n ₁
t ₄ intthehat\$theca a ₂		t ₄	a ₂

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

Which of these **e**'s is preceded by an **h**?

Locate **h₂** and **h₃** in F

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthecat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁ →	h ₂
e ₂ hat\$thecatint h ₃		e ₂ →	h ₃
h ₁ at\$thecatint e ₂		h ₁	e ₂
h ₂ ecatinthehat\$ t ₂		h ₂	t ₂
h ₃ ehat\$thecatint t ₃		h ₃	t ₃
i ₁ nthecat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecat i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatinthehat\$ s ₂		t ₂	s
t ₃ hehat\$thecatint n ₁		t ₃	n ₁
t ₄ inthehat\$theca a ₂		t ₄	a ₂

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

Which of these **e**'s is preceded by an **h**?

Locate **h₂** and **h₃** in F

Which of these **h**'s is preceded by a **t**?

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthecat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁ →	h ₂
e ₂ hat\$thecatint h ₃		e ₂ →	h ₃
h ₁ at\$thecatint e ₂		h ₁	e ₂
h ₂ ecatinthehat\$ t ₂		h ₂	t ₂
h ₃ ehat\$thecatint t ₃		h ₃	t ₃
i ₁ nthecat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecat i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatinthehat\$ s ₂		t ₂	s
t ₃ hehat\$thecatint n ₁		t ₃	n ₁
t ₄ inthehat\$theca a ₂		t ₄	a ₂

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

Which of these **e**'s is preceded by an **h**?

Locate **h₂** and **h₃** in F

Which of these **h**'s is preceded by a **t**?

Locate **t₂** and **t₃** in F

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthehat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁ →	h ₂
e ₂ hat\$thecatinth h ₃		e ₂ →	h ₃
h ₁ at\$thecatinthe e ₂		h ₁	e ₂
h ₂ ecatintthehat\$t t ₂		h ₂ →	t ₂
h ₃ ehat\$thecatin t ₃		h ₃ →	t ₃
i ₁ nthecat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecat i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatintthehat\$ \$ ₂		t ₂	\$ ₂
t ₃ hehat\$thecat n ₁		t ₃	n ₁
t ₄ intthehat\$theca a ₂		t ₄	a ₂

Does the substring **the** occur in the text?

We'll look for it, again starting with its last letter.

Which row starts with **e**, the last letter of our search pattern?

Which of these **e**'s is preceded by an **h**?

Locate **h**₂ and **h**₃ in F

Which of these **h**'s is preceded by a **t**?

Locate **t**₂ and **t**₃ in F

The substring **the** occurs in two locations in the text

F	L	F	L
\$ thecatinthehat t ₁		\$	t ₁
a ₁ t\$thecatinthe h ₁		a ₁	h ₁
a ₂ tinthehat\$the c ₁		a ₂	c ₁
c ₁ atinthehat\$the e ₁		c ₁	e ₁
e ₁ catinthehat\$th h ₂		e ₁ →	h ₂
e ₂ hat\$thecatinth h ₃		e ₂ →	h ₃
h ₁ at\$thecatinthe e ₂		h ₁	e ₂
h ₂ ecatintthehat\$ t ₂		h ₂ →	t ₂
h ₃ ehat\$thecatin t ₃		h ₃ →	t ₃
i ₁ nthecat\$thecat t ₄		i ₁	t ₄
n ₁ thehat\$thecat i ₁		n ₁	i ₁
t ₁ \$thecatinthe a ₁		t ₁	a ₁
t ₂ hecatintthehat\$ s ₂		t ₂	s
t ₃ hehat\$thecat n ₁		t ₃	n ₁
t ₄ intthehat\$theca a ₂		t ₄	a ₂