

Fundamentals–LINEARSEARCH and Binary Search

COMPSCI 2CO3: Data Structures and Algorithms

Jelle Hellings

Department of Computing and Software
McMaster University



Fall 2024

Reflection on CONTAINS

Is CONTAINS a *good* algorithm?

CONTAINS is correct and has a runtime complexity of $\Theta(|L|)$ \longrightarrow Sounds good to me!

Reflection on CONTAINS

Is CONTAINS a *good* algorithm?

CONTAINS is correct and has a runtime complexity of $\Theta(|L|)$ \longrightarrow Sounds good to me!

Critique: CONTAINS is *too specialized* \longrightarrow .

We cannot use CONTAINS for anything else than the contains problem!

Example

- ▶ Searching in only *part* of the list?
- ▶ Finding where v is in the list?

Reflection on CONTAINS

Critique: CONTAINS is *too specialized* \longrightarrow .

We cannot use CONTAINS for anything else than the contains problem!

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: " $o \leq r \leq |L|$ and $v \notin L[o, r)$ ", bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

Reflection on CONTAINS

Critique: CONTAINS is *too specialized* \longrightarrow .

We cannot use CONTAINS for anything else than the contains problem!

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

Algorithm LSCONTAINS(L, v):

1: **return** LINEARSEARCH($L, v, 0$) $\neq |L|$.

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: " $o \leq r \leq |L|$ and $v \notin L[o, r)$ ", bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

- With respect to worst case inputs ($v \notin L$): $\Theta(|L|)$.

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

- ▶ With respect to worst case inputs ($v \notin L$): $\Theta(|L|)$.
- ▶ With respect to best case inputs ($v = L[o]$): $\Theta(1)$.

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

- ▶ With respect to worst case inputs ($v \notin L$): $\Theta(|L|)$.
- ▶ With respect to best case inputs ($v = L[o]$): $\Theta(1)$.

Problem: Modeling runtime complexity in terms of *only the input* limits us!

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

Problem: Modeling runtime complexity in terms of *only the input* limits us!

Assume: $L[i] = v$ and i is the *first* offset after o equivalent to v .

The runtime complexity of LINEARSEARCH is $\Theta(i - o)$.

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

Problem: Modeling runtime complexity in terms of *only the input* limits us!

Assume: $L[i] = v$ and i is the *first* offset after o equivalent to v .

The runtime complexity of LINEARSEARCH is $\Theta(i - o)$ with $i = \text{LINEARSEARCH}(L, v, o)$.

Problem: What if we search often?

`LINEARSEARCH(L, v, o)` can read all of array L : *potentially-high cost*.

Can we do better?

Problem: What if we search often?

LINEARSEARCH(L, v, o) can read all of array L : *potentially-high cost*.

Can we do better?

No: We do not know anything about L to help us!

→ we have to look at all elements in L .

Problem: What if we search often?

LINEARSEARCH(L, v, o) can read all of array L : *potentially-high cost*.

Can we do better?

No: We do not know anything about L to help us!

→ we have to look at all elements in L .

Maybe: If we know more about L .

An example of a list

Consider a list *enrolled* with schema

enrolled(*dept*, *code*, *sid*, *date*)

that models a list of all students enrolled for a course.

What if...

We add enrollment data to *the end of the list*.

Question: What do we know about enrolled?

An example of a list

Consider a list *enrolled* with schema

enrolled(*dept*, *code*, *sid*, *date*)

that models a list of all students enrolled for a course.

What if...

We add enrollment data to *the end of the list*.

Question: What do we know about enrolled? → enrolled is *ordered* on *date*!

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

$L[i] < v$ As the list L is ordered, every value in $L[0, i]$ is smaller than v .

→ $v \in L$ if and only if $v \in L[i + 1, |L|)$.

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

$L[i] < v$ As the list L is ordered, every value in $L[0, i]$ is smaller than v .

→ $v \in L$ if and only if $v \in L[i+1, |L|)$.

$L[i] > v$ As the list L is ordered, every value in $L[i, |L|)$ is larger than v .

→ $v \in L$ if and only if $v \in L[0, i)$.

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

$L[i] < v$ As the list L is ordered, every value in $L[0, i]$ is smaller than v .

→ $v \in L$ if and only if $v \in L[i+1, |L|)$.

$L[i] > v$ As the list L is ordered, every value in $L[i, |L|)$ is larger than v .

→ $v \in L$ if and only if $v \in L[0, i)$.

$L[i] = v$ We found v !

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

$L[i] < v$ As the list L is ordered, every value in $L[0, i]$ is smaller than v .

→ $v \in L$ if and only if $v \in L[i+1, |L|)$.

$L[i] > v$ As the list L is ordered, every value in $L[i, |L|)$ is larger than v .

→ $v \in L$ if and only if $v \in L[0, i)$.

$L[i] = v$ We found v !

One comparison can remove a large portion of the array.

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

$L[i] < v$ As the list L is ordered, every value in $L[0, i]$ is smaller than v .

→ $v \in L$ if and only if $v \in L[i + 1, |L|)$.

$L[i] > v$ As the list L is ordered, every value in $L[i, |L|)$ is larger than v .

→ $v \in L$ if and only if $v \in L[0, i)$.

$L[i] = v$ We found v !

One comparison can remove a large portion of the array.

Binary Search: *Maximize potential* by comparing v with the middle of L .

The *recursive* Binary Search algorithm

Algorithm LOWERBOUNDREC($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

- Is LOWERBOUNDREC correct?
- What is the runtime and memory complexity of LOWERBOUNDREC?

Correctness of LOWERBOUNDREC

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```


Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Induction Hypothesis

For any L' , v' , and $0 \leq begin' \leq end' \leq |L'|$ with $0 \leq end' - begin' < m$,
LOWERBOUNDREC($L', v', begin', end'$) returns the correct result.

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

1: **if** $begin = end$ **then**

2: **return** $begin$.

3: **else**

4: $mid := (begin + end) \text{ div } 2$.

5: **if** $L[mid] < v$ **then**

6: **return** LOWERBOUNDREC($L, v, mid + 1, end$).

7: **else** $L[mid] \geq v$

8: **return** LOWERBOUNDREC($L, v, begin, mid$).

} *Base case:*
Inspecting $end - begin = 0$ elements.

} *Recursive case:*
Inspecting $end - begin > 0$ elements.

Induction Hypothesis

For any L' , v' , and $0 \leq begin' \leq end' \leq |L'|$ with $0 \leq end' - begin' < m$,

LOWERBOUNDREC($L', v', begin', end'$) returns the correct result.


Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

$begin \leq mid < end$



Induction Hypothesis

For any L' , v' , and $0 \leq begin' \leq end' \leq |L'|$ with $0 \leq end' - begin' < m$,
LOWERBOUNDREC($L', v', begin', end'$) returns the correct result.

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

$begin \leq mid < end$

Induction Hypothesis

Induction Hypothesis

Induction Hypothesis

For any L' , v' , and $0 \leq begin' \leq end' \leq |L'|$ with $0 \leq end' - begin' < m$,
LOWERBOUNDREC($L', v', begin', end'$) returns the correct result.

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Termination

Bound function: $end - begin$.

Intermezzo: Runtime complexity of LOWERBOUNDREC

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \div 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Complexity of LOWERBOUNDREC with $n = end - begin$

$$T(n) =$$

Intermezzo: Runtime complexity of LOWERBOUNDREC

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

} *Base case:*
1 operation.

Complexity of LOWERBOUNDREC with $n = end - begin$

$$T(n) = \begin{cases} 1 & \text{if } n = 0; \end{cases}$$

Intermezzo: Runtime complexity of LOWERBOUNDREC

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

} *Base case:*
1 operation.

} *Recursive case:*
1 operation and 1 recursive call.

Complexity of LOWERBOUNDREC with $n = end - begin$

$$T(n) = \begin{cases} 1 & \text{if } n = 0; \\ 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1 & \text{if } n \geq 1. \end{cases}$$

Intermezzo: Runtime complexity of LOWERBOUNDREC

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Complexity of LOWERBOUNDREC with $n = end - begin$

$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$

Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$

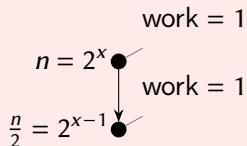
work = 1

$$n = 2^x \bullet$$

Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

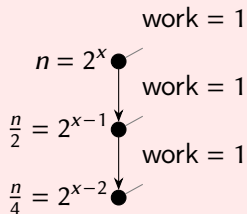
$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$



Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

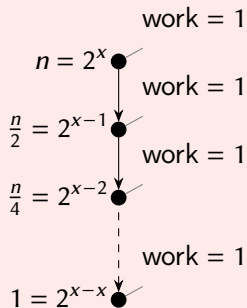
$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$



Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

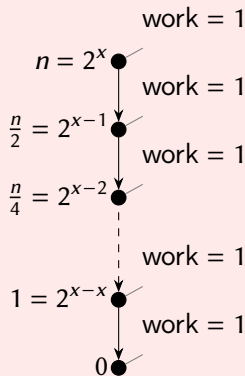
$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$



Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

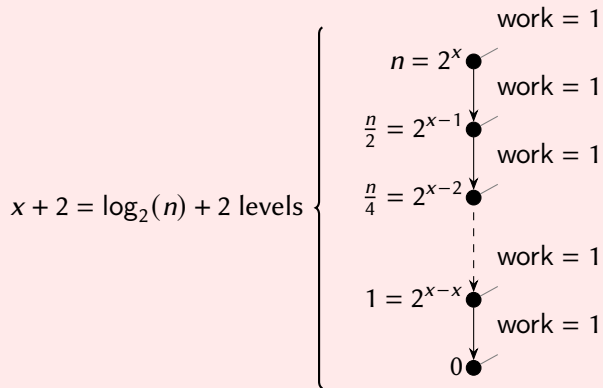
$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$



Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

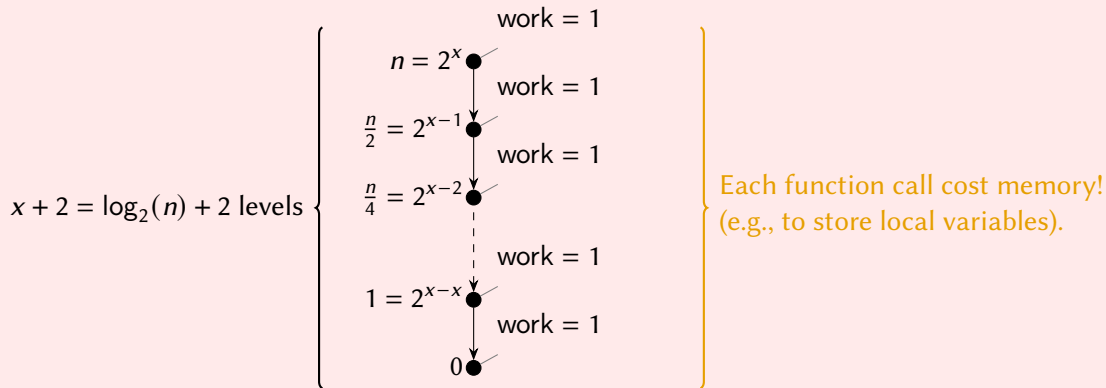
$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1.$$



Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1 = \Theta(\log_2(n)).$$



The *recursive* Binary Search algorithm

Algorithm LOWERBOUNDREC($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

Theorem

LOWERBOUNDREC is correct and has a runtime and memory complexity of $\Theta(\log_2(|L|))$.

The *non-recursive* Binary Search algorithm

Algorithm LOWERBOUND($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

```
1: while  $begin \neq end$  do  
2:    $mid := (begin + end) \text{ div } 2$ .  
3:   if  $L[mid] < v$  then  
4:      $begin := mid + 1$ .  
5:   else  
6:      $end := mid$ .  
7: return  $begin$ .
```

Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

The *non-recursive* Binary Search algorithm

Algorithm LOWERBOUND($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

```
1: while  $begin \neq end$  do
2:    $mid := (begin + end) \text{ div } 2$ .
3:   if  $L[mid] < v$  then
4:      $begin := mid + 1$ .
5:   else
6:      $end := mid$ .
7: return  $begin$ .
```

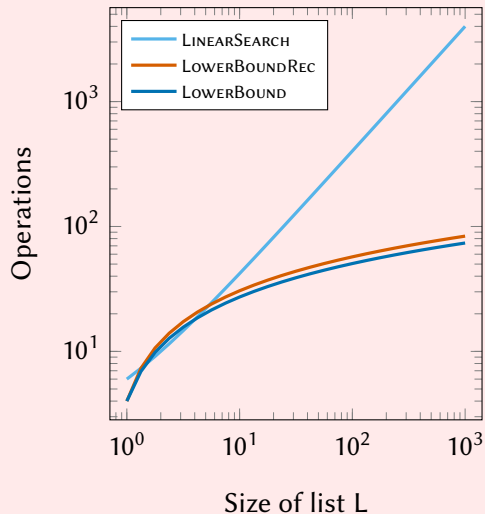
Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

Theorem

LOWERBOUND is correct, has a runtime complexity of $\Theta(\log_2(|L|))$, and a memory complexity of $\Theta(1)$.

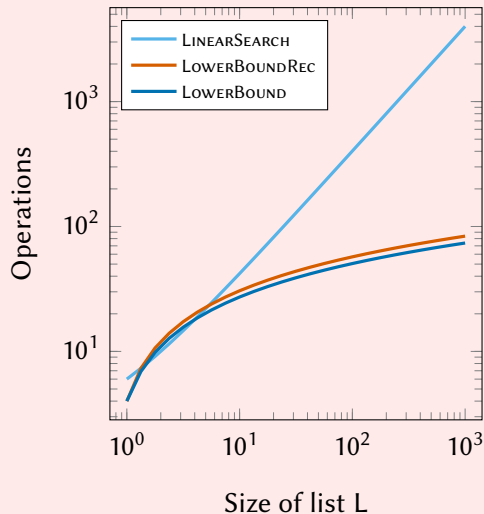
Comparing the complexity of searching

Theoretical complexity

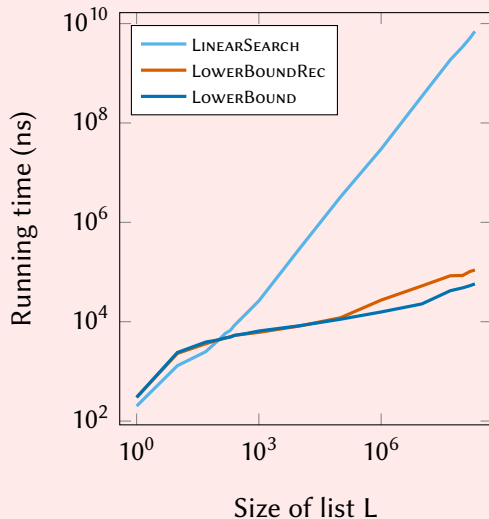


Comparing the complexity of searching

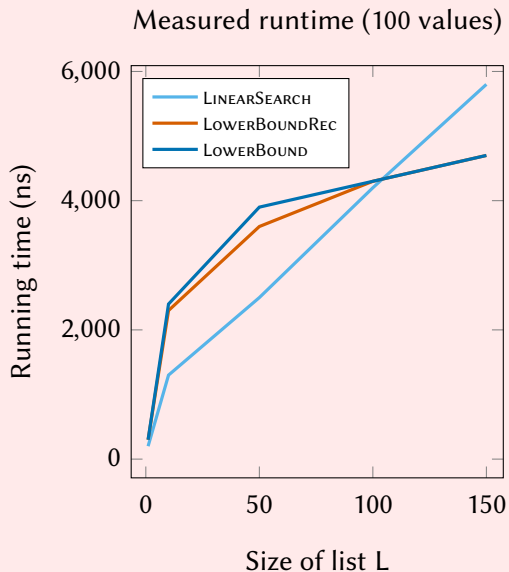
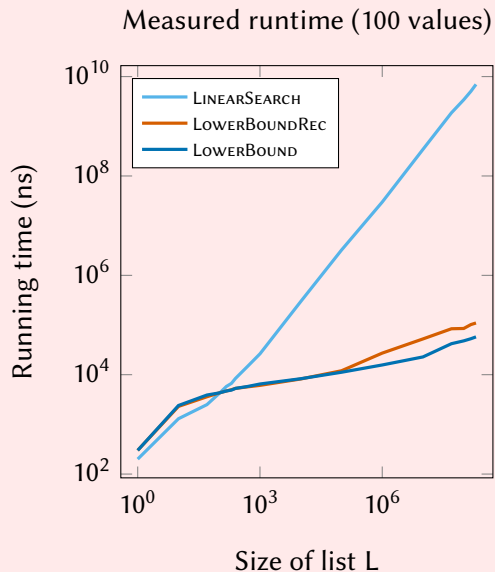
Theoretical complexity



Measured runtime (100 values)



Comparing the complexity of searching



Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Example

Consider a list *enrolled* with schema $\text{enrolled}(\text{dept}, \text{code}, \text{sid}, \text{date})$.

Query: All students enrolled in 2023

Range query on *enrolled* with $[(' ', ' ', -1, 2023), (' ', ' ', -1, 2024)]$.

Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Example

Consider a list *enrolled* with schema $\text{enrolled}(\text{dept}, \text{code}, \text{sid}, \text{date})$.

Query: All students enrolled in 2023

Range query on *enrolled* with $[('', '', -1, 2023), ('', '', -1, 2024)]$.

We add enrollment data to *the end of the list* \rightarrow enrolled is *ordered* on *date*!

Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Algorithm RANGEQUERY($L, [v, w]$):

Input: L is an ordered *array*, v, w are values, and $v \leq w$.

- 1: $i := \text{LOWERBOUND}(L, v, 0, |L|)$.
- 2: $j := i$.
- 3: **while** $j \leq |L|$ **and also** $L[j] \leq w$ **do**
- 4: $j := j + 1$.
- 5: **return** $L[i, j)$.

Result: return the list $L[m, n)$, $0 \leq m \leq n \leq |L|$, such that $L[m, n)$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Algorithm RANGEQUERY($L, [v, w]$):

Input: L is an ordered *array*, v, w are values, and $v \leq w$.

- 1: $i := \text{LOWERBOUND}(L, v, 0, |L|)$.
- 2: $j := i$.
- 3: **while** $j \leq |L|$ **and also** $L[j] \leq w$ **do**
- 4: $j := j + 1$.
- 5: **return** $L[i, j)$.

Result: return the list $L[m, n)$, $0 \leq m \leq n \leq |L|$, such that $L[m, n)$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Theorem

RANGEQUERY is correct and has worst case runtime complexity $\Theta(|L|)$.

Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Algorithm RANGEQUERY($L, [v, w]$):

Input: L is an ordered *array*, v, w are values, and $v \leq w$.

- 1: $i := \text{LOWERBOUND}(L, v, 0, |L|)$.
- 2: $j := i$.
- 3: **while** $j \leq |L|$ **and also** $L[j] \leq w$ **do**
- 4: $j := j + 1$.
- 5: **return** $L[i, j)$.

Result: return the list $L[m, n)$, $0 \leq m \leq n \leq |L|$, such that $L[m, n)$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Theorem

RANGEQUERY is correct and has *all case* runtime complexity $\Theta(\log_2(|L|) + |\text{result}|)$.

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Example

'apple'	'pear'	'orange'
---------	--------	----------

$INSPECT(L, 0) = \text{true}$

$INSPECT(L, 2) = \text{true}$

$INSPECT(L, 1) = \text{true}$

$INSPECT(L, 3) = \text{false}.$

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Solving the list-length problem

We have ordered list $0, 1, \dots$ of possible values for $|L|$.

Conclusion of $INSPECT(L, i)$?

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Solving the list-length problem

We have ordered list $0, 1, \dots$ of possible values for $|L|$.

Conclusion of $INSPECT(L, i)$?

$INSPECT(L, i) = \text{true}$ $|L| > i$ (list L has more than i values).

$INSPECT(L, i) = \text{false}$ $|L| \leq i$ (list L has at-most i values).

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Solving the list-length problem

We have ordered list $0, 1, \dots$ of possible values for $|L|$.

Conclusion of $INSPECT(L, i)$?

$INSPECT(L, i) = \text{true}$ $|L| > i$ (list L has more than i values).

$INSPECT(L, i) = \text{false}$ $|L| \leq i$ (list L has at-most i values).

Issue: no upper bound on the ordered list $0, 1, \dots$ of possible values for $|L|$

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $\text{INSPECT}(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Issue: no upper bound on the ordered list $0, 1, \dots$ of possible values for $|L|$

Guess repeatedly with exponentially-growing guesses.

Algorithm LISTLENGTHUB(L):

Input: L is an *array* of unknown length.

- 1: $n := 1$.
- 2: **while** $\text{INSPECT}(L, n)$ **do**
- 3: $n := 2 \cdot n$.
- 4: **return** n .

Result: return N , $|L| \leq N = 1$ or $|L| \leq N < 2|L|$.

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $\text{INSPECT}(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Algorithm $\text{LBLISTLENGTH}(L, N)$ with $N := \text{LISTLENGTHUB}(L)$:

```
1: begin, end := 0,  $N$ .
2: while begin  $\neq$  end do
3:    $\text{mid} := (\text{begin} + \text{end}) \text{ div } 2$ .
4:   if  $\text{INSPECT}(L, \text{mid})$  then
5:      $\text{begin} := \text{mid} + 1$ .
6:   else
7:      $\text{end} := \text{mid}$ .
8: return begin.
```

Result: return the length $|L|$ of array L .

Optimizing joins using range queries

Definition

A join of two lists L and M results in a list A in which each list value is computed from a combination of values $u \in L$ and $v \in M$ according to some *join condition*.

Example (Return pairs (p, r) of product name p and related category r)

products		categories	
<i>name</i>	<i>category</i>	<i>category</i>	<i>related</i>
Apple	Fruit	Fruit	Food
Bok choy	Vegetable	Fruit	Produce
Canelé	Pastry	Pastry	Food
Donut	Pastry	Vegetable	Food
		Vegetable	Produce

Optimizing joins using range queries

Definition

A join of two lists L and M results in a list A in which each list value is computed from a combination of values $u \in L$ and $v \in M$ according to some *join condition*.

Example (Return pairs (p, r) of product name p and related category r)

products	
<i>name</i>	<i>category</i>
Apple	Fruit
Bok choy	Vegetable
Canelé	Pastry
Donut	Pastry

categories	
<i>category</i>	<i>related</i>
Fruit	Food
Fruit	Produce
Pastry	Food
Vegetable	Food
Vegetable	Produce

Join Result	
<i>name</i>	<i>related</i>
Apple	Food
Apple	Produce
Bok choy	Food
Bok choy	Produce
Canelé	Food
Donut	Food

Optimizing joins using range queries

Algorithm NESTEDLOOPPC(products, categories):

Input: relations products(*name, category*) and categories(*category, related*).

1: *output* := \emptyset .

2: **for** (*p.n, p.c*) \in products **do**

3: **for** (*c.c, c.r*) \in categories **do**

4: **if** *p.c* = *c.c* **then**

5: add (*p.n, c.r*) to *output*.

Result: return $\{(p.n, c.r) \mid ((p.n, p.c) \in \text{products}) \wedge ((c.c, c.r) \in \text{categories})\}$.

Optimizing joins using range queries

Algorithm NESTEDLOOPPC(products, categories):

Input: relations products(*name, category*) and categories(*category, related*).

1: *output* := \emptyset .

2: **for** (*p.n, p.c*) \in products **do**

3: **for** (*c.c, c.r*) \in categories **do**

4: **if** *p.c* = *c.c* **then**

5: add (*p.n, c.r*) to *output*.

} $\Theta(|categories|)$.

Result: return $\{(p.n, c.r) \mid ((p.n, p.c) \in \text{products}) \wedge ((c.c, c.r) \in \text{categories})\}$.

Optimizing joins using range queries

Algorithm NESTEDLOOPPC(products, categories):

Input: relations products(*name, category*) and categories(*category, related*).

1: *output* := \emptyset .

2: **for** (*p.n, p.c*) \in products **do**

3: **for** (*c.c, c.r*) \in categories **do**

4: **if** *p.c* = *c.c* **then**

5: add (*p.n, c.r*) to *output*.

$\left. \begin{array}{l} \text{3: } \mathbf{for} (c.c, c.r) \in \text{categories } \mathbf{do} \\ \text{4: } \mathbf{if} p.c = c.c \mathbf{ then} \\ \text{5: } \text{add } (p.n, c.r) \text{ to } output. \end{array} \right\} \Theta(|categories|).$
 $\left. \begin{array}{l} \text{2: } \mathbf{for} (p.n, p.c) \in \text{products } \mathbf{do} \end{array} \right\} |products| \text{ times.}$

Result: return $\{(p.n, c.r) \mid ((p.n, p.c) \in \text{products}) \wedge ((c.c, c.r) \in \text{categories})\}$.

Theorem

The NESTEDLOOPPC algorithm is correct and has a runtime complexity of $\Theta(|product| \cdot |categories|)$.

Optimizing joins using range queries

Algorithm NESTEDBINARYPC(products, categories):

Input: relations products(*name, category*) and categories(*category, related*),
relation categories ordered.

- 1: *output* := \emptyset .
- 2: **for** (*p.n, p.c*) \in products **do**
- 3: *i* := LOWERBOUND(categories, (*p.c*, ‘’), 0, |categories|).
- 4: **while** *i* < |categories| **and also** categories[*i*].*category* = *p.c* **do**
- 5: add (*p.n*, categories[*i*].*related*) to *output*.
- 6: *i* := *i* + 1.

Result: return $\{(p.n, c.r) \mid ((p.n, p.c) \in \text{products}) \wedge ((c.c, c.r) \in \text{categories})\}$.

Theorem

The NESTEDBINARYPC algorithm is correct and has a runtime complexity of $\Theta(|\text{product}| \cdot \log_2(|\text{categories}|) + |\text{result}|)$.

CONTAINS, LINEARSEARCH, and LOWERBOUND in practice

Algorithm	C++	Java
CONTAINS	<code>std::ranges::contains</code>	<i>collection</i> .contains ^a
LINEARSEARCH	<code>std::find</code>	<i>collection</i> .indexOf ^a
LINEARPREDSEARCH	<code>std::find_if</code>	java.util.stream::filter ^b
LOWERBOUND	<code>std::lower_bound</code> <code>std::upper_bound</code> ^d	java.util.Arrays:: ^c binarySearch
Related libraries	<algorithm>, <ranges>	java.util.Arrays, java.util.ArrayList,...

^aHere, *collection* is a standard Java data collection such as `java.util.ArrayList`.

^bUsing the stream library supported by standard Java data collections.

^cDoes not guarantee to return the offset of the *first* occurrence of a value.

^dReturns the offset of the first element in the list that is strictly larger than the searched-for value.

Print-friendly summary slides

The following slides are the “final” pages of each slide (with intermediate animation steps removed).

Reflection on CONTAINS

Is CONTAINS a *good* algorithm?

CONTAINS is correct and has a runtime complexity of $\Theta(|L|)$ \longrightarrow Sounds good to me!

Critique: CONTAINS is *too specialized* \longrightarrow .

We cannot use CONTAINS for anything else than the contains problem!

Example

- ▶ Searching in only *part* of the list?
- ▶ Finding where v is in the list?

Reflection on CONTAINS

Critique: CONTAINS is *too specialized* \longrightarrow .

We cannot use CONTAINS for anything else than the contains problem!

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

Algorithm LSCONTAINS(L, v):

1: **return** LINEARSEARCH($L, v, 0$) $\neq |L|$.

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

- ▶ With respect to worst case inputs ($v \notin L$): $\Theta(|L|)$.
- ▶ With respect to best case inputs ($v = L[o]$): $\Theta(1)$.

Problem: Modeling runtime complexity in terms of *only the input* limits us!

Reflection on CONTAINS

Algorithm LINEARSEARCH(L, v, o):

Input: L is an *array*, v a value, $0 \leq o \leq |L|$.

1: $r := o$.

/* invariant: “ $o \leq r \leq |L|$ and $v \notin L[o, r)$ ”, bound function: $|L| - r$ */

2: **while** $r \neq |L|$ **and also** $L[r] \neq v$ **do**

3: $r := r + 1$.

4: **return** r .

Result: return the first offset r , $o \leq r < |L|$, with $L[r] = v$ or,
if no such offset exists, $r = |L|$.

What is the runtime complexity of LINEARSEARCH?

Problem: Modeling runtime complexity in terms of *only the input* limits us!

Assume: $L[i] = v$ and i is the *first* offset after o equivalent to v .

The runtime complexity of LINEARSEARCH is $\Theta(i - o)$ with $i = \text{LINEARSEARCH}(L, v, o)$.

Problem: What if we search often?

LINEARSEARCH(L, v, o) can read all of array L : *potentially-high cost*.

Can we do better?

No: We do not know anything about L to help us!

→ we have to look at all elements in L .

Maybe: If we know more about L .

An example of a list

Consider a list *enrolled* with schema

enrolled(*dept*, *code*, *sid*, *date*)

that models a list of all students enrolled for a course.

What if...

We add enrollment data to *the end of the list*.

Question: What do we know about enrolled? → enrolled is *ordered* on *date*!

Searching in an ordered list

3	14	16	18	30	37	44	49	66	76	84	90
---	----	----	----	----	----	----	----	----	----	----	----

↑
 v

Conclusion of comparing $L[i]$ and v ?

$L[i] < v$ As the list L is ordered, every value in $L[0, i]$ is smaller than v .

→ $v \in L$ if and only if $v \in L[i + 1, |L|)$.

$L[i] > v$ As the list L is ordered, every value in $L[i, |L|)$ is larger than v .

→ $v \in L$ if and only if $v \in L[0, i)$.

$L[i] = v$ We found v !

One comparison can remove a large portion of the array.

Binary Search: *Maximize potential* by comparing v with the middle of L .

The *recursive* Binary Search algorithm

Algorithm LOWERBOUNDREC($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

- Is LOWERBOUNDREC correct?
- What is the runtime and memory complexity of LOWERBOUNDREC?

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

1: **if** $begin = end$ **then**

2: **return** $begin$.

3: **else**

4: $mid := (begin + end) \text{ div } 2$.

5: **if** $L[mid] < v$ **then**

6: **return** LOWERBOUNDREC($L, v, mid + 1, end$).

7: **else** $L[mid] \geq v$

8: **return** LOWERBOUNDREC($L, v, begin, mid$).

} *Base case:*

Inspecting $end - begin = 0$ elements.

} *Recursive case:*

Inspecting $end - begin > 0$ elements.

Induction Hypothesis

For any L' , v' , and $0 \leq begin' \leq end' \leq |L'|$ with $0 \leq end' - begin' < m$,

LOWERBOUNDREC($L', v', begin', end'$) returns the correct result.

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

$begin \leq mid < end$

Induction Hypothesis

Induction Hypothesis

Induction Hypothesis

For any L' , v' , and $0 \leq begin' \leq end' \leq |L'|$ with $0 \leq end' - begin' < m$,
LOWERBOUNDREC($L', v', begin', end'$) returns the correct result.

Correctness of LOWERBOUNDREC

Recursion is repetition \longrightarrow induction.

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Termination

Bound function: $end - begin$.

Intermezzo: Runtime complexity of LOWERBOUNDREC

Algorithm LOWERBOUNDREC($L, v, begin, end$):

```
1: if  $begin = end$  then
2:   return  $begin$ .
3: else
4:    $mid := (begin + end) \text{ div } 2$ .
5:   if  $L[mid] < v$  then
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).
7:   else  $L[mid] \geq v$ 
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

} *Base case:*
1 operation.

} *Recursive case:*
1 operation and 1 recursive call.

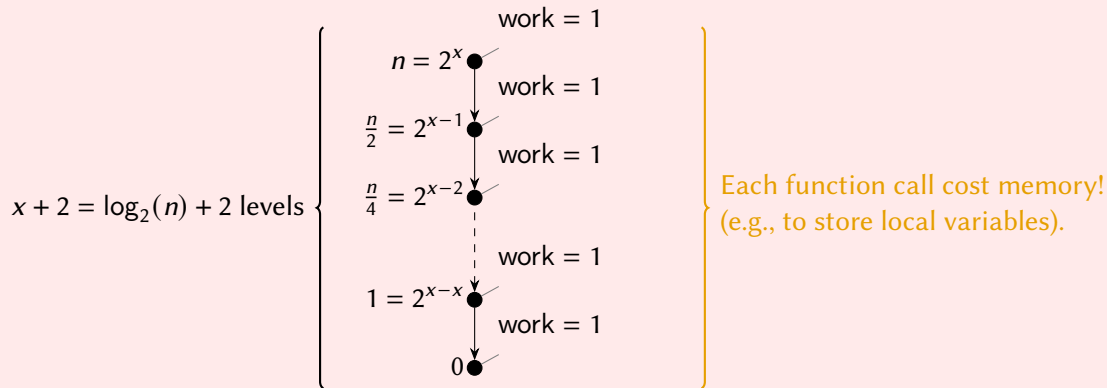
Complexity of LOWERBOUNDREC with $n = end - begin$

$$T(n) = \begin{cases} 1 & \text{if } n = 0; \\ 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1 & \text{if } n \geq 1. \end{cases}$$

Intermezzo: Runtime complexity of LOWERBOUNDREC

Complexity of LOWERBOUNDREC with $n = \text{end} - \text{begin}$ (assume: $n = 2^x$)

$$T(n) = 1 \cdot T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + 1 = \Theta(\log_2(n)).$$



The *recursive* Binary Search algorithm

Algorithm LOWERBOUNDREC($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

```
1: if  $begin = end$  then  
2:   return  $begin$ .  
3: else  
4:    $mid := (begin + end) \text{ div } 2$ .  
5:   if  $L[mid] < v$  then  
6:     return LOWERBOUNDREC( $L, v, mid + 1, end$ ).  
7:   else  $L[mid] \geq v$   
8:     return LOWERBOUNDREC( $L, v, begin, mid$ ).
```

Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

Theorem

LOWERBOUNDREC is correct and has a runtime and memory complexity of $\Theta(\log_2(|L|))$.

The *non-recursive* Binary Search algorithm

Algorithm LOWERBOUND($L, v, begin, end$):

Input: L is an ordered *array*, v a value, and $0 \leq begin \leq end \leq |L|$.

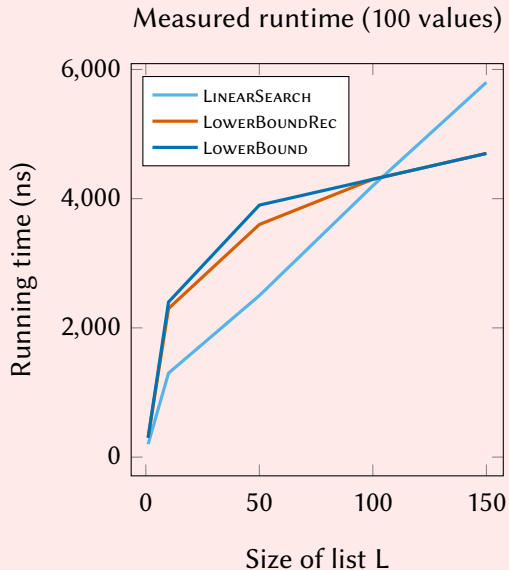
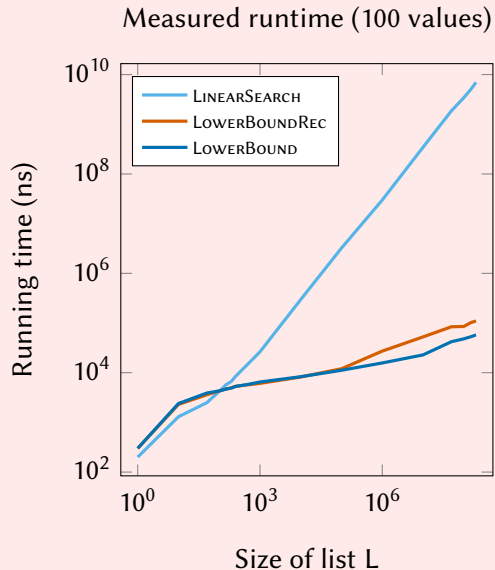
```
1: while  $begin \neq end$  do
2:    $mid := (begin + end) \text{ div } 2$ .
3:   if  $L[mid] < v$  then
4:      $begin := mid + 1$ .
5:   else
6:      $end := mid$ .
7: return  $begin$ .
```

Result: return the first offset r , $begin \leq r < end$, with $L[r] = v$ or,
if no such offset exists, $r = end$.

Theorem

LOWERBOUND is correct, has a runtime complexity of $\Theta(\log_2(|L|))$, and a memory complexity of $\Theta(1)$.

Comparing the complexity of searching



Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Example

Consider a list *enrolled* with schema $\text{enrolled}(\text{dept}, \text{code}, \text{sid}, \text{date})$.

Query: All students enrolled in 2023

Range query on *enrolled* with $[(" ", " ", -1, 2023), (" ", " ", -1, 2024)]$.

We add enrollment data to *the end of the list* \rightarrow enrolled is *ordered* on *date*!

Using Binary Search as a *building block*

Problem

Let L be a list and $[v, w]$ be a range query with $v \leq w$. The solution of the *range query problem* for L and $[v, w]$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Algorithm RANGEQUERY($L, [v, w]$):

Input: L is an ordered *array*, v, w are values, and $v \leq w$.

- 1: $i := \text{LOWERBOUND}(L, v, 0, |L|)$.
- 2: $j := i$.
- 3: **while** $j \leq |L|$ **and also** $L[j] \leq w$ **do**
- 4: $j := j + 1$.
- 5: **return** $L[i, j)$.

Result: return the list $L[m, n)$, $0 \leq m \leq n \leq |L|$, such that
 $L[m, n)$ is the list of all values $e \in L$ with $v \leq e \leq w$.

Theorem

RANGEQUERY is correct and has *all case* runtime complexity $\Theta(\log_2(|L|) + |\text{result}|)$.

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Example

'apple'	'pear'	'orange'
---------	--------	----------

$INSPECT(L, 0) = \text{true}$

$INSPECT(L, 2) = \text{true}$

$INSPECT(L, 1) = \text{true}$

$INSPECT(L, 3) = \text{false}.$

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $INSPECT(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Solving the list-length problem

We have ordered list $0, 1, \dots$ of possible values for $|L|$.

Conclusion of $INSPECT(L, i)$?

$INSPECT(L, i) = \text{true}$ $|L| > i$ (list L has more than i values).

$INSPECT(L, i) = \text{false}$ $|L| \leq i$ (list L has at-most i values).

Issue: no upper bound on the ordered list $0, 1, \dots$ of possible values for $|L|$

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $\text{INSPECT}(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Issue: no upper bound on the ordered list $0, 1, \dots$ of possible values for $|L|$

Guess repeatedly with exponentially-growing guesses.

Algorithm LISTLENGTHUB(L):

Input: L is an *array* of unknown length.

- 1: $n := 1$.
- 2: **while** $\text{INSPECT}(L, n)$ **do**
- 3: $n := 2 \cdot n$.
- 4: **return** n .

Result: return N , $|L| \leq N = 1$ or $|L| \leq N < 2|L|$.

Using the *idea* of Binary Search

Problem

Let L be list of values of *unknown length* and assume we have a function $\text{INSPECT}(L, i)$ that returns true if the list has an i -th value and returns false otherwise.

The *list-length problem* is the problem of finding the length of list L .

Algorithm $\text{LBLISTLENGTH}(L, N)$ with $N := \text{LISTLENGTHUB}(L)$:

```
1: begin, end := 0,  $N$ .  
2: while begin  $\neq$  end do  
3:    $\text{mid} := (\text{begin} + \text{end}) \text{ div } 2$ .  
4:   if  $\text{INSPECT}(L, \text{mid})$  then  
5:      $\text{begin} := \text{mid} + 1$ .  
6:   else  
7:      $\text{end} := \text{mid}$ .  
8: return begin.
```

Result: return the length $|L|$ of array L .

Optimizing joins using range queries

Definition

A join of two lists L and M results in a list A in which each list value is computed from a combination of values $u \in L$ and $v \in M$ according to some *join condition*.

Example (Return pairs (p, r) of product name p and related category r)

products	
<i>name</i>	<i>category</i>
Apple	Fruit
Bok choy	Vegetable
Canelé	Pastry
Donut	Pastry

categories	
<i>category</i>	<i>related</i>
Fruit	Food
Fruit	Produce
Pastry	Food
Vegetable	Food
Vegetable	Produce

Join Result	
<i>name</i>	<i>related</i>
Apple	Food
Apple	Produce
Bok choy	Food
Bok choy	Produce
Canelé	Food
Donut	Food

Optimizing joins using range queries

Algorithm NESTEDLOOPPC(products, categories):

Input: relations products(*name, category*) and categories(*category, related*).

1: *output* := \emptyset .

2: **for** (*p.n, p.c*) \in products **do**

3: **for** (*c.c, c.r*) \in categories **do**

4: **if** *p.c* = *c.c* **then**

5: add (*p.n, c.r*) to *output*.

$\left. \begin{array}{l} \text{3: } \mathbf{for} (c.c, c.r) \in \text{categories } \mathbf{do} \\ \text{4: } \mathbf{if } p.c = c.c \mathbf{ then} \\ \text{5: } \text{add } (p.n, c.r) \text{ to } output. \end{array} \right\} \Theta(|categories|).$
 $\left. \begin{array}{l} \text{2: } \mathbf{for} (p.n, p.c) \in \text{products } \mathbf{do} \\ \text{3: } \mathbf{for} (c.c, c.r) \in \text{categories } \mathbf{do} \\ \text{4: } \mathbf{if } p.c = c.c \mathbf{ then} \\ \text{5: } \text{add } (p.n, c.r) \text{ to } output. \end{array} \right\} |products| \text{ times.}$

Result: return $\{(p.n, c.r) \mid ((p.n, p.c) \in \text{products}) \wedge ((c.c, c.r) \in \text{categories})\}$.

Theorem

The NESTEDLOOPPC algorithm is correct and has a runtime complexity of $\Theta(|product| \cdot |categories|)$.

Optimizing joins using range queries

Algorithm NESTEDBINARYPC(products, categories):

Input: relations products(*name, category*) and categories(*category, related*),
relation categories ordered.

- 1: *output* := \emptyset .
- 2: **for** (*p.n, p.c*) \in products **do**
- 3: *i* := LOWERBOUND(categories, (*p.c*, ‘’), 0, |categories|).
- 4: **while** *i* < |categories| **and also** categories[*i*].category = *p.c* **do**
- 5: add (*p.n*, categories[*i*].related) to *output*.
- 6: *i* := *i* + 1.

Result: return $\{(p.n, c.r) \mid ((p.n, p.c) \in \text{products}) \wedge ((c.c, c.r) \in \text{categories})\}$.

Theorem

The NESTEDBINARYPC algorithm is correct and has a runtime complexity of $\Theta(|\text{product}| \cdot \log_2(|\text{categories}|) + |\text{result}|)$.

CONTAINS, LINEARSEARCH, and LOWERBOUND in practice

Algorithm	C++	Java
CONTAINS	<code>std::ranges::contains</code>	<i>collection</i> .contains ^a
LINEARSEARCH	<code>std::find</code>	<i>collection</i> .indexOf ^a
LINEARPREDSEARCH	<code>std::find_if</code>	<code>java.util.stream::filter</code> ^b
LOWERBOUND	<code>std::lower_bound</code> <code>std::upper_bound</code> ^d	<code>java.util.Arrays::</code> ^c <code>binarySearch</code>
Related libraries	<code><algorithm></code> , <code><ranges></code>	<code>java.util.Arrays</code> , <code>java.util.ArrayList</code> ,...

^aHere, *collection* is a standard Java data collection such as `java.util.ArrayList`.

^bUsing the stream library supported by standard Java data collections.

^cDoes not guarantee to return the offset of the *first* occurrence of a value.

^dReturns the offset of the first element in the list that is strictly larger than the searched-for value.

CONTAINS, LINEARSEARCH, and LOWERBOUND in practice

Algorithm	C++	Java
CONTAINS	<code>std::ranges::contains</code>	<i>collection</i> .contains ^a
LINEARSEARCH	<code>std::find</code>	<i>collection</i> .indexOf ^a
LINEARPREDSEARCH	<code>std::find_if</code>	<code>java.util.stream::filter</code> ^b
LOWERBOUND	<code>std::lower_bound</code> <code>std::upper_bound</code> ^d	<code>java.util.Arrays::</code> ^c <code>binarySearch</code>
Related libraries	<code><algorithm></code> , <code><ranges></code>	<code>java.util.Arrays</code> , <code>java.util.ArrayList</code> ,...

^aHere, *collection* is a standard Java data collection such as `java.util.ArrayList`.

^bUsing the stream library supported by standard Java data collections.

^cDoes not guarantee to return the offset of the *first* occurrence of a value.

^dReturns the offset of the first element in the list that is strictly larger than the searched-for value.

CONTAINS, LINEARSEARCH, and LOWERBOUND in practice

Algorithm	C++	Java
CONTAINS	<code>std::ranges::contains</code>	<i>collection</i> .contains ^a
LINEARSEARCH	<code>std::find</code>	<i>collection</i> .indexOf ^a
LINEARPREDSEARCH	<code>std::find_if</code>	<code>java.util.stream::filter</code> ^b
LOWERBOUND	<code>std::lower_bound</code> <code>std::upper_bound</code> ^d	<code>java.util.Arrays::</code> ^c <code>binarySearch</code>
Related libraries	<code><algorithm></code> , <code><ranges></code>	<code>java.util.Arrays</code> , <code>java.util.ArrayList</code> ,...

^aHere, *collection* is a standard Java data collection such as `java.util.ArrayList`.

^bUsing the stream library supported by standard Java data collections.

^cDoes not guarantee to return the offset of the *first* occurrence of a value.

^dReturns the offset of the first element in the list that is strictly larger than the searched-for value.