

# [ Searching ]

- Basic Points
  - Search returns
    - found
    - not found
  - Related terms
    - table lookup
    - find operation
  - Performance
    - may vary with size of collection
- Search Methods
  - internal (memory)
  - external (file/disk)
  - static (content fixed)
  - dynamic (changeable)
  - key comparison
  - key transformation
  - **Searching is closely linked to sorting**

# Searching: Examples

- Sequence lookup
  - hash tables
  - compiler symbol tables
  - database systems
  - key searching  
(info = key + data)
- Sequence
  - array / file / table / database
  - linked list / strings (patterns)

- **Sequential Searching**

- **$O(n)$**

collection of e where

e = key + data (+ ref)

**ref = first(C)**

**while (not eoc(C))**

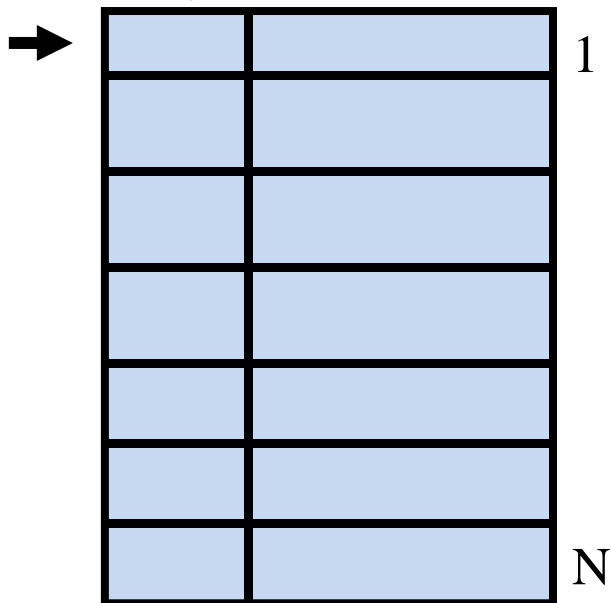
**if (V = value(e)) return T**

**else ref = next(ref)**

# [ Abstract Model ]

## ■ Collection

key data

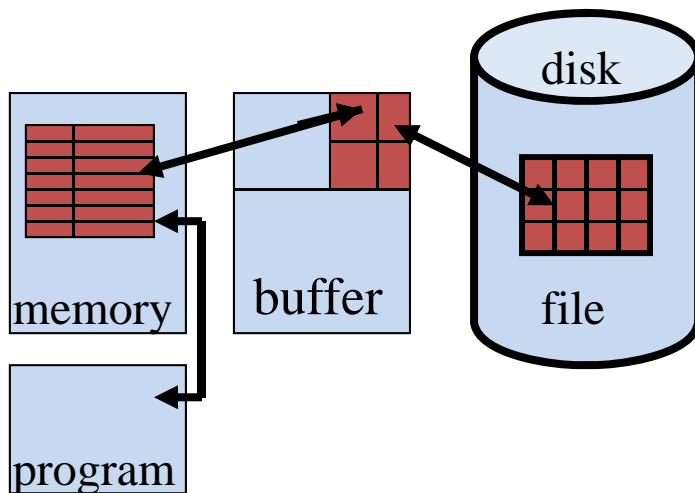


## ■ Operations

- size:  $C \rightarrow \text{int}$
- is\_empty:  $C \rightarrow \text{Bool}$
- get\_key:  $\text{ref} \rightarrow \text{key}$
- get\_data:  $\text{ref} \rightarrow \text{data}$
- find:  $C \times \text{key} \rightarrow \text{ref}$
- add:  $C \times e \rightarrow C$
- remove:  $C \times \text{ref} \rightarrow C$
- display:  $C \rightarrow C$

# Implementation Issues

## ■ Example



## ■ Performance

- size of C
- disk swapping/buffering
- search methods
  - hashing / B-trees
- key / data ratio
- number of key collections
  - primary (one)
  - secondary ( $\geq 0$ )

# [ Sorted Sequence ]

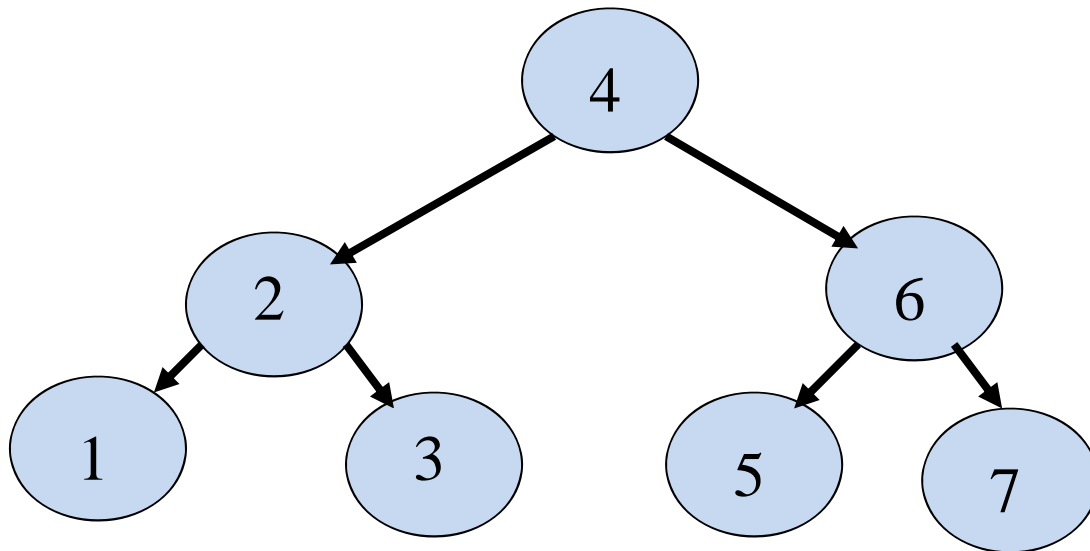
- Collection (sorted)

	key	data	
→	A		1
	C		
	D		
	E		
	F		
	G		
	H		N

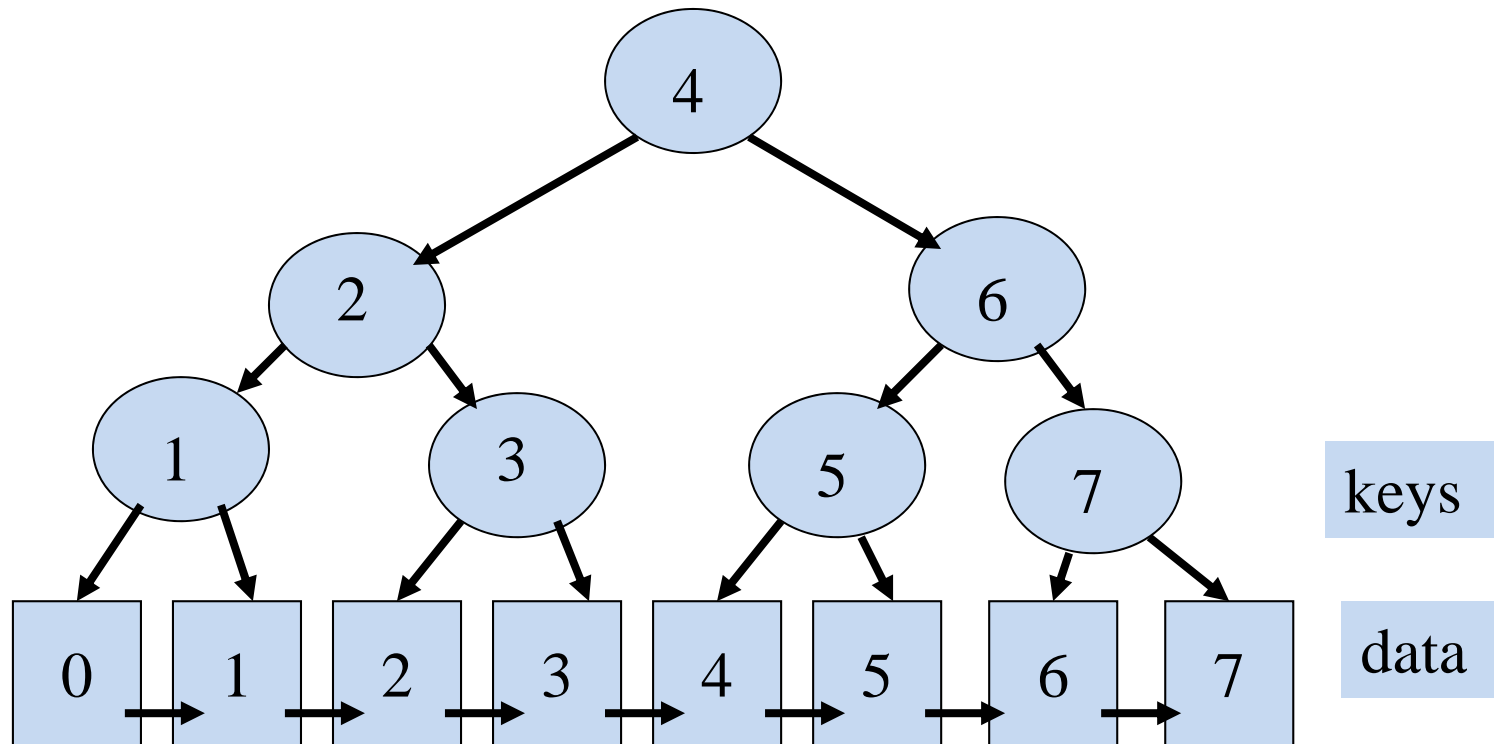
- Keys sorted

- binary search
  - initialise lo\_ref / hi\_ref  
 $K_{lo} \leq K \leq K_{hi}$
  - get mid-point  $K_{mid}$
  - compare
    - $K = K_{mid} \Rightarrow$   
**found**
    - $>$  **search upper**
    - $<$  **search lower**
  - repeat process

# [ Binary Search Tree ]

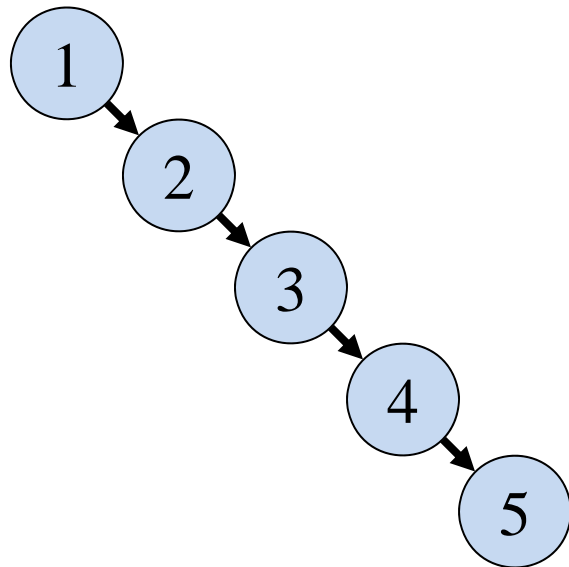


# Binary Search Tree + Sequence

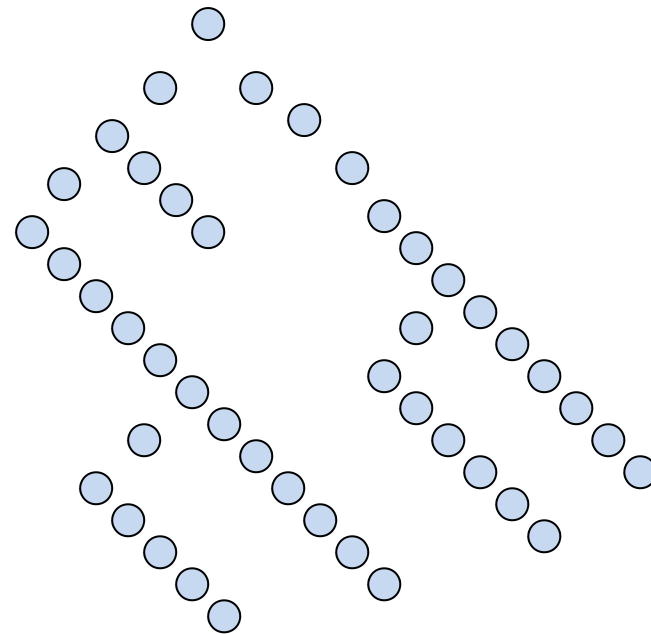


# [ Problem (pathological) cases ]

- bst created from  
1 2 3 4 5 (sorted)

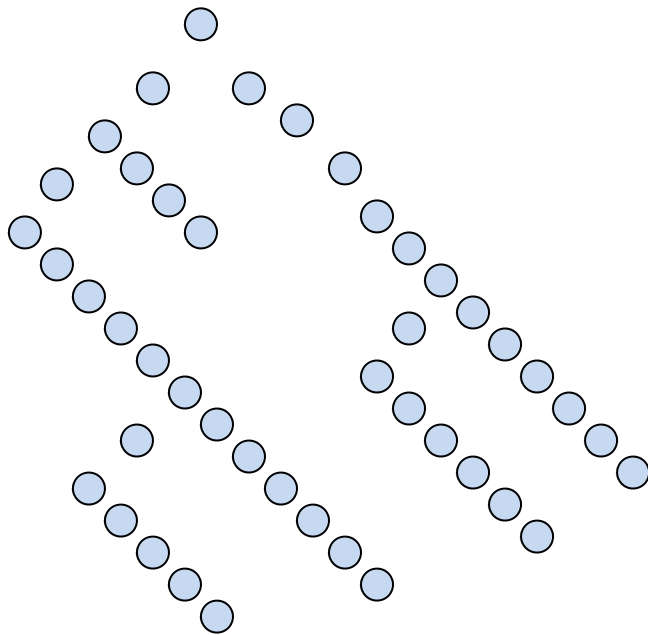


- Deep trees





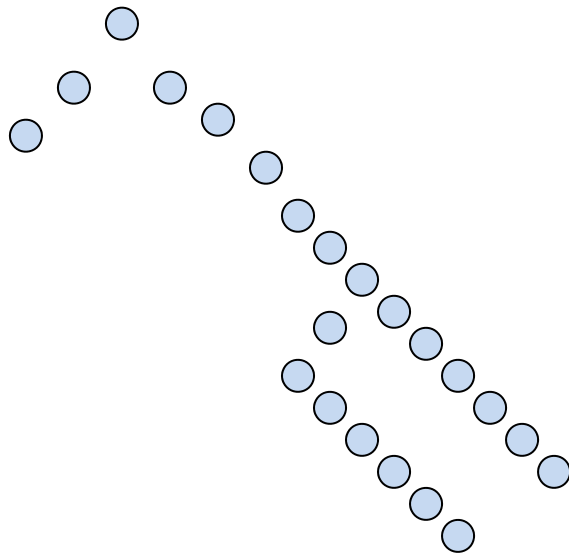
# [ Deep Trees ]



## ■ Disadvantages

- path length (# accesses)
- if node = disk page => high cost for swapping
- BTs can easily become deep as N increases
- $N(\text{max}) = 2^{d+1} - 1$

# [ Unbalanced Trees ]

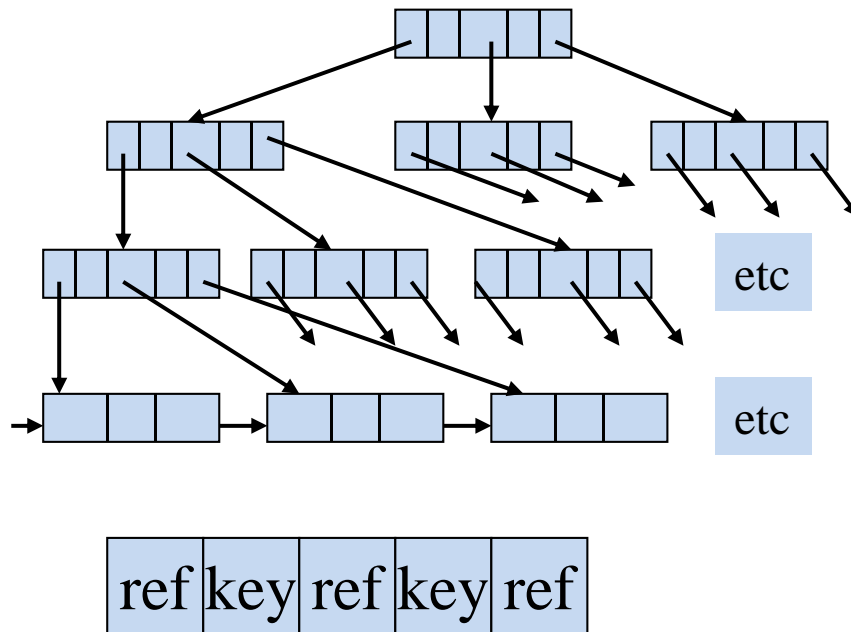


## ■ Disadvantages

- path length (# accesses)
- if node = disk page => high cost for swapping
- BTs can easily become deep as N increases
- depends on distribution of the input data
- static/dynamic analysis

# Solution: Balanced Shallow Trees

## Knuth's B+ variant of a B-tree



## Advantages

- shorter search paths
- index + sequence set
- index =  $p_1 k_1 p_2 k_2 \dots$
- sequence =  $d_1 d_2 d_3 \dots$
- sequence set may still be accessed linearly

## Issues

- add/remove algorithms

# Summary

- Search → found / not found
- performance depends on size(C) -  $O(N)$ ,  $O(\log N)$ , ...
- internal / external searching
- source data distribution
- unsorted => linear search -  $O(n)$
- sorted => binary search -  $O(\log N)$
- -ve : deep / unbalanced trees
- +ve : shallow / balanced trees
- b-tree family of trees
- based on key comparison

## Applications

- file / DB searching
- table lookup
  - symbol / function tables
  - primary / secondary key indexes
- lexicographical search
  - natural language
  - pattern matching
- fundamental to CS