

Fundamentals of Data Structure

— — By Rukawa Yuan

1 Basis

Time complexity & Space complexity

When analyzing the time complexity, generally $T_{avg}(N)$ and $T_{worst}(N)$ are analyzed.

1. $T(N) = O(f(N))$ if there are positive constants c and n_0 such that $T(N) \leq c \times f(N)$ for all $N \geq n_0$.
2. $T(N) = \Omega(g(N))$ if there are positive constants c and n_0 such that $T(N) \geq c \times g(N)$ for all $N \geq n_0$.
3. $T(N) = \Theta(h(N))$ if and only if $T(N) = O(h(N))$ and $T(N) = \Omega(h(N))$.
4. $T(N) = o(h(N))$ if $T(N) = O(p(N))$ and $T(N) \neq \Theta(p(N))$.

Always pick the smallest upper bound and the largest lower bound.

- If $T_1(N) = O(f(N))$ and $T_2(N) = O(g(N))$, then:

$$T_1(N) + T_2(N) = \max(O(f(N)), O(g(N)))$$

$$T_1(N) \times T_2(N) = O(f(N) \times g(N))$$

- $\log^k N = O(N)$ for all constant k (Logarithms grow very slowly)

Example: Fibonacci number

```
1  int Fib(int N){
2      if(N<=1){
3          return 1;
4      }else{
5          return Fib(N-1)+Fib(N-2);
6      }
7  }
```

$$T(N) = T(N-1) + T(N-2) + 2$$

$$\left(\frac{3}{2}\right)^N \leq Fib(N) \leq \left(\frac{5}{3}\right)^N$$

When $T(N) = O(f(N))$, it is guaranteed that $\lim_{N \rightarrow \infty} \frac{f(N)}{T(N)} = \text{constant}$ or $\lim_{N \rightarrow \infty} \frac{f(N)}{T(N)} = +\infty$

2 Examples

Max Subsequence Sum

- **Problem Description:** Given (possibly negative) integers A_1, A_2, \dots, A_N , find the maximum value of $\sum_{k=i}^j A_k$.

- **Algorithm 1**

```
1  int MaxSubsequenceSum(int A[],int N){
2      int ThisSum,MaxSum,i,j,k;
3      MaxSum=0;
4      for(i=0;i<N;i++){
5          for(j=i;j<N;j++){
6              ThisSum=0;
7              for(k=i;k<=j;k++){
8                  ThisSum+=A[k];
9              }
10             if(ThisSum>MaxSum){
11                 MaxSum=ThisSum;
12             }
13         }
14     }
15     return MaxSum;
16 }
```

This algorithm is simple and easy to understand: Find all the possible subsequences and compare them to find the maximum value.

Nevertheless, it is a typically bad algorithm with $T(N) = O(N^3)$.

- **Algorithm 2**

```
1  int MaxSubsequenceSum(int A[],int N){
2      int ThisSum,MaxSum,i,j;
3      MaxSum=0;
4      for(i=0;i<N;i++){
5          ThisSum=0;
6          for(j=i;j<N;j++){
7              ThisSum+=A[j];
8              if(ThisSum>MaxSum){
9                  MaxSum=ThisSum;
10             }
11         }
12     }
13     return MaxSum;
14 }
```

The idea of this algorithm is basically simple to the previous one. However, it does an improvement because it scans through all the subsequences in a smarter way and reduce time complexity to $O(N^2)$.

- **Algorithm 3 Divide and Conquer**

```
1  int MaxSubSum(int A[],int left,int right){
2      int MaxLeftSum,MaxRightSum;
3      int MaxLeftBorderSum,MaxRightBorderSum;
4      int LeftBorderSum,RightBorderSum;
5      int center,i;
```

```

6
7     if(left==right){
8         if(A[left]>0){
9             return A[left];
10        }else{
11            return 0;
12        }
13    }
14
15    center=(left+right)/2;
16    MaxLeftSum=MaxSubSum(A, left, center);
17    MaxRightSum=MaxSubSum(A, center+1, right);
18
19    MaxLeftBorderSum=0;
20    LeftBorderSum=0;
21    for(i=center; i>=left; i--){
22        LeftBorderSum+=A[i];
23        if(LeftBorderSum>MaxLeftBorderSum){
24            MaxLeftBorderSum=LeftBorderSum;
25        }
26    }
27
28    MaxRightBorderSum=0;
29    RightBorderSum=0;
30    for(i=center+1; i<=right; i++){
31        RightBorderSum+=A[i];
32        if(RightBorderSum>MaxRightBorderSum){
33            MaxRightBorderSum=RightBorderSum;
34        }
35    }
36
37    return MaxThree(MaxLeftSum, MaxRightSum, MaxLeftBorderSum+MaxRightBorderSum);
38 }
39
40 int MaxSubsequenceSum(int A[], int N){
41     return MaxSubSum(A, 0, N-1);
42 }

```

A complicated algorithm which involves recursions. The idea of the algorithm is to divide the whole set into two equal parts, calculate the max length of left part and right part, calculate the max length of subsequence which begins from center and extend leftward and rightward and return the maximum value of the three as the result.

Time complexity analysis:

$$T(N) = 2T\left(\frac{N}{2}\right) + N$$

$$T(N) = O(N \log N)$$

which is much better than $O(N^3)$ and $O(N^2)$.

■ Algorithm 4 Online Algorithm

```

1  int MaxSubSequenceSum(int A[], int N){
2      int ThisSum, MaxSum, j;
3      ThisSum=0;
4      MaxSum=0;
5      for(j=0; j<N; j++){
6          ThisSum+=A[j];
7          if(ThisSum>MaxSum){
8              MaxSum=ThisSum;

```

```

9      }
10     if(ThisSum<0){
11         ThisSum=0;
12     }
13 }
14 return MaxSum;
15 }

```

As is shown above, the algorithm is an online algorithm, which means that it can always return the correct result currently whenever it stops.

The idea of the algorithm is easy to understand. It simply scans through the whole array, add each value to *ThisSum* and update it if necessary. If $ThisSum < 0$, it is set 0 because we can exclude the negative part in the max subsequence.

The time complexity of this algorithm is $O(N)$, the best algorithm of all.

Binary Search

- **Problem Description:** Given a sorted array A and its size N , we need to find a particular element X . If found, return the index, otherwise return -1 .
- **Algorithm**

```

1  int BinarySearch(int A[],int N,int X){
2      int Low,Mid,High;
3      Low=0;
4      High=N-1;
5      while(Low<=High){
6          Mid=(Low+High)/2;
7          if(A[Mid]==X){
8              return Mid;
9          }else if(A[Mid]<X){
10             Low=Mid+1;
11          }else{
12             High=Mid-1;
13          }
14      }
15      return -1;
16  }

```

Time complexity analysis:

$$T(N) = T\left(\frac{N}{2}\right) + C$$

$$T(N) = O(\log N)$$

3 List

■ Key Operations:

1. Finding the k -th item in the list, $0 \leq k < N$
2. Inserting a new item after the k -th item of the list, $0 \leq k < N$
3. Deleting an item from the list

List ADT can be implemented with array and linked list.

■ Array Implementation:

1. Finding operation is quick and simple.
2. Does not involve pointers and the allocation of memory.
3. Maxsize has to be estimated in advance.
4. Insertion and Deletion take $O(N)$ time and involves numerous movements of data.

■ Linked List Implementation:

■ Initialization:

Singly Linked List:

```
1  typedef struct node* List;
2  struct node{
3      ElementType Data;
4      List Next;
5  };
```

Doubly Linked Circular List:

```
1  typedef struct node* List;
2  struct node{
3      List Leftptr;
4      ElementType Data;
5      List Rightptr;
6  }
```

■ Insertion:

```
1  List temp=(List)malloc(sizeof(struct node));
2  temp->Data=1;
3  temp->Next=Current->Next;
4  Current->Next=temp;
```

Always connect temp to the next node first.

Time complexity: $O(1)$

■ Deletion:

```
1  temp=Pre->Next;
2  Pre->Next=temp->Next;
3  free(temp);
```

Always connect the previous node to next node first.

Time complexity: $O(1)$

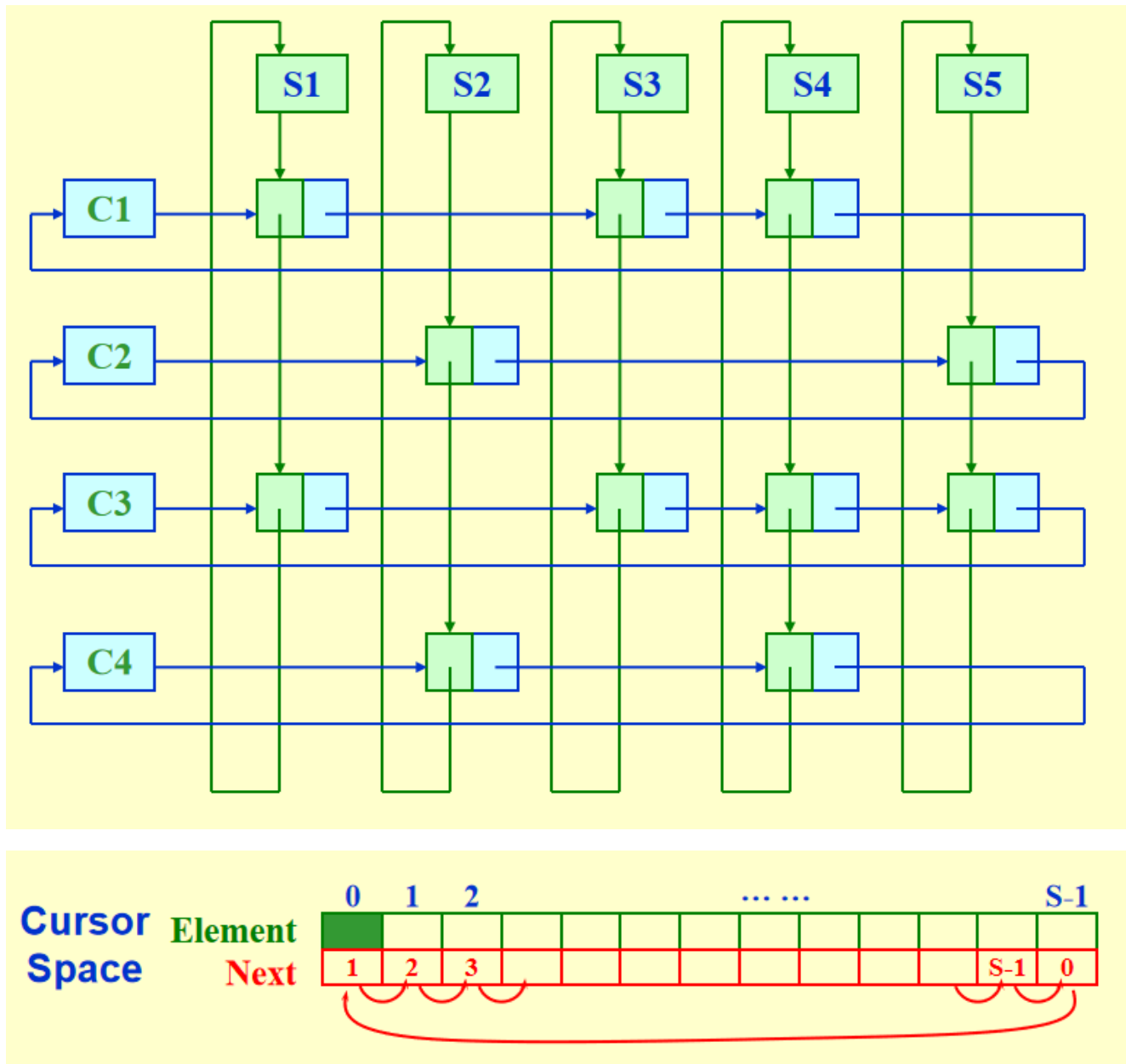
- Application: Polynomial ADT

Objects: $P(x) = a_1x^{e_1} + a_2x^{e_2} + \dots + a_nx^{e_n}$

Implementation: Linked List

```
1 typedef struct node* Polynomial;
2 struct node{
3     int Coefficient;
4     int Exponent;
5     Polynomial Next;
6 }
```

Other ADT: Multilist and Cursor Implementation of Linked List.



4 Stack

A stack is a Last-In-First-Out (LIFO) list, namely, an ordered list in which insertions and deletions are made at the top only.

- Key Operations:

```
1 Push(ElementType X,Stack S);
2 Pop(Stack S);
3 ElementType Top(Stack S);
```

A pop on an empty stack results in ADT error, a push on a full stack results in implementation error.

- Linked List Implementation:

A header node is used (dummy header) .

- Push:

```
1 TmpCell->Next=S->Next;
2 S->Next=TmpCell;
```

- Top:

```
1 return S->Next->Element;
```

- Pop:

```
1 Temp=S->Next;
2 S->Next=S->Next->Next;
3 free(Temp);
```

- Array Implementation:

The stack model is supposed to be well encapsulated.

```
1 struct Stack{
2     int Size;
3     int Pointer;
4     ElementType* Array;
5 };
```

Applications

Balancing Symbols

- **Problem Description:** Check if Parenthesis (), Brackets [] and Braces {} are balanced.

- **Algorithm**

```
1 bool BalancingSymbols(char Expression[],int N){
2     char Stack[MAXN];
```

```

3     int i,temp;
4     for(i=0;i<N;i++){
5         temp=Expression[i];
6         if(temp=='('||'['||'{'){
7             Push(temp,Stack);
8         }else{
9             if(Stack is Empty){
10                return false;
11            }else if(Top(Stack)!=Corresponding Symbol){
12                return false;
13            }else{
14                Pop(Stack);
15            }
16        }
17    }
18    if(Stack is Empty){
19        return true;
20    }else{
21        return false;
22    }
23 }

```

$T(N) = O(N)$, an online algorithm.

Calculating Postfix Expression

While a computer is calculating an expression, it firstly transform infix expression to postfix. Postfix expressions have advantages, like no parenthesis is needed.

- **Problem Description:** Given a postfix expression (Reverse Polish Notation) , calculate its value.
- **Example:** $6\ 2\ /\ 3\ -\ 4\ 2\ *\ +\ =\ 8$
- **Algorithm**

```

1     int PostfixCalculation(char Expression[],int N){
2         char Stack[MAXN];
3         int i;
4         char temp,temp1,temp2;
5         for(i=0;i<N;i++){
6             temp=Expression[i];
7             if(temp>='0'&&temp<='9'){
8                 Push(temp,Stack);
9             }else{
10                temp1=Top(Stack);
11                Pop(S);
12                temp2=Top(Stack);
13                Pop(S);
14                Push(Calculate(temp1,temp2,temp),Stack);/*Details Omitted*/
15            }
16        }
17        return Stack[0];
18    }

```

$T(N) = O(N)$

Infix to Postfix Conversion

- **Problem Description:** Transform an infix expression to postfix expression.
- **Example:** $a + b * c - d = a b c * + d -$
- **Algorithm**

```
1 char* InfixtoPostfix(char Infix[],int N){
2     Stack S;
3     char* Result;
4     for(i=0;i<N;i++){
5         temp=Infix[i];
6         if(temp is an Operand){
7             Append temp to Result;
8         }else{
9             if(temp=='('){
10                Pop and append until '(';
11            }
12            while(Priority(Top(S))>=Priority(temp)){
13                Append Top(S) to Result;
14                Pop(S);
15            }
16            Push(temp,S);
17        }
18    }
19    Pop and append until S is empty;
20    return Result;
21 }
```

Priority:

'(' (Before Push) > '*' '/' > '+' '-' > '(' (After Push)

$$T(N) = O(N)$$

System Stack

While running recursive programs, Return Address, Stack Frame, Local Variables are all pushed into the stack. If too many recursions are called, the system stack will be overflowed and the system will crash.

Tail recursions can always be removed. Non-recursive programs are generally faster than equivalent recursive programs.

When you write a recursive program, the compiler will automatically transformed the recursive program to non-recursive program. (Goto top may be used)

Example:

```
1 void PrintList(List L){
2     if(L!=NULL){
3         printf("%d",L->Element);
4         PrintList(L->Next);
5     }
6 }
```

Will be transformed to:

```
1 void Printlist(List L){
2   top:if(L!=NULL){
3       printf("%d",L->Element);
4       L=L->Next;
5       goto top;
6   }
7 }
```

5 Queue

A queue is a First-In-First-Out (FIFO) list, namely, an ordered list in which insertions take place at one end and deletions take place at the opposite end.

■ Key Operations:

```
1 Enqueue(ElementType X,Queue Q);
2 ElementType Front(Queue Q);
3 Dequeue(Queue Q);
```

Generally queue ADT is implemented by arrays, because linked list implementation is way too trivial. For arrays, there are 2 categories: simple array and circular queue.

```
1 struct Queue{
2     int Front;
3     int Rear;
4     int Size;
5     int Capacity;
6     ElementType* Array;
7 };
```

■ Circular Queue:

- Initially $\text{Front} = \text{Rear} + 1$.
- If there is only 1 element, $\text{Front} = \text{Rear}$
- There is bound to be 1 place that is empty. If all spaces are filled with elements, we cannot distinguish empty and full. (If a circular queue is full, $\text{Front} = \text{Rear} + 2$.)
- Arithmetic like *mod* will be included to implement the function of "Circular". (Like $(\text{Rear}++) \% \text{Size}$)

6 Tree

Terminologies:

Root	Subtree	Edge	Node
Degree of a Node	Degree of a Tree (Max Node Degree)	Parent	Children (Left, Right)
Siblings	Leaf Node	Path (Sequence of Nodes)	Length of Path
Depth (D(root)=0)	Height (H(Leaf)=0)	Height (Depth) of Tree (Max)	Ancestor&Descendant

- Question: How to change a general tree to a binary tree?
- Using **First Child-Next Sibling** Representation

Expression Tree

An expression tree is composed of operators and operands, where operands are leaf nodes and the calculation of the expression can be conducted from bottom to top.

- **Problem Description:** Given an infix/postfix expression, constructing an expression tree.
- **Algorithm (In Pseudocode)**

```

1 Transform the expression to POSTFIX expression
2 Create a Stack
3   Scan through the postfix expression:
4       if Current is Operand:
5           Push(Current,Stack)
6       else:
7           Get two elements from Stack, temp1 and temp2
8           Construct a tree:
9               Current as the root, temp1 and temp2 be left and right subtrees
10          Push the tree into Stack
11 return Stack[0]
```

Tree Traversals

Preorder Traversal, Inorder Traversal, Postorder Traversal and Levelorder Traversal

Preorder Traversal

Visit Current Node, Left Subtree and Right Subtree sequentially.

```

1 void Preorder(Tree_ptr Tree){
2     visit(Tree);
3     Preorder(Tree->Left);
4     Preorder(Tree->Right);
5 }
```

Postorder Traversal

Visit Left Subtree, Right Subtree and Current Node sequentially.

```
1 void Postorder(Tree_ptr Tree){
2     Postorder(Tree->Left);
3     Postorder(Tree->Right);
4     visit(Tree);
5 }
```

Levelorder Traversal

Scan through each line sequentially.

```
1 void Levelorder(Tree_ptr Tree){
2     Queue Q;
3     Enqueue(Tree);
4     while(Q is not empty){
5         visit(Q[Front]);
6         Dequeue(Q);
7         for(Each child C of Q[Front]){
8             Enqueue(C);
9         }
10    }
11 }
```

Inorder Traversal

Visit Left Subtree, Current Node and Right Subtree sequentially.

Inorder traversal can be implemented both iteratively and recursively.

■ Recursive Implementation

```
1 void Inorder(Tree_ptr Tree){
2     Inorder(Tree->Left);
3     visit(Tree);
4     Inorder(Tree->Right);
5 }
```

■ Iterative Implementation

```
1 void Inorder(Tree_ptr Tree){
2     Stack S;
3     while(1){
4         while(Tree!=NULL){
5             Push(Tree,S);
6             Tree=Tree->Left;
7         }
8         Tree=Top(S);
9         Pop(S);
10        if(Tree==NULL){
11            break;
12        }
13        Visit(Tree->Element);
14    }
```

```

14         Tree=Tree->Right;
15     }
16 }

```

That is a little bit complicated. Let's describe the process in human language:

1. Keep going leftward if possible and put all nodes into stack.
2. When meeting the left end (leaf node), visit current node and its right subtree.
3. Go back a step and visit the node and its right subtree, repeatedly.

For an given expression tree, the Inorder Traversal corresponds to Infix Expression, the Preorder Traversal corresponds to Prefix Expression, and the Postorder Traversal corresponds to Postfix Expression.

Threaded Binary Tree

Take the Inorder threaded binary tree as an example:

If $Tree \rightarrow Left$ is NULL, replace it with a pointer to the inorder **predecessor** of the Tree.

If $Tree \rightarrow Right$ is NULL, replace it with a pointer to the inorder **successor** of the Tree.

There must not be any loose threads. Therefore a **head node** is needed such that the left pointer of the first item in traversal and the right pointer of the last item in traversal point to it.

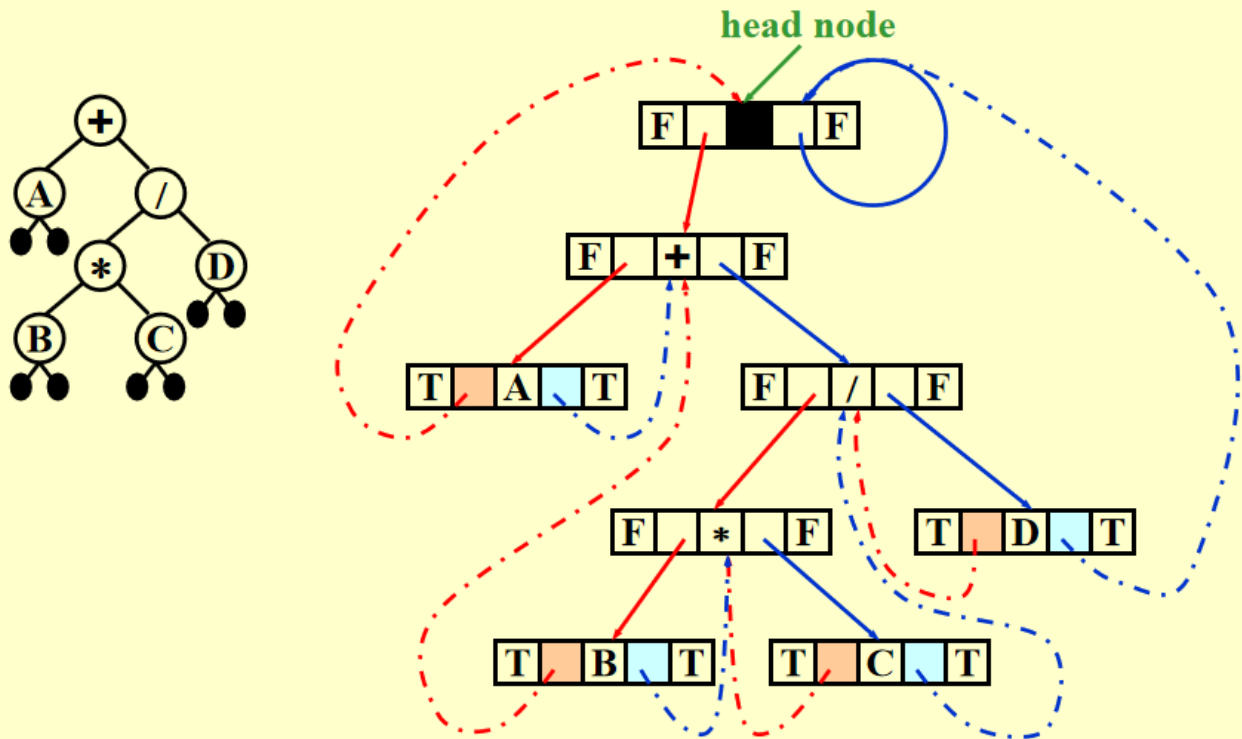
```

1  typedef struct ThreadedTreeNode* PtrToThreadedNode;
2  struct ThreadedTreeNode{
3      bool LeftThread;
4      PtrToThreadedNode LeftPtr;
5      ElementType Element;
6      bool RightThread;
7      PtrToThreadedNode RightPtr;
8  };

```

[[Example]] Given the syntax tree of an expression (infix)

$$A + B * C / D$$



Some Terminologies:

Skewed Binary Tree: Trees that degenerate to linked list

Complete Binary Tree: All the leaf nodes are on two adjacent levels. On the last level, leaf nodes always begin from left.

Full Binary Tree: Each node excluding leaf node has exactly m children.

Perfect Binary Tree: A tree of height h has $2^{h+1} - 1$ nodes.

Properties of Binary Tree

1. For trees with N nodes, they have $N - 1$ edges.

2. The maximum number of nodes on level i is 2^i .

The maximum number of nodes in a binary tree of depth k is $2^{k+1} - 1$.

3. Let n_i denotes a node of degree i , we have: $n_0 = n_2 + 1$.

Prove: $n = n_0 + n_1 + n_2$, $E = n - 1 = n_0 + n_1 + n_2 - 1 = n_1 + 2n_2$

$\therefore n_0 = n_2 + 1$

An Important Problem

Given a **preorder traversal** sequence and an **inorder traversal** sequence, build a binary from those two arrays. (Similarly we can define a problem whose input is the inorder and postorder traversal sequences and output is a binary tree.)

Here is a problem on **LeetCode**:

Structure and function interface:

```

1  /**
2   * Definition for a binary tree node.
3   * struct TreeNode {
4   *     int val;
5   *     struct TreeNode *left;
6   *     struct TreeNode *right;
7   * };
8   */
9  struct TreeNode* buildTree(int* preorder, int preorderSize, int* inorder, int inorderSize);

```

Solution:

```

1  struct TreeNode* buildTree(int* preorder, int preorderSize, int* inorder, int inorderSize) {
2      if(preorderSize==0){
3          return NULL;
4      }
5      struct TreeNode* T=(struct TreeNode*)malloc(sizeof(struct TreeNode));
6      int root=preorder[0];
7      int i;
8      for(i=0;i<inorderSize;i++){
9          if(inorder[i]==root){
10             break;
11         }
12     }
13     T->val=root;
14     T->left=buildTree(preorder+1,i,inorder,i);
15     T->right=buildTree(preorder+1+i,preorderSize-1-i,inorder+i+1,inorderSize-1-i);
16     return T;
17 }

```

The idea is simple, while the implementation is a little bit complicated. Given the preorder traversal sequentially visit the root node, the left subtree and the right subtree, we know that the first element in preorder traversal is the root of the current tree. Then we find the root node in the inorder traversal. Because the inorder traversal visits the left subtree, the current node and the right subtree sequentially, we can divide the inorder traversal sequence into 3 parts: Left subtree, current node and right subtree, while the current node is the first element in preorder sequence. Then we recursively conduct the procedure and eventually build a corresponding binary tree.

Search Tree ADT

A binary search tree is a binary tree with distinct integers as its keys.

The keys in the left subtree are smaller than root node, while the keys in the right subtree are larger than root node. Such property makes searching and finding easier.

■ Key Operations:

1. Find a particular element in BST
2. Insert a new element into BST
3. Delete a certain element in BST

■ Find

Recursive:

```

1 Tree_ptr Find(ElementType X, Tree_ptr T){
2     if(T==NULL){
3         return NULL;
4     }
5     if(T->Element==X){
6         return T;
7     }else if(T->Element<X){
8         Find(X, T->Right);
9     }else{
10        Find(X, T->Left);
11    }
12 }

```

Iterative:

```

1 Tree_ptr Find(ElementType X, Tree_ptr T){
2     while(T!=NULL){
3         if(T->Element==X){
4             return T;
5         }else if(T->Element>X){
6             T=T->Left;
7         }else{
8             T=T->Right;
9         }
10    }
11    return NULL;
12 }

```

$T(N) = O(\text{Depth}) = O(\log N)$ (Best Case) = $O(N)$ (Worst Case)

FindMin and **FindMax** are basically simple and identical:

For FindMin: Keep going leftward and return the value of the terminal node.

■ Insert

```

1 Tree_ptr Insert(ElementType X, Tree_ptr T){
2     if(T==NULL){
3         T=(Tree_ptr)malloc(sizeof(struct TreeNode));
4         T->Element=X;
5         T->Left=NULL;
6         T->Right=NULL;
7     }
8     if(T->Element>X){
9         T->Left=Insert(X, T->Left);
10    }else if(T->Element<X){
11        T->Right=Insert(X, T->Right);
12    }
13    return T;
14 }

```

$T(N) = O(\text{Depth}) = O(\log N)$ (Best Case) = $O(N)$ (Worst Case)

■ Delete

There are 3 cases in total:

1. Deleting a **leaf node**: Simply set the pointer of its parent node to NULL.

2. Deleting a node of **degree 1**: Replace the node by its child.

3. Deleting a node of **degree 2**:

Picking the largest element in left subtree or the smallest element in right subtree and replace the current node with picked one.

Deleting the node from the subtree. (*Degree = 0 or 1*)

That is, we transform the problem to delete a node of degree 0 or 1.

```
1  Tree_ptr Delete(ElementType X, Tree_ptr T){
2      if(X>T->Element){
3          T->Right=Delete(X, T->Right);
4      }else if(X<T->Element){
5          T->Left=Delete(X, T->Left);
6      }else{
7          if(T->Left==NULL&&T->Right==NULL){
8              return NULL;
9          }else if(T->Left==NULL){
10             return T->Right;
11          }else if(T->Right==NULL){
12             return T->Left;
13          }else{
14              TmpCell=FindMin(T->Right);    /*Or FindMax(T->Left)*/
15              T->Element=TmpCell->Element;
16              T->Right=Delete(T->Element, T->Right);
17              return T;
18          }
19      }
20 }
```

$$T(N) = O(\text{Depth}) = O(\log N) (\text{Best Case}) = O(N) (\text{Worst Case})$$

Lazy Deletion: Does not actually delete an node. Instead, add a mark to each node that determines whether the node is legitimate or empty. Lazy deletions can significantly improve the efficiency of operations.

7 Priority Queue (Heap)

A heap is a ADT such that the element on the top always has the highest or lowest priority.

■ Key Operations:

1. Insert an element into the heap
2. Delete the element on the top (DeleteMin)
3. Get the value of the element on the top (FindMin)

Binary Heap

A binary tree with n nodes and height h is complete iff its nodes correspond to the nodes numbered from 1 to n in the perfect binary tree of height h .

A binary heap is a complete binary tree.

When using array implementation, $H[0]$ will not be used, elements begin from index 1.

Index

Index of $Parent(i)$: $\lfloor \frac{i}{2} \rfloor$

Index of $Left_Child(i)$: $2 \times i$

Index of $Right_Child(i)$: $2 \times i + 1$

Note that a **sentinel** is placed in index 0, namely, a number that is smaller than all the elements in the heap (MinHeap) or larger than all the elements in the heap (MaxHeap). Usually -1 or $MAXN$.

Max and Min Heap

- A max heap is a complete binary tree and a max tree, which means that the key value of each node \geq the key value of its children.
- A min heap is a complete binary tree and a min tree, which means that the key value of each node \leq the key value of its children.

Percolate Up and Down (MinHeap)

Percolate up and down are frequently used operations in heaps. Percolate up handles insertions, while percolate down handles deletions.

- Percolate Up

```
1 void PercolateUp(int Heap[],int Location){
2     int temp=Heap[Location];
3     int i=Location;
4     while(Heap[i/2]>temp){
5         Heap[i]=Heap[i/2];
6         i=i/2;
7     }
8     Heap[i]=temp;
9 }
```

- Percolate Down

```
1 void PercolateDown(int Heap[],int Size,int Location){
2     int temp=Heap[Location];
3     int current,child;
4     current=Location;
5     child=2*current;
6     while(child<=Size){
7         if(child+1<=Size&&Heap[child+1]<Heap[child]){
8             child++;
9         }
10        if(Heap[child]<temp){
11            Heap[current]=Heap[child];
12            current=child;
13            child=2*current;
14        }else{
15            break;
16        }
17    }
18    Heap[current]=temp;
```

The time complexity of both Percolate Up and Percolate Down is $O(\log N)$.

Insert (MinHeap)

The key part of insertion is actually Percolate Up. We append the new element to the end of the heap and then adjust the heap to a MinHeap.

"Heap" is a structure that contains two parts: Size and Elements[].

```
1 void Insert(Heap H,int X){
2     H->Size++;
3     H->Elements[H->Size]=X;
4     PercolateUp(H->Elements,H->Size);
5 }
```

DeleteMin (MinHeap)

Similarly, the key part of deletion is Percolate Down. We replace the root node by the last element in heap, deleting the last element and adjusting the heap to a MinHeap.

```
1 int DeleteMin(Heap H){
2     int temp=H->Elements[1];
3     H->Elements[1]=H->Elements[H->Size];
4     H->Size--;
5     PercolateDown(H->Elements,H->Size,1);
6     return temp;
7 }
```

BuildHeap (MinHeap)

A naive way would be conducting N insertions. $T(N) = O(N \log N)$ this way, obviously not fast.

A clever way is firstly putting all the N elements into the heap ignoring the property of heap. Then for $i = \frac{N}{2}, i > 0$, $i --$, conducting PercolateDown(i). $T(N) = O(N)$ this way.

```
1 void BuildHeap(int A[],int N){
2     int i;
3     for(i=0;i<N;i++){
4         Heap[i+1]=A[i];
5     }
6     for(i=N/2;i>0;i--){
7         PercolateDown(Heap,N,i);
8     }
9 }
```

Proof of linear time complexity:

[Theorem] For the perfect binary tree of height h containing $2^{h+1} - 1$ nodes, the sum of the heights of the nodes is:

$$2^{h-1} \times 1 + 2^{h-2} \times 2 + \dots + 2^0 \times h = 2^h \left(\frac{1}{2^1} + \frac{2}{2^2} + \frac{3}{2^3} + \dots + \frac{h}{2^h} \right) = 2^{h+1} - 1 - (h + 1).$$

\therefore The time complexity equals to the sum of heights, namely, $O(N)$.

[Example] Given a list of N elements and an integer k , find **the $k - th$ largest element**.

This problem can be solved in numerous ways. One is building a MaxHeap and conducting k deletions, whose time complexity is $O(N + k \log N) = O(N)$.

Another more efficient way is using **QuickSelection**, a variator of QuickSort. The time complexity will be $O(N)$, namely, linear.

d -Heap

A d -heap is a heap whose nodes have d children.

DeleteMin will take $d - 1$ comparisons to find the smallest child. Hence the total time complexity would be $O(d \log_d N)$.

8 Disjoint Set

Relation

A relation R is defined on a set S if for every pair of elements (a, b) , $a, b \in S$, $a R b$ is either true or false. If $a R b$ is true, then we say that a is related to b .

A relation, \sim , over a set, S , is said to be an equivalence relation over S if and only if it is symmetric, reflexive and transitive over S . If x and y are in the same equivalence class, then $x \sim y$.

Operations

[Example] $S_1 = \{6, 7, 8, 10\}$, $S_2 = \{1, 4, 9\}$, $S_3 = \{2, 3, 5\}$

Note that all the sets are **disjoint** ($S_i \cap S_j = \emptyset$).

In the array representation, $A[7] = 6$, $A[8] = 6$, $A[10] = 6$, $A[6] = -4$, etc.

Common operations:

$Union(i, j) ::=$ Replace S_i and S_j by $S = S_i \cup S_j$.

$Find(i) ::=$ Find the set S_k which contains the element i and return the index of the root.

Array Implementation

- $S[Element] =$ the element's parent.
- $S[root] = 0$ or $-x$ (x means the number of elements in S).

Union

The idea is to make S_i a subtree of S_j , or vice versa.

Basic version:

```
1 void Union(int set[],int num1,int num2){
2     int root1,root2;
3     root1=Find(Set,num1);
4     root2=Find(Set,num2);
5     if(root1==root2){
6         return;
7     }
8     set[root2]=root1;
9 }
```

Notice that in the worst case, the set will deteriorate to a linked list, $O(N) = N$ this way.

So we have two strategies: Union by size and Union by rank (height).

Union by Size

```
1 void Union(int set[],int num1,int num2){
2     int root1,root2;
3     root1=Find(Set,num1);
4     root2=Find(Set,num2);
5     if(root1==root2){
6         return;
7     }
8     if(set[root1]<set[root2]){
9         set[root1]=set[root1]+set[root2];
10        set[root2]=root1;
11    }else{
12        set[root2]=set[root1]+set[root2];
13        set[root1]=root2;
14    }
15 }
```

In the process, $set[root]$ is a negative number, whose absolute value denotes the number of elements in the set.

If $set[root1] < set[root2]$, then $set1$ has more elements than $set2$, so we union the smaller set $set2$ into $set1$ by setting $set[root2] = root1$.

- Let T be a tree created by union-by-size with N nodes, then:

$$height(T) \leq \lfloor \log_2 N \rfloor + 1$$

- Time complexity of N union and M find operations is $O(N + M \log_2 N)$.

Union by Rank (Height)

```
1 void Union(int set[],int rank[],int num1,int num2){
2     int root1=Find(set,num1);
3     int root2=Find(set,num2);
4     if(root1==root2){
5         return;
6     }
7     if(rank[root1]<rank[root2]){
```

```

8         set[root1]=root2;
9     }else if(rank[root1]>rank[root2]){
10         set[root2]=root1;
11     }else{
12         set[root1]=root2;
13         rank[root2]++;
14     }
15 }

```

Similarly, we always change the shallow tree.

Find

Continue to find the current element's parent until root is found.

Basic version:

```

1 int Find(int set[],int num){
2     while(set[num]>0){
3         num=set[num];
4     }
5     return num;
6 }

```

However, if the size of S is extremely large, find operation will take much time. So **Path Compression** is needed to make a sequence of find operations faster.

The idea is to let all the elements from current node to root node point to root node. So the height of the set will be reduced.

```

1 int Find(int set[],int num){
2     int root,current,parent;
3     root=num;
4     while(root>0){
5         root=set[root];
6     }
7     current=num;
8     while(current!=root){
9         parent=set[current];
10        set[current]=root;
11        current=parent;
12    }
13    return root;
14 }

```

Let $T(M, N)$ be the maximum time required to process an intermixed sequence of $M \geq N$ finds and $N - 1$ unions, then:

$$k_1 M \alpha(M, N) \leq T(M, N) \leq k_2 M \alpha(M, N), k_1, k_2 \text{ are constants, } \alpha(M, N) \leq 4.$$

9 Graph

Definitions and terminologies

Self loop and Multigraph are not considered.

$G(V, E)$	Undirected graph(v_i, v_j)	Directed graph (Digraph) $< v_i, v_j >$
Complete Graph	Adjacent (to/from) Incident on	Degree (in/out): $\sum degree = 2 \times E$
Subgraph	Path: a set of vertices	Length of path (number of edges)
Simple path (No cycle)	v_i and v_j are connected	G is connected
Connected components	DAG (Directed acyclic graph)	Strongly connected components (SCC): Exist path from v_i to v_j and v_j to v_i

Representation of Graph

Adjacency Matrix

- $A[x][y] = 1$ if $(v_i, v_j) \in E(G)$, 0 otherwise.
- If G is undirected, the adjacency matrix is symmetric. We can store half of the matrix to save space. $A = \{a_{11}, a_{21}, a_{22}, a_{31}, a_{32}, a_{33}, \dots, a_{n1}, \dots, a_{nn}\}$, the index for a_{ij} is $\frac{i \times (i-1)}{2} + j$.
- Adjacency matrix may waste space if the graph is sparse. So we choose adjacency matrix when the graph has many edges.

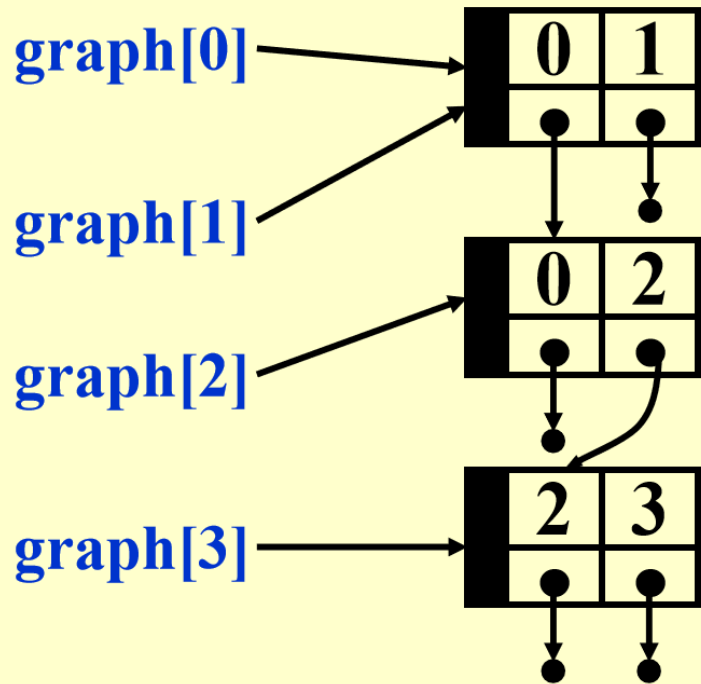
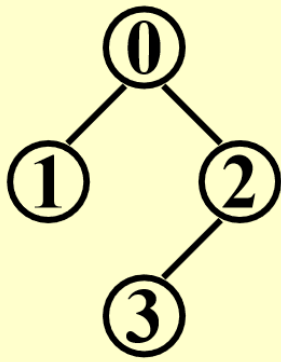
Adjacency List

- Good when the graph is sparse.
- Replace each row by a linked list that stores its adjacent nodes.

Adjacency Multilist

- Sometimes we need to mark the edge after examine it and then find the next edge. Adjacency multilist makes it easy to do so.
- Focus on the edge and let vertices points to the edge.

【 Example 】



Topological Sort

[Example] Courses with prerequisites need to be arranged in a university. As an illustration, "Programming 1" and "Discrete Mathematics" are required before learning "Data Structure", so a student is supposed to take the two courses in semester 1 and "Data Structure" in semester 2.

Topological sort is based on an AOV (Activity on Vertex) Network, with vertices representing activities (courses) and edges representing precedence relations.

An topological order (not unique) is a linear ordering of vertices of a graph such that for any two vertices i, j , if i is a predecessor of j in the network then i precedes j in the linear ordering.

Implementation:

```

1 void Topsort(Graph G){
2     Queue Q;
3     int count=0;
4     Vertex V,W;
5     for(each vertex V){
6         if(Indegree[V]==0){
7             Enqueue(V,Q);
8         }
9     }
10    while(Q is not empty){
11        V=Dequeue(Q);
12        TopNum[V]=count;
13        count++;
14        for(All the adjacent vertices W of V){
15            if(--Indegree[W]==0){
16                Enqueue(W,Q);
17            }
18        }
19    }
20    if(count!=NumberOfVertex){

```



```

21         Error("Graph has a cycle");
22     }
23 }

```

The data structure used in the algorithm is a queue, which stores the vertices of degree 0. When a vertex is dequeued and marked, the indegree of all the adjacent vertices will -1.

Time complexity equals to $O(V + E)$.

Unweighted Shortest Path (BFS)

As shown in the title, the idea is Breadth-First-Search (BFS), which means we begin from a particular vertex and scan through the graph level by level.

```

1  void BFS(Graph G,Vertex S){
2      Queue Q;
3      Vertex V,W;
4      Enqueue(S,Q);
5      distance[S]=0;
6      while(Q is not empty){
7          V=Dequeue(Q);
8          for(All the adjacent vertices W of V){
9              if(distance[W]==Infinity){
10                 distance[W]=distance[V]+1;
11                 Enqueue(W,Q);
12             }
13         }
14     }
15 }

```

The algorithm is quite similar to topological sort. Both need a queue as a container and involves a procedure of checking a vertex's adjacent vertices. Time complexity equals to $O(V + E)$.

Weighted Shortest Path (Dijkstra's Algorithm)

Dijkstra's algorithm cannot handle graphs with negative cost edge.

Three arrays need to be maintained in Dijkstra's algorithm: *Known*[], *Distance*[] and *Path*[].

Here is the implementation of the algorithm:

```

1  void Dijkstra(Graph G,Vertex S){
2      for(i=0;i<NumV;i++){
3          Known[i]=0;
4          Distance[i]=Infinity;
5          Path[i]=0;
6      }
7      Known[S]=1;
8      Distance[S]=0;
9      while(1){
10         V=unknown vertex with smallest distance;
11         Known[V]=1;
12         for(i=0;i<NumV;i++){

```

```

13         if(G[V][i]!=0&&Known[i]==0&&Distance[i]<Distance[V]+G[V][i]){
14             Distance[i]=Distance[V]+G[V][i];
15             Path[i]=V;
16         }
17     }
18 }
19 }

```

We can divide the algorithm into 4 steps:

1. Find the unknown vertex with smallest distance.
2. Mark the vertex as "Known".
3. Find the vertex's all the adjacent vertices that have not been visited.
4. Compare the original distance and new distance, update it if necessary.

For step 1, to find the vertex with smallest distance, there are 2 implementations:

- Simply scan the graph: $O(V)$, overall time complexity $O(V^2 + E)$.
- Keep distances in a heap and call DeleteMin: $O(\log V)$. Considering the adjustment of the heap after conducting $Distance[i] = Distance[V] + G[V][i]$, the overall time complexity equals to $O(V \log V + E \log V) = O(E \log V)$.

Graph with Negative Edge Cost

A special type of graph. Negative edge cost means that Dijkstra's algorithm can no longer be used, instead we keep a queue as the container and **visit each edge ≥ 1 time** ($Known[]$ is meaningless).

```

1 void Negative(Graph G,Vertex S){
2     Queue Q;
3     Vertex V,W;
4     Enqueue(S,Q);
5     Inqueue[S]=1;
6     while(Q is not empty){
7         V=Dequeue(Q);
8         Inqueue[V]=0;
9         for(All the adjacent vertices W of V){
10             if(distance[V]+G[V][W]<distance[W]){
11                 distance[W]=distance[V]+G[V][W];
12                 if(Inqueue[W]==0){
13                     Enqueue(W,Q);
14                     Inqueue[W]=1;
15                 }
16             }
17         }
18     }
19 }

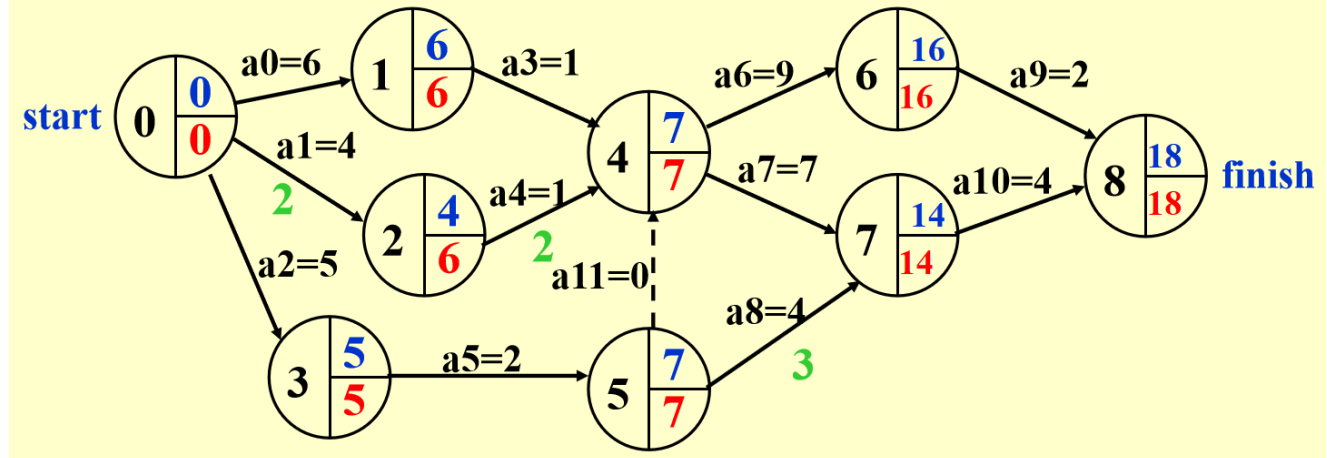
```

AOE (Activity on Edge) Network

AOE network imitates the process of scheduling a project. We will discuss the Earliest Completion time (EC), Latest Completion time (LC) and Critical Path Method (CPM).

- Calculation of EC: Start from V_0 , for any $a_i = \langle v, w \rangle$, we have $EC[w] = \max\{EC[v] + Cost_{\langle v, w \rangle}\}$.
- Calculation of LC: Start from last vertex, for any $a_i = \langle v, w \rangle$, we have $LC[v] = \min\{LC[w] - Cost_{\langle v, w \rangle}\}$.
- Slack time of $\langle v, w \rangle = LC[w] - EC[v] - Cost_{\langle v, w \rangle}$.
- Critical Path ::= the path consisting entirely of zero-slack edges.

【Example】 AOE network of a hypothetical project



Network Flow Problem

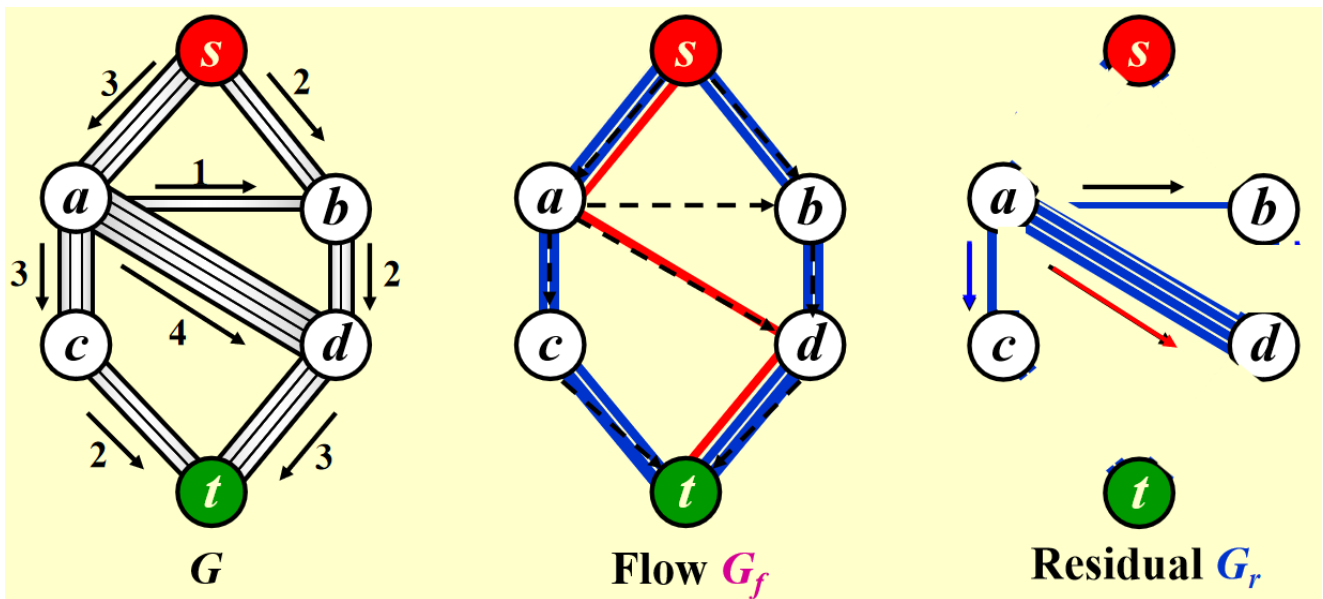
Given a network of pipes, we need to determine the maximum amount of flow that can pass from s (source) to t (sink).

G is the original graph, Flow G_f is the graph represents the result, while G_r is the graph that shows the remaining flow.

Simple algorithm

We can divide the algorithm into 4 parts:

1. Find any path from s to t in G .
2. Take the minimum edge on this path as the amount of flow and add to G_f .
3. Update G_r and remove the 0 flow edge.
4. If there is still a path from s to t , go back to step 1.

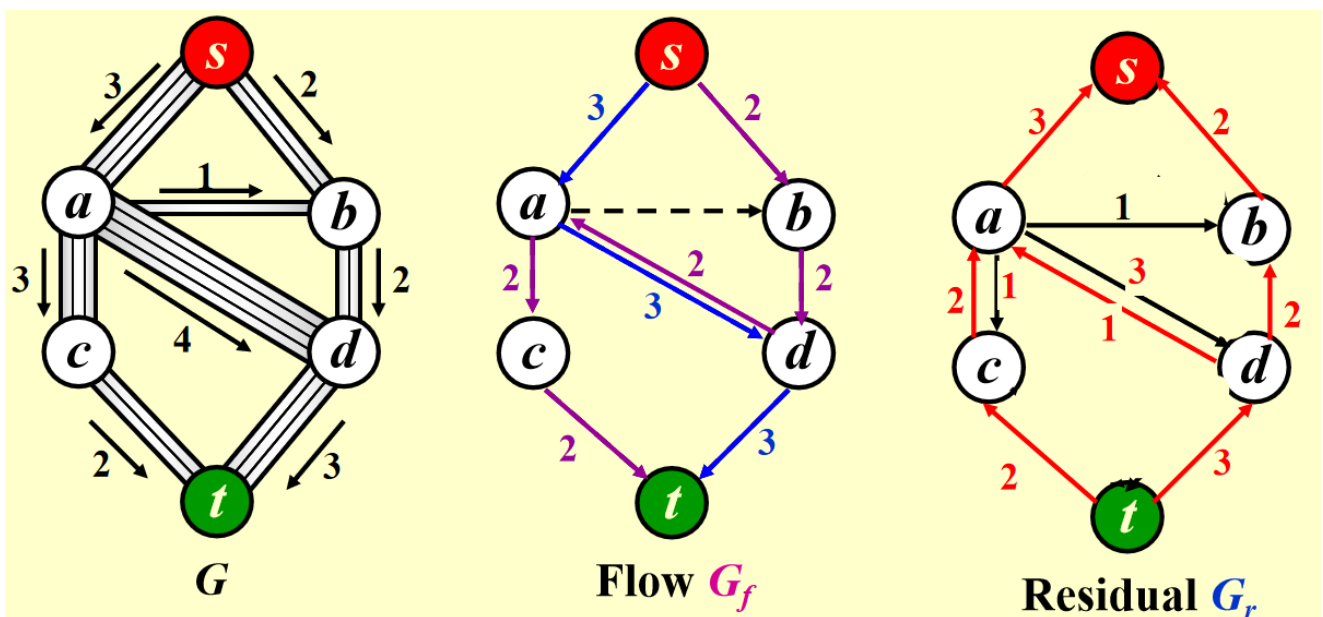


We should always pick the smallest edge first (CANNOT be greedy), and the algorithm does not always produce the correct answer.

Ford-Fulkerson's Algorithm

For each edge (v, w) with flow $f_{v,w}$ in G_f , add an edge (w, v) with flow $f_{v,w}$ in G_r .

This algorithm allows the computer to **undo** its decisions.



Analysis

An augmenting path is a path from source to sink.

An augmenting path can be found by an unweighted shortest path algorithm.

Always choose the augmenting path that allows the largest increase in flow.

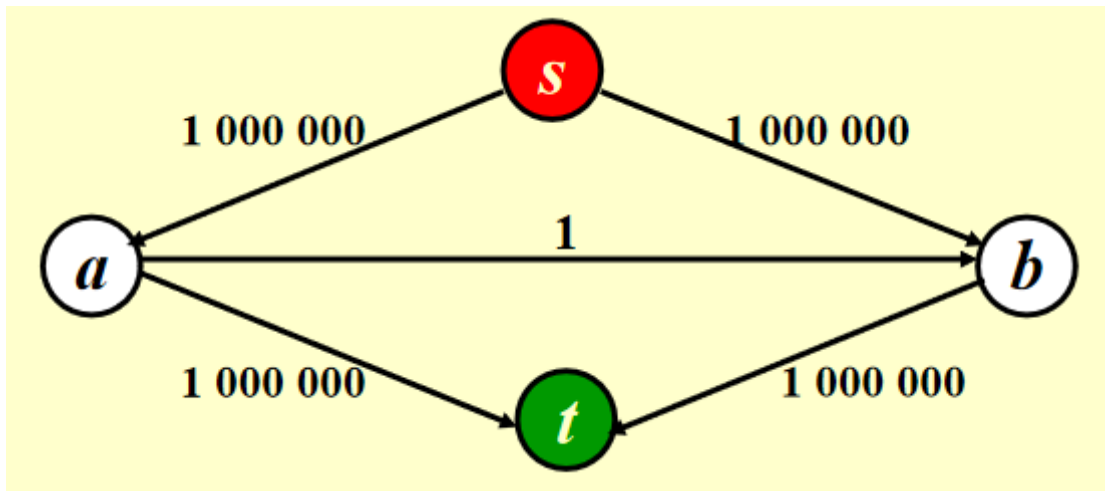
If we always choose the augmenting path with the largest flow, then:

$$T(N) = T_{\text{augmentation}} \times T_{\text{find a path}} = O(E \log \text{cap}_{\max}) \times O(E \log V) = O(E^2 \log V)$$

If we always choose the augmenting path that has the least number of edges, then: (BFS is used to find the path with least number of edges)

$$T(N) = T_{\text{augmentation}} \times T_{\text{find a path}} = O(E) \times O(EV) = O(E^2V)$$

A Special Case



In this case, the maximum flow is obviously 2000000. However, if we use the simple algorithm to find the maximum flow, we will get a wrong answer of 1999999. If we use Ford-Fulkerson's algorithm to find the maximum flow but find a path with minimum flow each time, we have to do it 2000000 times!

Minimum Spanning Tree

A spanning tree of a graph G is a tree which consists of $V(G)$ and a subset of $E(G)$.

- Minimum: the sum of weight of all edges is minimized.
- Tree: acyclic, $E = V - 1$.
- Spanning: the tree contains every vertex in the original graph.

There are two algorithms to find the minimum spanning tree, Prim's Algorithm and Kruskal's Algorithm. Both are greedy algorithms because they find the edge with minimum weight at each step.

Prim's Algorithm

Here I provide a pseudocode:

```

1  Tree Prim(Graph G){
2      Tree={};
3      Find the edge(V,W) with smallest weight;
4      Add vertex V,W and edge(V,W) to the tree;
5      for(i=1;i<=NumV-2;i++){
6          Find the minimum cost edge(X,Y) deriving from the vertices of the current tree that
          does not form a cycle;
7          Add vertex Y and edge(X,Y) to the tree;
8      }
9      return Tree;
10 }
  
```

To be more specific, at each step we find a edge satisfying the following 3 requirements:

1. One vertex of the edge is in the current tree.
2. After adding the edge to the tree, there is no cycle in the tree.
3. The edge is a minimum cost edge.

We can use an analogy to describe the algorithm: **Let a small tree grow up.**

Kruskal's Algorithm

Disjoint set ADT may be used in this algorithm to determine whether there is a cycle.

```
1 Tree Kruskal(Graph G){
2     Tree={};
3     for(i=1;i<=NumV-1;i++){
4         Find a minimum cost edge(V,W) in G that does not form a cycle after added to Tree;
5         Add V,W(if necessary) and edge(V,W) to the tree;
6     }
7     return Tree;
8 }
```

Q: How to determine whether there is a cycle after adding edge(V, W)?

A: We use disjoint set ADT to handle the problem. Each time we add an edge to Tree, we Union the two vertices. If $Find(V)$ and $Find(W)$ shows that V and W are in the same equivalence class, then adding edge(V, W) will produce a cycle.

Similarly an analogy can be used to describe the algorithm: **Union trees into forests.**

Depth First Search (DFS)

Depth First Search is a traversal of all vertices that makes depth its priority. We search deeper and deeper until we cannot go further. Then we return to the previous vertices (retrospect) and recursively conduct the procedure.

```
1 void DFS(Graph G,Vertex S){
2     visited[S]=true;
3     print(S);
4     for(All the adjacent vertices W of S){
5         if(visited[W]==false){
6             DFS(G,W);
7         }
8     }
9 }
```

The algorithm is simple and easy to understand. We can regard this algorithm as a generalization of preorder traversal. Recursions are used to conduct DFS.

There is a lot of applications of DFS, here I list two important ones: Biconnectivity and Euler Circuit.

Biconnectivity of Undirected Graph

- Articulation point: after deleting the vertex, there are at least 2 connected components in the new graph.

A graph G is biconnected if G is connected and has no articulation points. Similarly we can define a biconnected component.

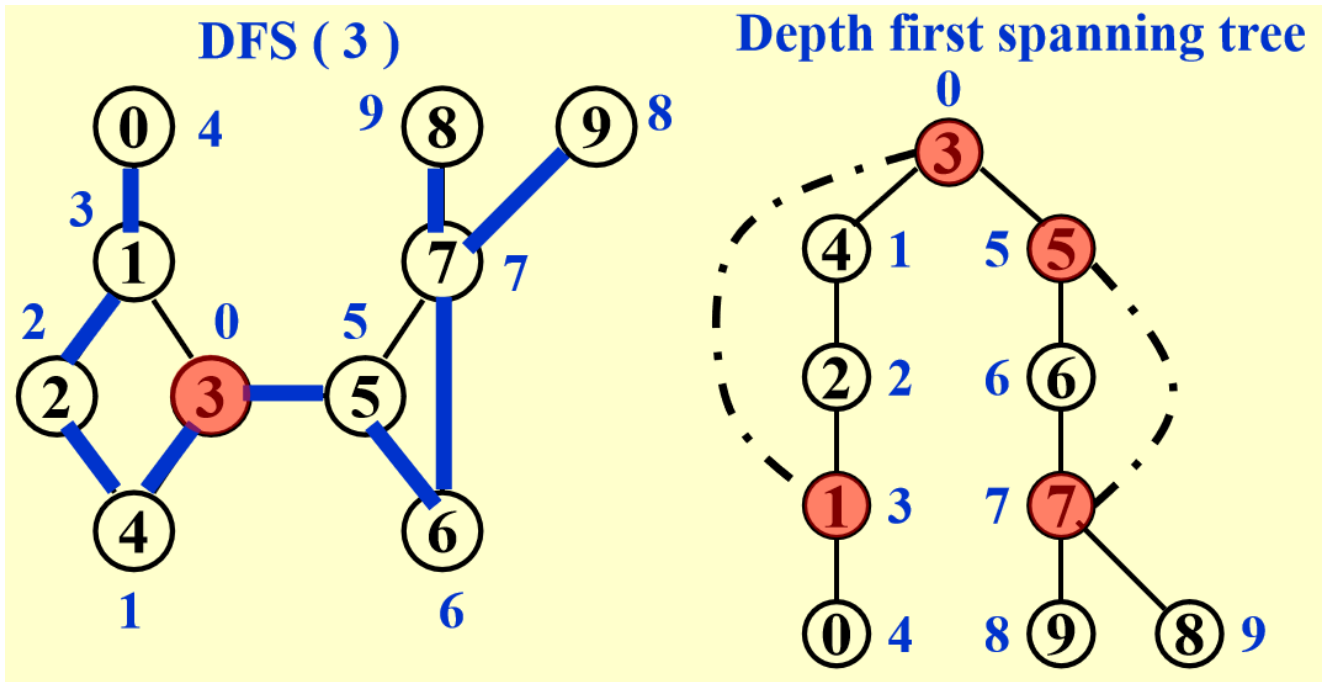
We can divide the algorithm into 3 steps:

1. Use DFS to obtain a spanning tree of G . Each vertex is assigned a number that denotes the order we visit it.
2. Construct a DFS spanning tree and add **back edge** to the tree. Note that if A is the parent of B , then $DFSNum(A) < DFSNum(B)$.

3. Determine articulation point:

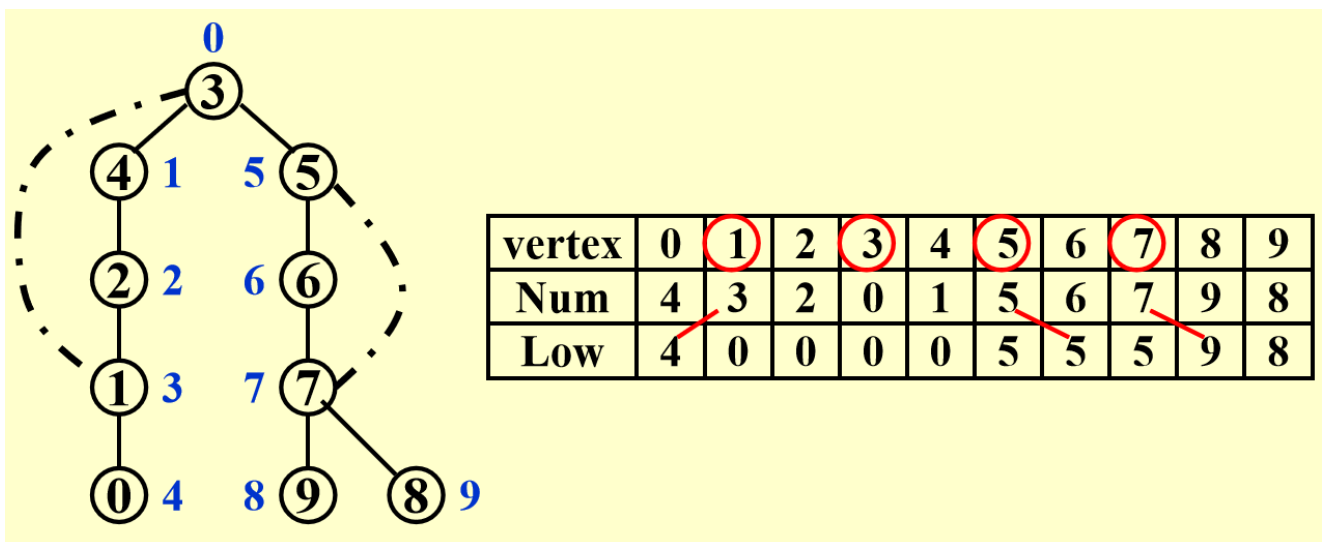
- If the root has at least 2 children, then it is an articulation point.
- Leaf nodes are not articulation point.
- For other nodes, if it is possible to move down at least 1 step and jump to its ancestor, then it is not an articulation point.

Here is a detailed example:



However, such method is not feasible on computers. So we devise the algorithm a little:

1. For each vertex, assign $Num(v)$ and $Low(v)$ to it. $Num(v)$ equals to the DFS number, $Low(v) = \min\{Num(v), \min\{Low(w) | w \text{ is a child of } v\}, \min\{Num(w) | (u, w) \text{ is a back edge}\}\}$.
2. u is an articulation point iff:
 - u is the root and has at least 2 children.
 - u has a **child** such that $Low(child) \geq Num(u)$.



Euler Circuit

- Euler Circuit: a path that traverses through each edge exactly once, and finish at the starting point.
- Euler Path: a path that traverses through each edge exactly once.

An Euler Circuit is possible iff all the vertices have an even degree.

An Euler Path is possible iff 2 vertices have odd degree, while other vertices have even degree.

Hamilton cycle: a cycle that visits every **vertex** exactly once.

Strongly Connected Components (SCC, Tarjan's Algorithm)

To be supplemented.....

10 Sorting

Function interface: **void X_sort(int A[],int N);**

For simplicity, we only consider internal sorting and ascending sorting algorithms.

Insertion Sort

```
1 void Insertion_Sort(int A[],int N){
2     int i,j,temp;
3     for(i=1;i<N;i++){
4         temp=A[i];
5         for(j=i;j>0&&A[j-1]>temp;j--){
6             A[j]=A[j-1];
7         }
8         A[j]=temp;
9     }
10 }
```

The algorithm is simple. We can use an analogy to describe it: Assume that you are playing cards. When you have a coming card, you need to insert the new card into your cards. So you just keep searching backward until you find a card whose value is smaller than the new card. Then it is the proper position to conduct insertion.

Time complexity analysis:

- Worst case (reversed order): $T(N) = O(N^2)$.
- Best case (sorted order): $T(N) = O(N)$.

An inversion in an array is any ordered pair (i, j) such that $i < j$ but $A[i] > A[j]$.

Each swap of adjacent elements eliminates exactly 1 inversion. Thus $T(N) = O(Inversions + N)$.

There are $\frac{N(N-1)}{4}$ inversions in an array of N elements on average. Thus any algorithm sorts by exchanging adjacent elements has a time complexity of $O(N^2)$ on average.

Bubble Sort and Selection Sort

They are so simple that I just list the code without further explanation.

```
1 void Bubble_Sort(int A[],int N){
2     int i,j;
3     for(i=N-2;i>=0;i--){
4         for(j=0;j<=i;j++){
5             if(A[j]>A[j+1]){
6                 swap(&A[j],&A[j+1]);
7             }
8         }
9     }
10 }
```

```
1 void Selection_Sort(int A[],int N){
2     int i,j,MinIndex;
3     for(i=0;i<N;i++){
4         MinIndex=i;
5         for(j=i+1;j<N;j++){
6             if(A[j]<A[MinIndex]){
7                 MinIndex=j;
8             }
9         }
10        swap(&A[i],&A[MinIndex]);
11    }
12 }
```

Shell Sort

Shell sort is a generalization of insertion sort. In shell sort, an increment sequence $h_1 < h_2 < \dots < h_t$ is defined and we conduct h_k -Insertion_Sort each time.

```
1 void Shell_Sort(int A[],int N){
2     int i,j,k,Increment,temp;
3     for(k=X,k>=1;k--){
4         Increment=h[k];
5         for(i=Increment;i<N;i++){
6             temp=A[i];
7             for(j=i;j>=Increment&&A[j-Increment]>temp;j=j-Increment){
8                 A[j]=A[j-Increment];
9             }
10            A[j]=temp;
11        }
12    }
13 }
```

Worst case: If Shell's increment sequence $\{1, 2, 4, 8, 16, 32, \dots\}$ is used, for the array $A = [1, 9, 2, 10, 3, 11, 4, 12, 5, 13, 6, 14, 7, 15, 8, 16]$, 8-sort, 4-sort, 2-sort will be useless. Only 1-sort can sort the array.

So choosing an appropriate increment sequence is important.

- Hibbard's Increment Sequence: $h_k = 2^k - 1$, worst-case time complexity $T(N) = \Theta(N^{\frac{3}{2}})$, average time complexity $T(N) = O(N^{\frac{5}{4}})$.
- Sedgewick's Increment Sequence (Best known increment sequence): $h = \{1, 5, 19, 41, 109, \dots\}$, $h_k = 9 \times 4^i - 9 \times 2^i + 1$ for h_1, h_3, \dots , $h_k = 4^i - 3 \times 2^i + 1$ for h_2, h_4, \dots . $T_{avg}(N) = O(N^{\frac{7}{6}})$ and $T_{worst}(N) = O(N^{\frac{4}{3}})$.

Shell sort is good for sorting small-scale input, medium-scale input and up to moderately large input (tens of thousands).

Heap Sort

A MaxHeap is used to conduct heap sort. First we use the linear-time-complexity algorithm to build a maxheap, and conduct $N - 1$ DeleteMax operations. Note that the deleted element is placed at the end of the array.

```

1 void Heap_Sort(int A[],int N){
2     int i;
3     for(i=N/2;i>=0;i--){
4         PercolateDown(A,i,N);
5     }
6     for(i=N-1;i>=1;i--){
7         swap(&A[0],&A[i]);
8         PercolateDown(A,0,i);
9     }
10 }
```

Heap sort sorts the elements from end to beginning.

Although heap sort gives the best average time, in practice it is slower than shell sort that used Sedgewick's increment sequence because there are so many operations involve addresses. Heap sort may be used to find the k -th largest number.

Merge Sort

The key idea is to recursively merge two sorted lists.

```

1 void Merge_Sort(int A[],int N){
2     int* TmpArray;
3     TmpArray=(int*)malloc(sizeof(int)*N);
4     Msort(A,TmpArray,0,N-1);
5 }
6
7 void Msort(int A[],int TmpArray[],int left,int right){
8     int center;
9     if(left<right){
10         center=(left+right)/2;
11         Msort(A,TmpArray,left,center);
12         Msort(A,TmpArray,center+1,right);
13         Merge(A,TmpArray,left,center+1,right);
14     }
15 }
16
17 void Merge(int A[],int TmpArray[],int LeftPtr,int RightPtr,int RightEnd){
18     int i,LeftEnd,num,ptr;
19     LeftEnd=RightPtr-1;
20     ptr=LeftPtr;
```

```

21     num=RightEnd-LeftPtr+1;
22     while(LeftPtr<=LeftEnd&&RightPtr<=RightEnd){
23         if(A[LeftPtr]<A[RightPtr]){
24             TmpArray[ptr++]=A[LeftPtr++];
25         }else{
26             TmpArray[ptr++]=A[RightPtr++];
27         }
28     }
29     while(LeftPtr<=LeftEnd){
30         TmpArray[ptr++]=A[LeftPtr++];
31     }
32     while(RightPtr<=RightEnd){
33         TmpArray[ptr++]=A[RightPtr++];
34     }
35     for(i=0;i<num;i++,RightEnd--){
36         A[RightEnd]=TmpArray[RightEnd];
37     }
38 }

```

Time complexity analysis:

$$T(1) = 1, T(N) = 2T\left(\frac{N}{2}\right) + O(N) = O(N + N \log N) = O(N \log N).$$

However, merge sort requires linear extra memory, and copying an array is slow, so it is seldom used in internal sorting, but is quite useful for external sorting.

Replacement Selection

To be supplemented.....

Quick Sort

Quick sort is the fastest known sorting algorithm in practice.

The algorithm can be divided into 3 steps:

1. Picking a pivot in $A[]$.
2. Partition the set into 2 parts, $A_1 = \{a \in S | a \leq pivot\}$, $A_2 = \{a \in S | a \geq pivot\}$.
3. Recursively conduct quick sort in A_1 and A_2 .

The pivot is placed at the right place once and for all. That is, each run there is at least 1 element sorted.

- The strategy to pick the pivot: Median-of-Three Partitioning: $pivot = median(left, center, right)$.
- Partitioning strategy: i points to the beginning of the array, j points to the end of the array, while $A[i] < pivot$ conducting $i++$, while $A[j] > pivot$ conducting $j--$ until $i > j$, then the set is correctly partitioned. Note that if $A[i] = pivot$ we still need to stop to guarantee the set is **equally-divided**.
- If N is small, insertion sort is faster than quick sort. So we set a cutoff to decide which sorting algorithm we use.

```

1  int MedianThree(int A[],int left,int right){
2      int center=(left+right)/2;
3      if(A[left]>A[center]){
4          swap(&A[left],&A[center]);
5      }
6      if(A[left]>A[right]){

```

```

7         swap(&A[left],&A[right]);
8     }
9     if(A[center]>A[right]){
10         swap(&A[center],&A[right]);
11     }
12     swap(&A[center],&A[right]);
13     return A[right];
14 }
15
16 void Quick_Sort(int A[],int N){
17     Qsort(A,0,N-1);
18 }
19
20 void Qsort(int A[],int left,int right){
21     if(right-left+1<cutoff){
22         Insertion_Sort(A+left,right-left+1);
23     }else{
24         int i,j,pivot;
25         pivot=MedianThree(A,left,right);
26         i=left;
27         j=right;
28         while(1){
29             while(A[++i]<pivot);
30             while(A[--j]>pivot);
31             if(i<j){
32                 swap(&A[i],&A[j]);
33             }else{
34                 break;
35             }
36         }
37         swap(&A[i],&A[right]);
38         Qsort(A,left,i-1);
39         Qsort(A,i+1,right);
40     }
41 }

```

Time complexity analysis:

$$T(N) = T(i) + T(N - i - 1) + O(N).$$

In the worst case, the set is divided into two parts containing 0 and $N - 1$ elements, $T(N) = T(N - 1) + O(N) = O(N^2)$.

In the best case, the set is **equally divided**, $T(N) = 2T(\frac{N}{2}) + O(N) = O(N \log N)$. So does the average case.

Application

Given a list of N elements and an integer k , find the k -th largest element.

Recall that pivot is placed at the right position once and for all. Devise the algorithm:

```

1     .....
2     if(k<N-i){
3         Qselect(A,N,i+1,right,k);
4     }else if(k>N-i){
5         Qselect(A,N,left,i-1,k);
6     }
7     .....

```

We only need to focus on one side of the array. $T(N) = T(\frac{N}{2}) + O(N) = O(N)$, the fastest algorithm to solve k -th largest problem.

Sorting Large Structures (Table Sort)

For large structures, swapping will be very expensive. So we add a pointer to the structure and swap pointers instead. Rearrange the structures at last.

An array $table[]$ is needed to store pointers. Each time we compare $list[table[i]]$ and $list[table[j]]$ and swap $table[i]$ and $table[j]$. The sorted list is $list[table[0]], list[table[1]], \dots, list[table[N-1]]$.

list	[0]	[1]	[2]	[3]	[4]	[5]
key	d	b	f	c	a	e
table	4	1	3	0	5	2

Then place the structure at the right place by conducting $list[i] = list[table[i]]$. $T = O(mN)$ where m is the size of a structure.

Theorem: Any algorithm that sorts by comparisons only must have a worst case time complexity of $O(N \log N)$.

Proof: For an array of N elements, there are $N!$ possible permutations, which is the number of leaves in the result tree. Recall that $l \leq 2^h$, $h \geq \log l = \log N!$. Time complexity equals $h = O(\log N!) = \Theta(N \log N)$.

Bucket Sort and Radix Sort

Bucket Sort

If we have N students, whose scores are in the range 0 to 100, then we can sort their scores in linear time:

First we create 101 buckets $bucket[0]$, $bucket[1]$, \dots , $bucket[100]$. Then we read in each student's score X and insert it to $bucket[X]$. Lastly we scan through the buckets and output the list.

$T = O(M + N)$, namely, linear. However, to use bucket sort we need to know the distribution and range in advance, and it cannot handle float numbers.

Radix Sort

[Example] Given 10 integers in the range of 0 to 999, sort them in linear time.

Solution: Sort according to the Least Significant Digit first (LSD) algorithm. That is, first put the number to the bucket according to its last digit, and then second last, third last, \dots

The following picture shows the procedure to sort 64, 8, 216, 512, 27, 729, 0, 1, 343, 125:

0	1	2	3	4	5	6	7	8	9
0	1	512	343	64	125	216	27	8	729
0	512	125		343		64			
1	216	27							
8		729							
0	125	216	343		512		729		
1									
8									
27									
64									

$T(N) = O(P(N + B))$, where P is the number of passes, N is the number of elements, B is the number of buckets.

MSD (Most Significant Digit) and LSD (Least Significant Digit)

Suppose that the record R_i has r keys.

- K_j^i denotes the j -th key of record R_i .
- K_0^i is the most significant digit of R_i .
- K_{r-1}^i is the least significant digit of R_i .

[Example] For a deck of cards, K_0 denotes the suit, $Club < Diamond < Heart < Spade$ K_1 denotes the face value, $2 < 3 < 4 < 5 < 6 < 7 < 8 < 9 < 10 < J < Q < K < A$.

- MSD: First sort on K_0 , create 4 buckets for the suits, then sort according to face value in each bucket.
- LSD: First sort on K_1 , create 13 buckets for face values, then sort according to suit in each bucket.

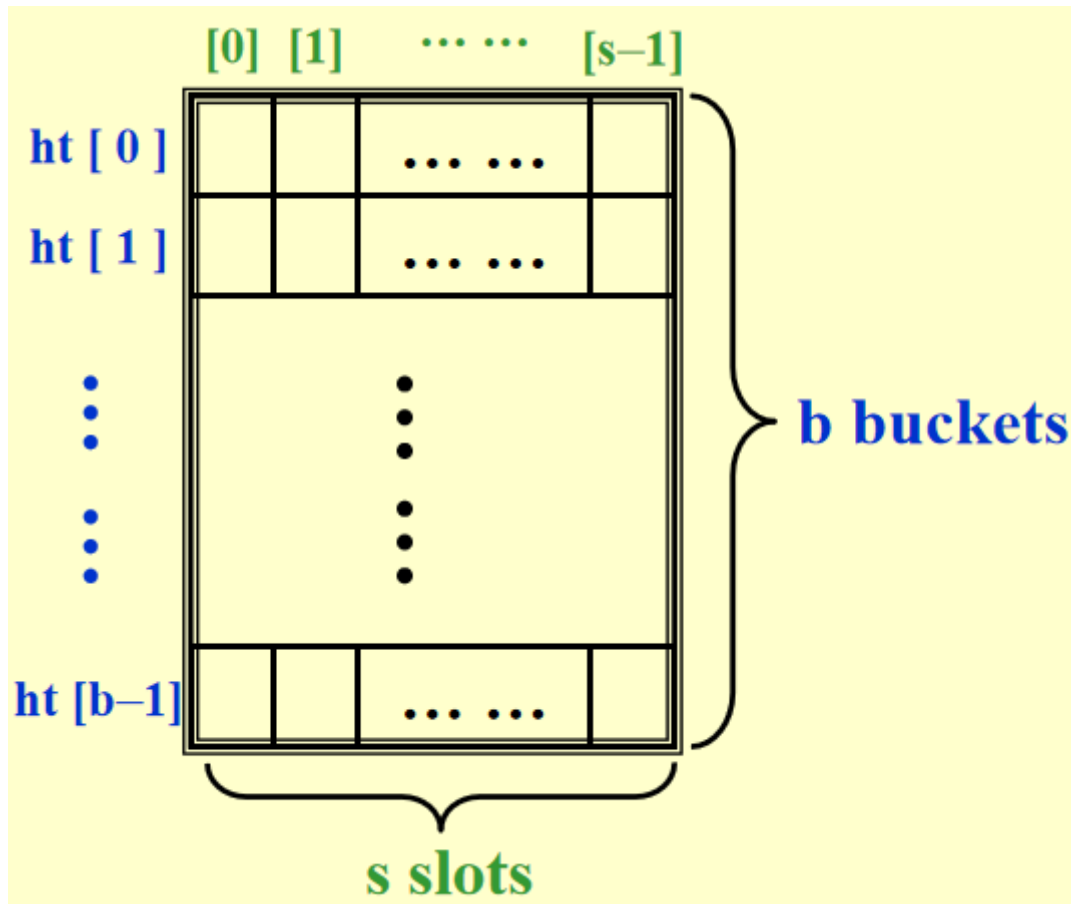
11 Hashing

Although we have learned sorting algorithms with time complexity $O(N \log N)$ and binary search with time complexity $O(\log N)$, it is still very slow to do a search. That is the problem that hashing solves.

Hashing is a method to make find, delete and insert easier (in constant time).

- General idea: create a symbol table (like dictionary in Python) of $\langle name, attribute \rangle$.
- Key operations: Find, Insert and Delete.
- The most significant part of hashing: Find a **hash function** $f(x)$ whose input is the identifier x and output is the index in the hash table.
- Identifier density $:= \frac{n}{T}$, n is the total number of identifiers, T is the total number of distinct possible values for x .
- Loading density $\lambda := \frac{n}{sb}$, s is the number of slots, b is the number of buckets in hash table.

Here is a hash table:



A collision occurs when $i_1 \neq i_2$ but $f(i_1) = f(i_2)$.

An overflow occurs when we hash a new identifier into a full bucket. If $s = 1$, then collision and overflow occurs simultaneously.

Hash Function

- f must be easy to compute and minimizes the number of collisions.
- f should be unbiased and uniform. That is, $P(f(x) = i) = \frac{1}{b}$.

TableSize is supposed to be **prime** to minimize collisions.

An possible hash function for a string: $f(x) = \sum str[N - i - 1] \times 32^i$.

```
1 int Hash(char* str,int TableSize){
2     int HashValue=0;
3     while(*str!='\0'){
4         HashValue=(HashValue<<5)+*x;
5         x++;
6     }
7     return HashValue%TableSize;
8 }
```

As is known to all, collisions are bad for hash table. We have two methods to solve collisions: Separate Chaining and Open Addressing.

Separate Chaining

The idea is to keep a list of all key that are hashed to the same value.

Structures and typedef:

```
1 typedef struct ListNode* List;
2 struct ListNode{
3     int element;
4     List next;
5 }
6 typedef struct HashTbl* HashTable;
7 struct HashTbl{
8     int TableSize;
9     List* TheLists;
10 }
```

TheLists is an array whose elements are pointers that points to the node that contains elements.

Here are some commonly used operations:

- Create an empty table:

```
1 HashTable CreateTable(int TableSize){
2     HashTable H;
3     int i;
4     H=(HashTable)malloc(sizeof(struct HashTbl));
5     H->TableSize=NextPrime(TableSize);
6     H->TheLists=malloc(sizeof(List)*H->TableSize);
7     for(i=0;i<H->TableSize;i++){
8         H->TheLists[i]=malloc(sizeof(struct ListNode));
9         H->TheLists[i]->next=NULL;
10    }
11    return H;
12 }
```

- Find a key from the hash table:

```
1 List Find(int key,HashTable H){
2     List ptr;
3     List L;
4     L=H->TheLists[Hash(key,H->TableSize)];
5     ptr=L->next;
6     while(ptr!=NULL&&ptr->element!=key){
7         ptr=ptr->next;
8     }
9     return ptr;
10 }
```

- Insert a key into a hash table:


```

1 void Insert(int key, HashTable H){
2     List ptr, NewCell;
3     List L;
4     ptr = Find(key, H);
5     if(ptr == NULL){
6         NewCell = malloc(sizeof(struct ListNode));
7         L = H->TheLists[Hash(key, H->TableSize)];
8         NewCell->next = L->next;
9         NewCell->element = key;
10        L->next = NewCell;
11    }
12 }

```

These algorithms are explicit and easy to understand. So no further explanation is provided.

Open Addressing

In extreme scenarios, separate chaining will produce a super long linked list and a lot of empty lists. Besides, there are so many operations involving pointers in separate chaining, which makes it slow and inefficient. Different from separate chaining, open addressing algorithm finds another empty cell to solve collision.

Overall algorithm:

```

1 int OpenAddressing(int key, HashTable H){
2     int index = Hash(key);
3     int i = 0;
4     if(H is full){
5         ERROR("No space left");
6     }
7     while(Collision at index){
8         i++;
9         index = (hash(key) + f(i)) % H->TableSize;
10    }
11    return index;
12 }

```

There are different $f(x)$ in different algorithms. All $f(x)$ satisfies that $f(0) = 0$.

Linear Probing

$f(i) = i$ in linear probing algorithm.

The expected number of probes $p = \frac{1}{2}(1 + \frac{1}{(1-\lambda)^2})$ for insertions and unsuccessful searches, $\frac{1}{2}(1 + \frac{1}{1-\lambda})$ for successful searches.

[Example] Mapping 11 strings into the hash table with $f(str) = str[0] - 'a'$. Strings: acos, atoi, char, define, exp, ceil, cos, float, atol, floor, ctime.

As we can verify, all the elements are crowded around index 0 to 10. If we need to map a new string "abc" the search time will be 12, which is extremely large.

Linear probing will cause **primary clustering**: any key that hashes into the cluster will add to the cluster after several attempts to resolve the collision, making searching extremely slow.

Quadratic Probing

$f(i) = i^2$ in quadratic probing.

Here are 2 important theorems:

- ┆ If the table size is prime and the table is at least half empty, then a new element can always be inserted.
- ┆ If the table size of prime and the form is $4k + 3$, then the quadratic probing with $f(i) = \pm i^2$ can probe the entire table.

Here are some commonly used operations:

- Find an element:

```
1  int Find(int key,HashTable H){
2      int index;
3      int i=0;
4      index=Hash(key,H->TableSize);
5      while(H->Array[index].info==legitimate&&H->Array[index].element!=key){
6          index=index+2*(++i)-1;
7          if(index>=H->TableSize){
8              index=index-TableSize;
9          }
10     }
11     return index;
12 }
```

- Note that $f(i) = i^2$, so $f(i) = f(i - 1) + 2i - 1$, which is the meaning of " $index = index + 2 * (+ + i) - 1$ ".
- A algorithm similar to "lazy deletions" is used in quadratic probing. That is, when deleting an element, we simply set $H \rightarrow Array[X].info = empty$ without actually deleting it. If info equals to *legitimate*, then there is an element in it. If we want to delete an element, just set the info to *empty*.
- Insert an element:

```
1  void Insert(int key,HashTable H){
2      int index;
3      index=Find(key,H);
4      if(H->Array[index].info==empty){
5          H->Array[index].info=legitimate;
6          H->Array[index].element=key;
7      }
8  }
```

┆ Although primary clustering is solved, secondary clustering occurs. Keys that hash to the same position will probe the same alternative cells.

Double Hashing

$f(i) = i \times hash_2(x)$ in double hashing, where $hash_2$ is another hash function.

┆ $Hash_2(x) = R - (x \% R)$ will work well, where R is a prime number smaller than TableSize.

Quadratic probing does not require the use of a second hash function and is thus likely to be simpler and faster in practice.

Rehashing

If there are more than half elements in the hash table and quadratic probing is used to resolve collisions, probing may fail. So rehashing is needed that builds another hash table that is twice as big and scan down the original table and use a new function to hash those elements into the new table.

Q: When to rehash?

A: 1. The table is half full. 2. When an insertion fails. 3. When the table reaches a particular loading factor.

| In an interactive system, the unfortunate user whose insertion causes a rehashing can witness an obvious slowdown.