

# 第二部分 分布式算法

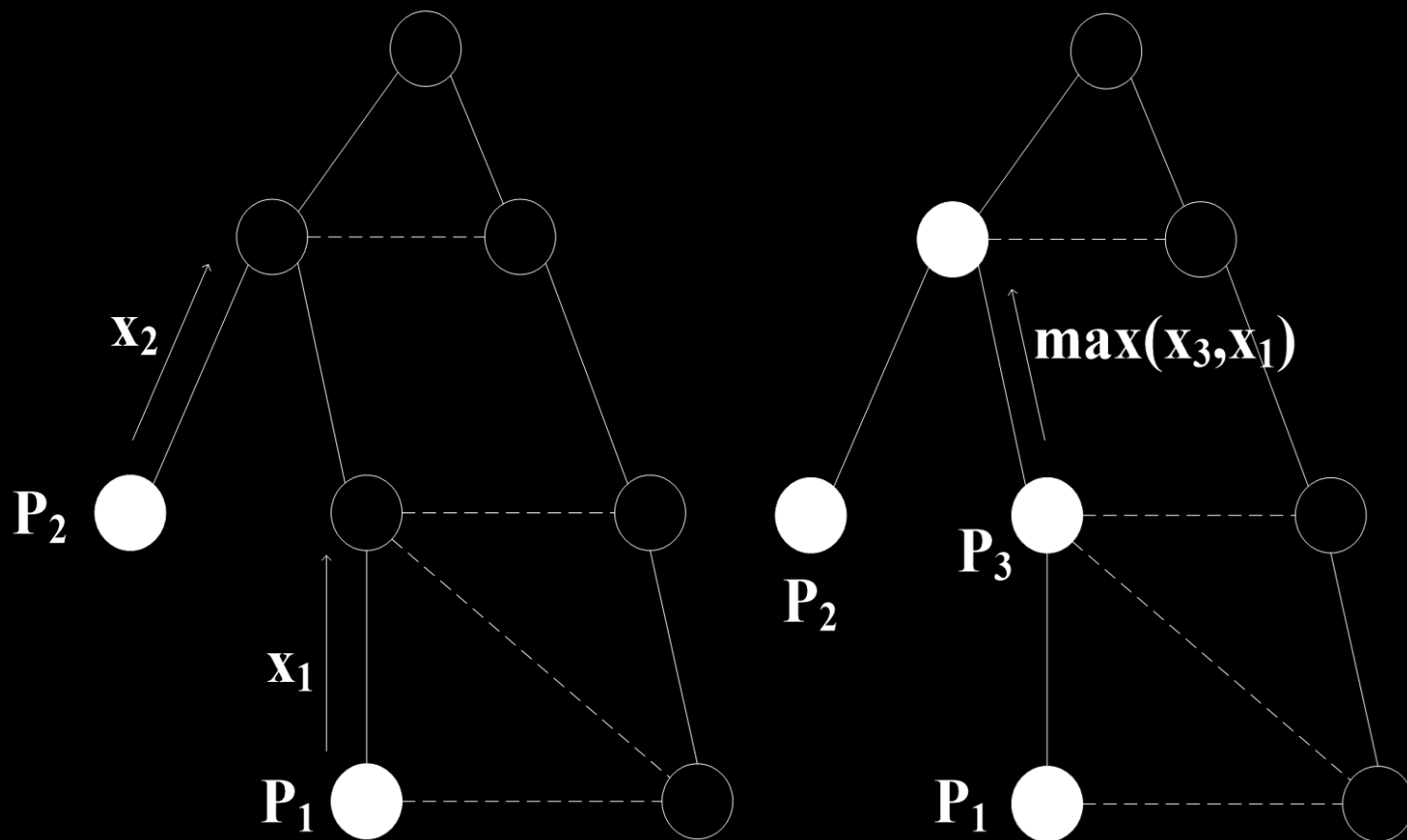
## 第三次课

中国科学技术大学软件学院

## § 2.2.2 convergecast(汇集, 敛播)

与广播问题相反, 汇集是从所有结点收集信息至根。  
为简单起见, 先考虑一个特殊的变种问题:

假定每个 $p_i$ 开始时有一初值 $x_i$ , 我们希望将这些值中最大者发送至根 $p_r$ 。



## § 2.2.2 convergecast(汇集, 敛播)

### ■ 算法

每个叶子结点 $p_i$ 发送 $x_i$ 至双亲。//启动者

对每个非叶结点 $p_j$ , 设 $p_j$ 有 $k$ 个孩子 $p_{i_1}, \dots, p_{i_k}$ ,  $p_j$ 等待 $k$ 个孩子的msg  $v_{i_1}, v_{i_2}, \dots, v_{i_k}$ , 当 $p_j$ 收到所有孩子的msg之后将 $v_j = \max\{x_j, v_{i_1}, \dots, v_{i_k}\}$ 发送到 $p_j$ 的双亲。

换言之: 叶子先启动, 每个处理器 $p_i$ 计算以自己为根的子树里的最大值 $v_i$ , 将 $v_i$ 发送给 $p_i$ 的双亲。

### ■ 复杂性

**Th2.5** 当生成树高为 $d$ 时, 存在一个异步的敛播方法, 其msg复杂性为 $n-1$ , 时间复杂度为 $d$ 。(与Th2.2相同)

■ **广播和敛播算法用途:** 初始化某一信息请求(广播发布), 然后用敛播响应信息至根。

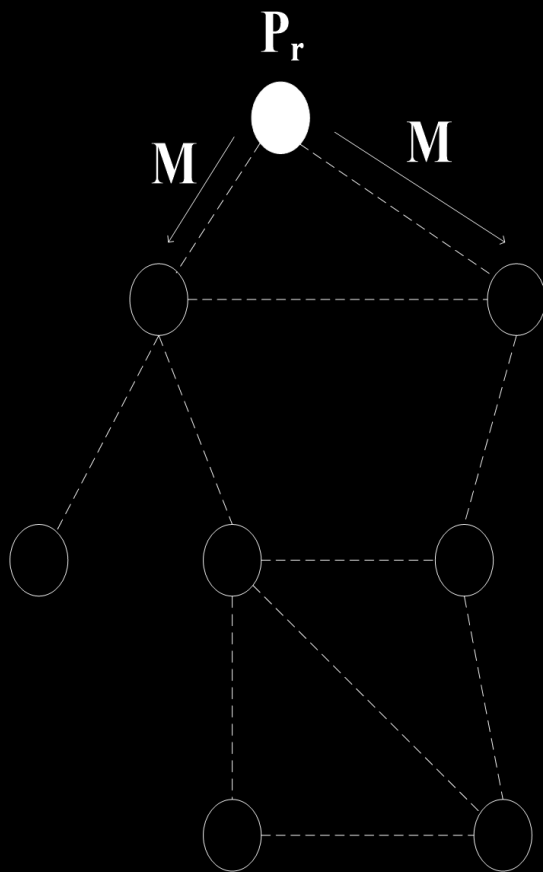
## § 2.3 构造生成树

上节算法均假设通信网的生成树已知。本节介绍生成树的构造问题。

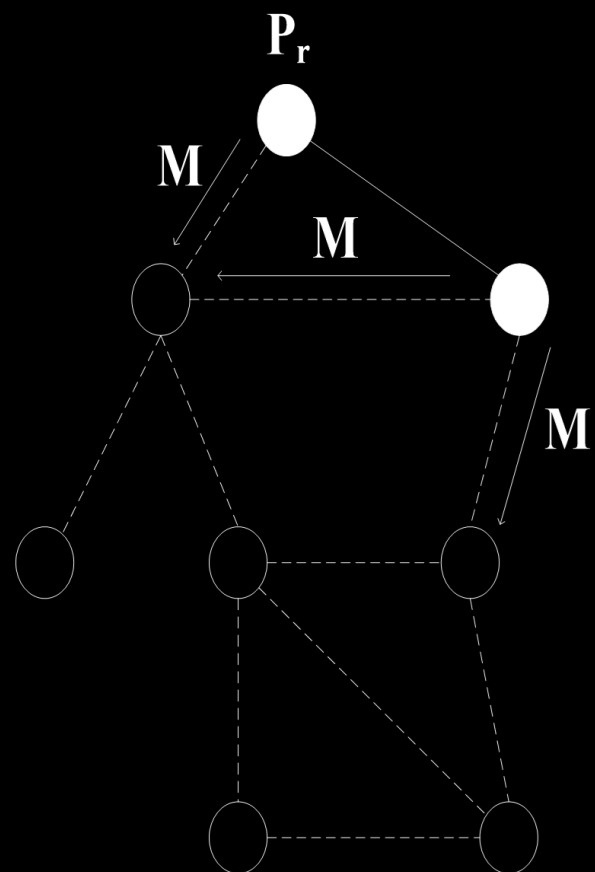
### 1. Flooding算法(淹没)

#### ■ 算法

设  $p_r$  是特殊处理器。从  $p_r$  开始，发送  $M$  到其所有邻居。当  $p_i$  第1次收到消息  $M$ （不妨设此 msg 来自于邻居  $p_j$ ）时， $p_i$  发送  $M$  到除  $p_j$  外的所有邻居。



(a)



(b)

## § 2.3 构造生成树

### ■ msg复杂性

因为每个结点在任一信道上发送M不会多于1次，所以每个信道上M至多被发送两次(使用该信道的每个处理器各1次)。

在最坏情况下：M除第1次接收的那些信道外，所有其他信道上M被传送2次。因此，有可能要发送 $2m-(n-1)$ 个msgs。这里m是系统中信道总数，它可能多达 $n(n-1)/2$ 。

### ■ 时间复杂性： $O(D)$ D—网直径

## 2.构造生成树

对于flooding稍事修改即可得到求生成树的方法。

## § 2.3 构造生成树

### ① 基本思想

- 首先， $p_r$ 发送M给所有邻居， $p_r$ 为根
- 当 $p_i$ 从某邻居 $p_j$ 收到的M是第1个来自邻居的msg时， $p_j$ 是 $p_i$ 的双亲；若 $p_i$ 首次收到的M同时来自多个邻居，则用一个comp事件处理自上一comp事件以来的所有已收到的msgs，故此时， $p_i$ 可在这些邻居中任选一个邻居 $p_j$ 做双亲。
- 当 $p_i$ 确定双亲是 $p_j$ 时，发送<parent>给 $p_j$ ，并向此后收到发来M的处理器发送<reject>msg

## § 2.3 构造生成树

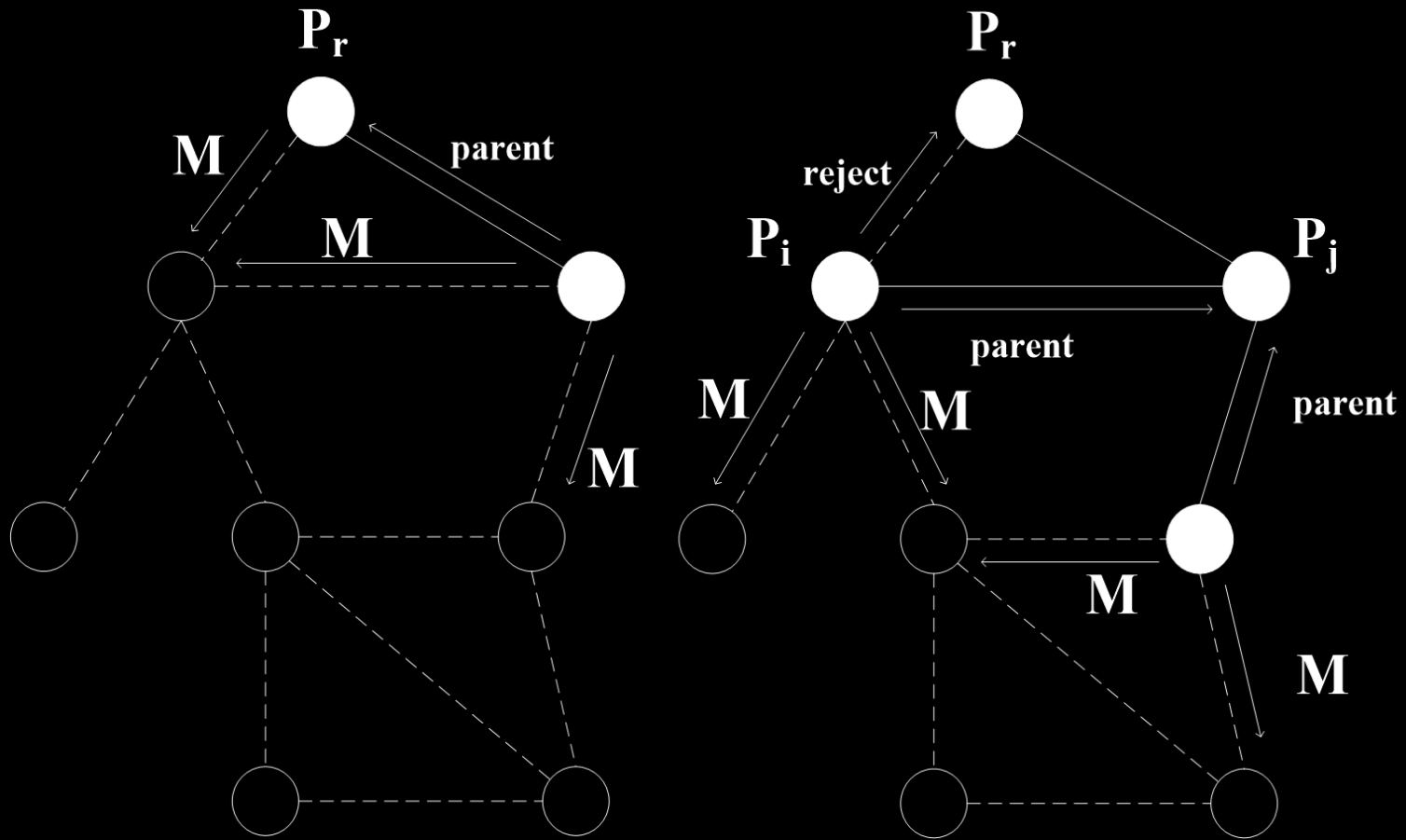
### ① 基本思想

- 因为 $p_i$ 收到 $p_j$ 的M是第1个M，就不可能已收到其他结点的M，当然可能同时收到(说明 $p_i$ 与这些邻居间不是父子关系，或说它们不是生成树中的边)；同时 $p_i$ 将M转发给其余邻居，这些邻居尚未发M给 $p_i$ ，或虽然已发M给 $p_i$ ，但 $p_i$ 尚未收到。
- $p_i$ 向那些尚未发M给 $p_i$ (或已发M但尚未到达 $p_i$ )的邻居转发M之后，等待这些邻居发回响应msg: <parent> 或<reject>。那些回应<parent>的邻居是 $p_i$ 的孩子。
- 当 $p_i$ 发出M的所有接收者均已回应 (<parent> 或<reject>), 则 $p_i$ 终止。将parent和children边保留即为生成树。

## § 2.3 构造生成树

### ②图示

$p_i$ 若认为 $p_j$ 是其双亲，则 $p_i$ 向 $p_r$ 发出M，而 $p_r$ 仍会向 $p_i$ 发reject，但因为此前 $p_r$ 向 $p_i$ 发出过M，故 $p_i$ 收到M时仍会向 $p_r$ 发reject。(可以改进?)





## § 2.3 构造生成树

③算法： Alg2.2 构造生成树 (code for  $p_i$   $0 \leq i \leq n-1$ )

初值： parent=nil; 集合children和other均为 $\emptyset$

1. upon receiving no message:
2.     if  $i=r$  and parent=nil then { //根尚未发送M
3.         send M to all neighbors;
4.         parent:=i;} //根的双亲置为自己
5. upon receiving M from neighbor  $p_j$ :
6.     if parent=nil then {// $p_i$ 此前未收到过M, M是 $p_i$ 收到的第1个msg
7.         parent:=j;
8.         send <parent> to  $p_j$ ; //  $p_j$ 是 $p_i$ 的双亲
9.         send M to all neighbors except  $p_j$ ;
10.         }else //  $p_j$ 不可能是 $p_i$ 的双亲,  $p_i$ 收到的M不是第1个msg
11.         send <reject> to  $p_j$ ;
12. upon receiving <parent> from neighbor  $p_j$ :
13.     children:=children  $\cup$  { j }; //  $p_j$ 是 $p_i$ 的孩子, 将j加入孩子集
14.     if children  $\cup$  other 包含了除parent外的所有邻居 then terminate;
15. upon receiving <reject> from neighbor  $p_j$ :
16.     other:=other  $\cup$  { j }; //将j加入other, 通过非树边发送的msg。
17.     if children  $\cup$  other包含了除 $p_i$ 的双亲外的所有邻居 then terminate

## § 2.3 构造生成树

### ④分析

**Lemma 2.6** 在异步模型的每个容许执行中，算法 2.2 构造一棵根为  $p_r$  的生成树。(正确性)

Pf: 算法代码告诉我们两个重要事实

- a) 一旦处理器设置了 `parent` 变量，它绝不改变，即它只有一个双亲
- b) 处理器的孩子集合决不会减小。

因此，最终由 `parent` 和 `children` 确定的图结构  $G$  是静止的，且 `parent` 和 `children` 变量在不同结点上是一致的，即若  $p_j$  是  $p_i$  的孩子，则  $p_i$  是  $p_j$  的双亲。

下述证明结果图  $G$  是根为  $p_r$  的有向生成树。

## § 2.3 构造生成树

### ■ 为何从根能到达每一结点？(连通)

**反证：**假设某结点在G中从 $p_r$ 不可达，因网络是连通的，若存在两个相邻的结点 $p_i$ 和 $p_j$ 使得 $p_j$ 在G中是从 $p_r$ 可达的(以下简称 $p_j$ 可达)，但 $p_i$ 不可达。因为G里一结点从 $p_r$ 可达当且仅当它曾设置过自己的parent变量(**Ex 证明**)，所以 $p_i$ 的parent变量在整个执行中仍为nil，而 $p_j$ 在某点上已设置过自己的parent变量，于是 $p_j$ 发送M到 $p_i$ (line9)，因该执行是容许的，此msg必定最终被 $p_i$ 接收，使 $p_i$ 将自己的parent变量设置为j。矛盾！

## § 2.3 构造生成树

### ■ 为何无环? (无环)

假设有一环,  $p_{i_1}, \dots, p_{i_k}, p_{i_1}$ , 若  $p_i$  是  $p_j$  的孩子, 则  $p_i$  在  $p_j$  第1次收到  $M$  之后第1次收到  $M$ 。因每个处理器在该环上是下一处理器的双亲, 这就意味着  $p_{i_1}$  在  $p_{i_1}$  第1次接收  $M$  之前第1次接收  $M$ 。矛盾!

### ■ 复杂性

显然, 此方法与淹没算法相比, 增加了  $\text{msg}$  复杂性, 但只是一个常数因子。在异步通信模型里, 易看到在时刻  $t$ , 消息  $M$  到达所有与  $p_r$  距离小于等于  $t$  的结点。因此有:

**Th2.7** 对于具有  $m$  条边和直径  $D$  的网络, 给定一特殊结点, 存在一个  $\text{msg}$  复杂性为  $O(m)$ , 时间复杂性为  $O(D)$  的异步算法找到该网络的一棵生成树。

## § 2.3 构造生成树

Alg2.2在同步模型下仍可工作。其分析类似于异步情形。然而，与异步不同的是，它所构造的生成树一定是一棵广度优先搜索(BFS)树。

**Lemma2.8** 在同步模型下，Alg2.2的每次容许执行均构造一棵根为 $p_r$ 的BFS树。

**Pf:** 对轮 $t$ 进行归纳。即要证明：在第 $t$ 轮开始时刻

①根据parent变量构造的图 $G$ 是一棵包括所有与 $p_r$ 距离至多为 $t-1$ 结点的BFS树；

②而传输中的消息 $M$ 仅来自于与 $p_r$ 距离恰为 $t-1$ 的结点。

由此构造的树是一棵根为 $p_r$ 的BFS

## § 2.3 构造生成树

当 $t=1$ 时，所有parent初值为nil，M从 $p_r$ 传出。

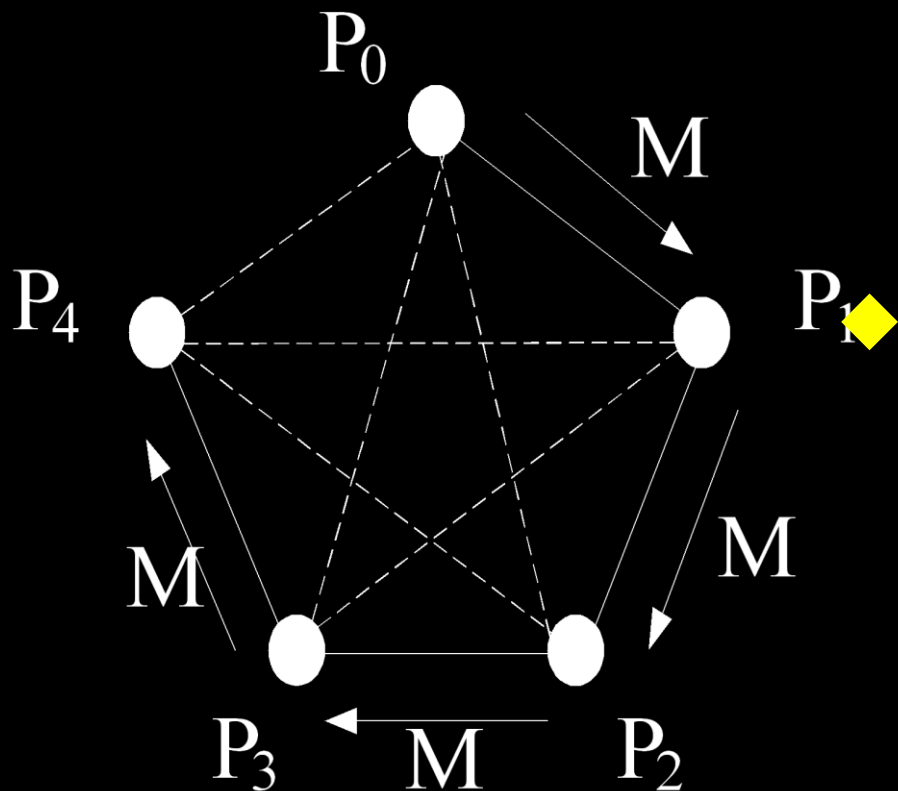
假设引理对第 $t-1 \geq 1$ 轮为真，在该轮里，从距离  $t-2$ 的结点传出的M被接收，任何接收M的结点与 $p_r$ 的距离不超过 $t-1$ (恰为 $t-1$ 或更短)，那些parent值非空的接收结点显然与 $p_r$ 的距离不超过 $t-2$ ，他们既不改变parent的值也不转发M；而与 $p_r$ 距离为 $t-1$ 的结点在 $t-1$ 轮里收到M，因为它们的parent为nil，故将其置为合适的双亲并转发M。距离 $p_r$ 大于 $t-1$ 的结点不会收到M，因此也不会转发M。因此有如下定理：

**Th2.9** 对于具有 $m$ 条边直径为 $D$ 的网络，给定一个特殊结点，存在一个同步算法在msg复杂性为 $O(m)$ ，时间复杂性为 $O(D)$ 内找到一棵BFS树。

## § 2.3 构造生成树

### ■ 异步系统里，Alg2.2能构造BFS树？

例如，考虑5个顶点的完全连通图



◆  $P_0$ 为根，假定M消息按 $P_0$ 到 $P_1$ ， $P_1$ 到 $P_2$ ， $P_2$ 到 $P_3$ ， $P_3$ 到 $P_4$ 的次序快速传播，而M在其它路径上传播较慢。结果生成树是从 $P_0$ 到 $P_4$ 的链，它不是BFS树。虽然此图直径 $D=1$ ，生成树的高度 $d=4$ ，但是算法的运行时间仍然为 $O(D)$ 而不是 $O(d)$ 。理解： $P_0$ 到 $P_4$ 的M在1个时间内到达，即 $P_0 \rightarrow P_1 \rightarrow P_2 \rightarrow P_3 \rightarrow P_4$ 的时间之和不超1。

## § 2.3 构造生成树

### ■ 信息的请求和收集

将算法2.2(求生成树)和汇集算法组合即可完成。组合算法的时间复杂性在同步和异步模型中不同,设网是**完全图**

❖ **求生成树算法**: 同步和异步均为  $\begin{cases} \text{消息复杂性 } O(m) \\ \text{时间复杂性 } O(D) \end{cases}$

❖ **汇集算法**: 同步和异步均有  $\begin{cases} \text{msg } n-1 \\ \text{time } d // \text{生成树高} \end{cases}$

### ❖ 组合算法

①同步: 组合算法的msg复杂性 $O(m+n)$ ; BFS树中,  $d=1, d \leq D$ , 故时间复杂性 $O(D+d)=O(D)=O(1)$ 。

②异步: 组合算法的msg复杂性 $O(m+n)$ ; **生成树高**  $d=n-1$ , 所以时间复杂性 $O(D+d)=O(d)=O(n)$ 。

1-time复杂性的组合算法 $T(n)=O(D)$ 。

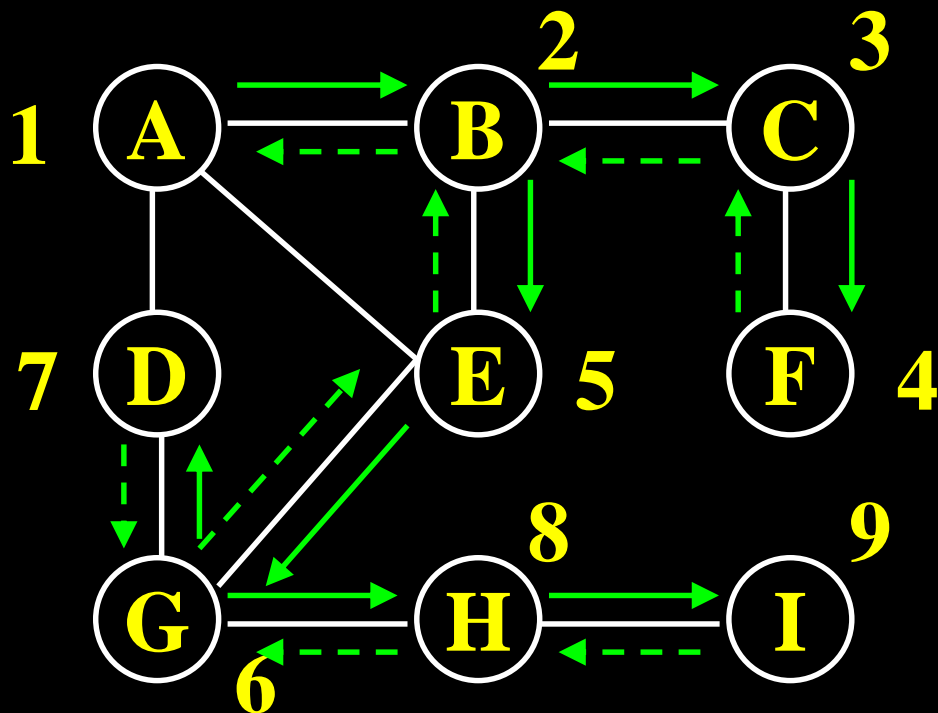


## § 2.4 构造DFS生成树(指定根)

回忆无向图的深度优先搜索问题：

■ 方法：

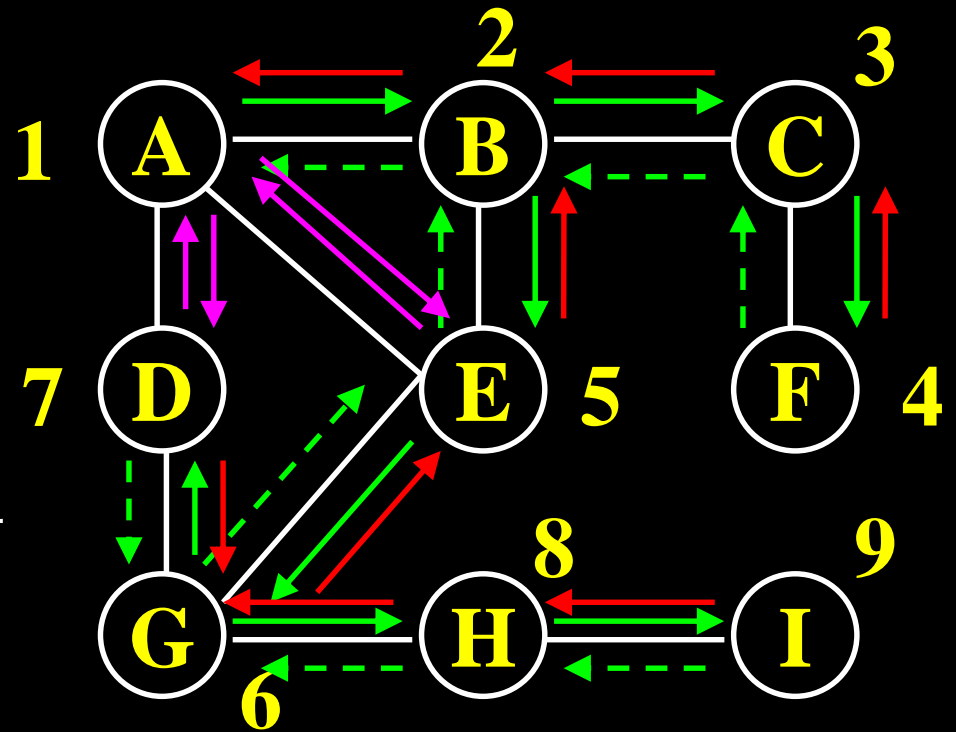
- ❖ 从任意顶点开始访问，例如A
- ❖ 然后访问它的一个没有访问过的邻接点，例如B
- ❖ 若无未访问过的邻接点,则后退寻找，直至全部被访问为止



## § 2.4 构造DFS生成树(指定根)

基本思想:

- 设定 $P_i$ 为指定的根节点,  $P_i$ 从还未向其发送消息 $\langle m \rangle$ 的邻接节点中任选一个节点发送消息 $\langle m \rangle$ 。
- 当 $P_i$ 从 $P_j$ 收到的消息 $\langle m \rangle$ 是第一个来自于邻接节点的消息时,  $P_j$ 为 $P_i$ 的双亲, 向其发送 $\langle \text{parent} \rangle$ 消息, 并向此后向自己发送消息的邻居发送 $\langle \text{reject} \rangle$ 消息。
- $P_i$ 从还未发送过消息的邻居中任选一个, 发送消息 $\langle m \rangle$ , 然后等待 $\langle \text{parent} \rangle$ 或者 $\langle \text{reject} \rangle$ 消息, 并将回应 $\langle \text{parent} \rangle$ 的节点加到自己的孩子集合中。
- 当 $P_i$ 向所有的邻居都转发过消息后,  $P_i$ 终止。



## § 2.4 构造DFS生成树(指定根)

构造DFS树时每次加一个结点，而不像Alg2.2那样，试图在树上同时增加同一层的所有结点。

### Alg2.3 构造DFS生成树，为 $P_r$ 为根

Code for processor  $P_i$ ,  $0 \leq i \leq n-1$

var parent: init nil;

children: init  $\phi$ ;

unexplored: init all the neighbors of  $P_i$   
//未访问过的邻居集

1: upon receiving no msg:

2: if  $(i=r)$  and  $(parent=nil)$  then { //当 $P_i$ 为根且未发送M时

3: parent :=  $i$ ; //将parent置为自身的标号

4:  $\forall P_j \in$  unexplored;

5: 将 $P_j$ 从unexplored中删去; //若 $P_r$ 是孤立结点,4-6应稍作修改

6: send M to  $P_j$ ;

}//endif

## § 2.4 构造DFS生成树(指定根)

```
7: upon receiving M from neighbor  $P_j$ :
8:   if parent=nil then { //  $P_i$ 此前未收到M
9:     parent :=  $j$ ; //  $P_j$ 是  $P_i$ 的父亲
10:    从unexplored中删  $P_j$ 
11:    if unexplored  $\neq \phi$  then {
12:       $\forall P_k \in$  unexplored;
13:      将  $P_k$ 从unexplored中删去;
14:      send M to  $P_k$ ;
15:    } else send <parent> to parent;
        //当  $P_i$ 的邻居均已访问过, 返回到父亲
16: }else send <reject> to  $P_j$ ; //当  $P_i$ 已访问过时
```

## § 2.4 构造DFS生成树(指定根)

```
17: upon receiving <parent> or <reject> from neighbor  $P_j$ ;  
18:   if received <parent> then add  $j$  to children;  
    //  $P_j$ 是 $P_i$ 的孩子  
19:   if unexplored =  $\phi$  then { //  $P_i$ 的邻居均已访问  
20:     if parent  $\neq i$  then send <parent> to parent;  
     //  $P_i$ 非根, 返回至双亲  
21:     terminate; // 以 $P_i$ 为根的DFS子树已构造好!  
22:   } else { // 选择 $P_i$ 的未访问过的邻居访问之  
23:      $\forall P_k \in$  unexplored;  
24:     将 $P_k$ 从unexplored中删去;  
25:     send  $M$  to  $P_k$ ;  
    }
```

## § 2.4 构造DFS生成树(指定根)

**引理2.10** 在异步模型里的每个容许执行, Alg2.3构造一棵以 $P_r$ 为根的DFS树。证明留作练习。

**Th2.11** 对于一个具有 $m$ 条边,  $n$ 个结点的网络, 以及给定的特殊顶点, 存在一个时间复杂性和消息复杂性均为 $O(m)$ 的异步算法找到一棵DFS树。

**Pf:** 每个结点在其邻接边上至多发送 $M$ 一次, 每个结点至多生成一个msg(<reject>或<parent>)作为对每个邻接边上收到的 $M$ 的响应。因此Alg2.3至多发送 $4m$ 个消息(其实大部分没有4倍), 即算法的msg复杂性为 $O(m)$ 。

时间复杂性证明留作练习。

如何改进使msg的复杂性不是 $4m$ ?

**注意:** 上述算法msg复杂性较好, 但时间复杂性太差。可降至 $O(n)$ 。

**Ex.**

- 2.1 分析在同步和异步模型下，convergecast算法的时间复杂性。
- 2.2 证明在引理2.6中，一个处理器在图G中是从 $P_r$ 可达的，当且仅当它的parent变量曾被赋过值
- 2.3 证明Alg2.3构造一棵以 $P_r$ 为根的DFS树。
- 2.4 证明Alg2.3的时间复杂性为 $O(m)$ 。
- 2.5 修改Alg2.3获得一新算法，使构造DFS树的时间复杂性为 $O(n)$ ，并证明。

**下次继续！**