

哈希与字符串

Trie

KMP算法

AC 自动机

失配指针

## 哈希与字符串

Hash 的核心思想在于，将输入映射到一个值域较小、可以方便比较的范围。

我们定义一个把字符串映射到整数的函数  $f$ ，这个  $f$  称为是 Hash 函数。

具体来说，我们希望在 Hash 函数值不一样的时候，两个字符串一定不一样。

通常我们采用的是多项式 Hash 的方法，即

$$f(s) = \sum s[i] \times b^i \pmod{M}$$

这里的  $b$  和  $M$  需要选取得足够合适才行，以使得 Hash 函数的值分布尽量均匀。

如果  $b$  和  $M$  互质，在输入随机的情况下，这个 Hash 函数在  $[0, M)$  上每个值概率相等。

令  $f_i(s)$  表示  $f(s[1..i])$ ，那么  $f(s[l..r]) = \frac{f_r(s) - f_{l-1}(s)}{b^{l-1}}$ ，其中  $\frac{1}{b^{l-1}}$  也可以预处理出来。

使用 unsigned long long 来表示  $f(s)$ ，unsigned long long 表示无符号数，也就是所表示的数全部为正数，这样在正数范围内表示的范围比 long long 更大。

### 1455: 【例题1】Oulipo

```
1  #include<bits/stdc++.h>
2  using namespace std;
3  int T;
4  int n,m;
5  int ans;
6  char s1[10000+10],s2[1000000+10];
7  unsigned long long power[1000000+10];
8  unsigned long long H[1000000+10];
9  unsigned long long s,b=27,h=1<<31;
10 int main(){
11     power[0]=1;
12     for(int i=1;i<=1000000;i++)
13         power[i]=power[i-1]*b;
14     scanf("%d",&T);
15     while(T--){
16     {
17         scanf("%s%s",s1+1,s2+1);
18         n=strlen(s1+1);m=strlen(s2+1);
19         H[0]=0;
20         for(int i=1;i<=m;i++)
21             H[i]=(H[i-1]*b+(unsigned long long)(s2[i]-'A'+1))%h;
```

```

22     s=0;
23     for(int i=1;i<=n;i++)
24         s=(s*b+(unsigned long long)(s1[i]-'A'+1))%h;
25     ans=0;
26     for(int i=0;i<=m-n;i++)
27         if(s==H[i+n]-H[i]*power[n])
28             ans++;
29     printf("%d\n",ans);
30 }
31 return 0;
32 }

```

#### 1456: 【例题2】图书管理

```

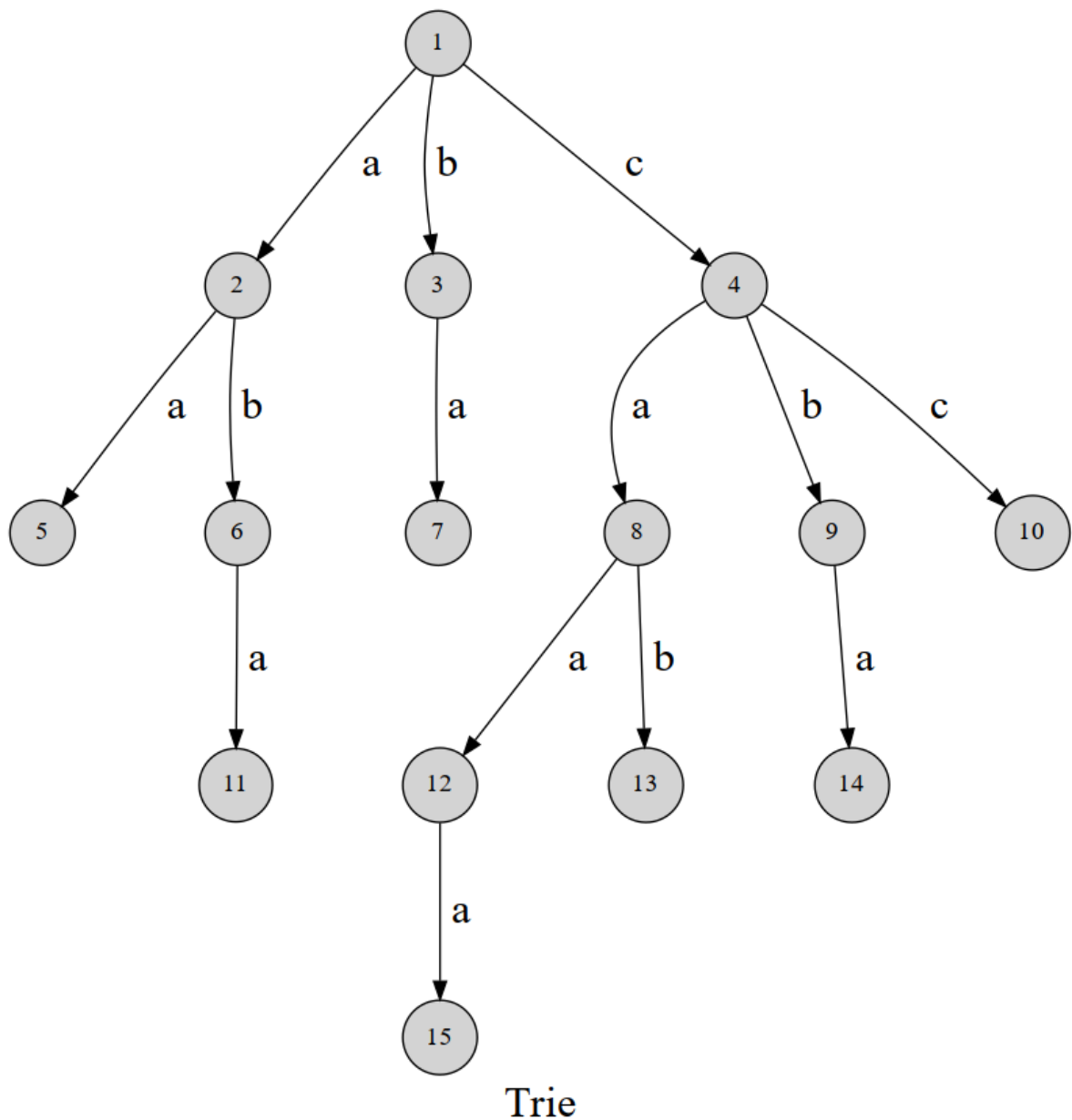
1  #include<bits/stdc++.h>
2  using namespace std;
3  typedef unsigned long long ull;
4  map<ull,int> mp;
5  ull b=27,mod=1<<31;
6  char op[205],str[205];
7  int main(){
8      int t;
9      cin>>t;
10     getchar();
11     while(t--){
12         scanf("%s",op);
13         cin.getline(str,205);
14         //gets(str);//必须用这个，不然最后一个测试点会超时
15         int len=strlen(str);
16         ull s=0;
17         for( int i=0; i<len; i++){
18             s=(s*b+(ull)(str[i]-'\\0'))%mod;
19         }
20         if(!strcmp(op,"find")){
21             if(mp[s]==1)
22                 printf("yes\n");
23             else
24                 printf("no\n");
25         }
26         else{
27             mp[s]=1;
28         }
29     }
30 }

```

遇到的问题：输入字符串时使用C语言语法才能通过测试用例。

## Trie

字典树，英文名 Trie。顾名思义，就是一个像字典一样的树。



可以发现，这棵字典树用边来代表字母，而从根结点到树上某一结点的路径就代表了一个字符串。

$1 \rightarrow 4 \rightarrow 8 \rightarrow 12$  表示的就是字符串 `caa`。

### 流程

#### 1. 初始化

一棵空Trie仅包含一个根节点，该点的字符指针均指向空。

```

1  int ch[N][Z], cnt;    //Z为字符集大小，N 是节点个数
2  bool val[N];
3  //若val=true则表示从根到该点经过的边上字母组成的字符串是实际字符串集合中的元素

```

2. 插入当需要插入一个字符串S时，我们令一个指针P起初指向根节点。然后，依次扫描S中的每个字符c：

a.若P的c字符指针指向一个已经存在的节点Q，则令P=Q。

a.若P的c字符指针指向空，则新建一个节点Q，令P的c字符指针指向Q，然后令P=Q。当S中的字符扫描完毕时，在当前节点P上标记它是一个字符串的末尾。

```

1 void insert(char *s) { //char *s表示一个字符数组
2     int len = strlen(s);
3     int u = 1; //1为根节点
4     for(int i = 0; i < len; ++i){
5         int c = s[i] - 'a';
6         if(!ch[u][c])
7             ch[u][c] = ++tot; //若不存在这条边则要新建一个节点与转移边
8         u = ch[u][c]; //tot为总点数
9     }
10    val[u] = true;
11    //在串的结尾处将val赋值，表示它代表一个实际字符串集合中的元素
12 }

```

3. 查询当需要查询一个字符串S在Trie中是否存在时，我们令一个指针P起初指向根节点，然后依次扫描S中的每个字符c：

a. 若P的c字符指针指向空，则说明S没有被插入过Trie，结束查询。

b. 若P的c字符指针指向一个已经存在的节点Q，则令P=Q。

当S中的字符扫描完毕时，若当前节点P被标记为一个字符串的末尾，则说明S在Trie中存在，否则说明S没有被插入过Trie。

```

1 bool find(char s[]) {
2     int len = strlen(s);
3     int u = 1;
4     for(int i = 0; i < len; ++i){
5         int c = s[i] - 'a';
6         if(!ch[u][c])
7             return false;
8         u = ch[u][c];
9     }
10    return true;
11 }

```

### 1471: 【例题1】Phone List

```

1 #include<bits/stdc++.h>
2 using namespace std;
3 const int maxn = 1e5+10;
4 int flag, t, n;
5 char str[15];
6 bool vis[maxn]; //vis是为了防止出现911, 91这种情况
7 int ch[maxn][15];
8 bool val[maxn];
9 int cnt;
10 void insert(char* s){
11     int len=strlen(s), u=0;
12     for( int i=0; i<len; i++){
13         int v=s[i]-'0';
14         if(!ch[u][v])
15             ch[u][v]=cnt++;
16         u=ch[u][v];
17         if(vis[u]&&i==len-1){
18             flag=1;
19             return;

```

```

20     }
21     vis[u]=true;
22     if(val[u]){
23         flag=1;
24         return;
25     }
26 }
27 val[u]=1;
28 }
29 int main(){
30     scanf("%d",&t);
31     while(t--){
32         flag=0;
33         memset(ch,0,sizeof(ch));
34         memset(vis,0,sizeof(vis));
35         memset(val,0,sizeof(val));
36         cnt=1;
37         scanf("%d",&n);
38         for( int i=0; i<n; i++){
39             scanf("%s",str);
40             insert(str);
41         }
42         if(flag)
43             printf("NO\n");
44         else
45             printf("YES\n");
46     }
47     return 0;
48 }

```

#### 1472: 【例题2】 The XOR Largest Pair

思路：对每一个数字 x 的二进制位进行查找，每一步都尽量找与当前位置相反的数字进行访问，如果有相反的数字，根据xor运算这一位就会留下一个 1，最后找出最大值就行。

```

1  #include <bits/stdc++.h>
2  using namespace std;
3  #define maxn 100005
4  int cnt;
5  int ch[maxn*32][2];
6  int n;
7  int a[maxn];
8  void insert(int x){
9     int u=0;
10    for(int i=31;i>=0;i--){
11        int id=(x>>i)&1;
12        if(!ch[u][id])
13            ch[u][id]=++cnt;
14        u=ch[u][id];
15    }
16 }
17 int find(int x){
18     int u=0;
19     int ans=0;
20     for(int i=31;i>=0;i--){
21         int id=(x>>i)&1;
22         if(ch[u][!id]){//如果存在与x这一位数字相反的边

```

```

23         u=ch[u][!id];
24         ans=(ans<<1)|1;//把这一位上的 1 给添加到ans上
25     }
26     else{//如果不存在
27         u=ch[u][id];
28         ans=ans<<1;//把0 给添加到ans上
29     }
30 }
31 return ans;
32 }
33 int main(){
34     scanf("%d",&n);
35     for(int i=1;i<=n;i++){
36         scanf("%d",&a[i]);
37         insert(a[i]);
38     }
39     int maxx=-1;
40     for(int i=1;i<=n;i++){
41         maxx=max(maxx,find(a[i]));
42     }
43     printf("%d\n",maxx);
44     return 0;
45 }

```

## KMP算法

### KMP求解什么问题？

字符串匹配。给你两个字符串，寻找其中一个字符串是否包含另一个字符串，如果包含，返回包含的起始位置。

```

1 char *str = "bacbababadababacambabacaddababacasdsd";
2 char *ptr = "ababaca";

```

### 算法说明

朴素的字符串匹配过程：

从目标字符串str（假设长度为n）的第一个下标选取和ptr长度（长度为m）一样的子字符串进行比较，如果一样，就返回开始处的下标值，不一样，选取str下一个下标，同样选取长度为n的字符串进行比较，直到str的末尾。这样的时间复杂度是 $O(n * m)$ 。

复杂度太高!!!

KMP：时间复杂度 $O(m + n)$

充分利用了目标字符串ptr的性质（比如里面部分字符串的重复性，即使不存在重复字段，在比较时，实现最大的移动量）。

next转移函数（pre前缀函数）

len表示最长前缀长度，i,以字符串ababcbabaa为例

1)ptr[len]==ptr[i] len++;next[i]=len;i++;

2)ptr[len]!=ptr[i] len=next[len-1];

```

1 void getNext(){

```

```

2     nextt[0]=0;
3     int len=0;
4     int i=1;
5     while(i<m){
6         if(s2[i]==s2[len]){
7             len++;
8             nextt[i]=len;
9             i++;
10        }
11        else{
12            if(len>0){
13                len=nextt[len-1];
14            }
15            else{//避免死循环
16                nextt[i]=len;
17                i++;
18            }
19        }
20    }
21 }

```

## KMP搜索

str n i

ptr m j

str[i]==ptr[j]    i++;j++;

str[i]!=ptr[j]    j=next[j];

```

1  int kmp(){
2      int i=0,j=0;
3      while(i<n){
4          if(s1[i]==s2[j]){
5              i++;j++;
6          }
7          else{
8              j=nextt[j];
9              if(j==-1){//j到了-1
10                 i++;
11                 j++;
12             }
13         }
14         if(j==m){
15             printf("找到了\n");
16             break;
17         }
18     }
19     return ans;
20 }

```

## 1465: 【例题1】剪花布条

```

1  #include<bits/stdc++.h>
2  using namespace std;
3  string s1,s2;
4  int nextt[1004],n,m;

```

```

5 void getNext(){
6     nextt[0]=0;
7     int len=0;
8     int i=1;
9     while(i<m){
10         if(s2[i]==s2[len]){
11             len++;
12             nextt[i]=len;
13             i++;
14         }
15         else{
16             if(len>0){
17                 len=nextt[len-1];
18             }
19             else{//避免死循环
20                 nextt[i]=len;
21                 i++;
22             }
23         }
24     }
25 }
26 void move_next(){
27     for(int i=m-1;i>0;i--){
28         nextt[i]=nextt[i-1];
29     }
30     nextt[0]=-1;
31 }
32 int kmp(){
33     int ans=0,j=0;
34     int i=0;
35     while(i<n){
36         if(s1[i]==s2[j]){
37             i++;j++;
38         }
39         else{
40             j=nextt[j];
41             if(j==-1){//j到了-1
42                 i++;
43                 j++;
44             }
45         }
46         if(j==m){
47             ans++;
48             j=0;
49         }
50     }
51     return ans;
52 }
53 int main(){
54     while(1){
55         cin>>s1;
56         n=s1.size();
57         if(n==1&&s1[0]=='#')
58             break;
59         cin>>s2;
60         m=s2.size();
61         getNext();
62         move_next();

```



```

63     printf("%d\n", kmp());
64 }
65 }

```

### 1466: 【例题2】Power Strings

m字符串长度，如果存在循环节，则m-next[m-1]表示每个循环节的长度，则

若 $m \%(m - \text{next}[m-1]) = 0$ ,  $\text{ans} = m / (m - \text{next}[m-1])$

若 $m \%(m - \text{next}[m-1]) \neq 0$ ,  $\text{ans} = 1$

```

1  #include <bits/stdc++.h>
2  using namespace std;
3  int m, nextt[1000005];
4  char s[1000005];
5  void getNext(){
6      nextt[0]=0;
7      int len=0;
8      int i=1;
9      while(i<m){
10         if(s[i]==s[len]){
11             len++;
12             nextt[i]=len;
13             i++;
14         }
15         else{
16             if(len>0){
17                 len=nextt[len-1];
18             }
19             else{//避免死循环
20                 nextt[i]=len;
21                 i++;
22             }
23         }
24     }
25 }
26 int main(){
27     while(scanf("%s", s)){
28         if (s[0]=='.')
29             break;
30         m=strlen(s);
31         getNext();
32         printf("%d\n", m % (m - nextt[m-1]) == 0 ? m / (m - nextt[m-1]) : 1 );
33     }
34 }

```

## AC 自动机

AC 自动机是以 **TRIE** 的结构为基础，结合 **KMP** 的思想 建立的。

### 失配指针

AC 自动机利用一个 fail 指针来辅助多模式串的匹配。

状态  $u$  的 fail 指针指向另一个状态  $v$ ，其中  $v \in Q$ ，且  $v$  是  $u$  的最长后缀（即在若干个后缀状态中取最长的一个作为 fail 指针）。

### fail 指针与 KMP 中的 next 指针的异同

- 共同点：两者同样是在失配的时候用于跳转的指针。
- 不同点：next 指针求的是最长 Border（即最长的相同前后缀），而 fail 指针指向所有模式串的前缀中匹配当前状态的最长后缀。

因为 KMP 只对一个模式串做匹配，而 AC 自动机要对多个模式串做匹配。有可能 fail 指针指向的结点对应着另一个模式串，两者前缀不同。

### 构建 fail 指针

考虑字典树中当前的结点  $u$ ， $u$  的父结点是  $p$ ， $p$  通过字符  $c$  的边指向  $u$ ，即  $trie[p, c] = u$ 。假设深度小于  $u$  的所有结点的 fail 指针都已求得。

- 如果  $trie[fail[p], c]$  存在：则让  $u$  的 fail 指针指向  $trie[fail[p], c]$ 。相当于在  $p$  和  $fail[p]$  后面加一个字符  $c$ ，分别对应  $u$  和  $fail[u]$ 。
- 如果  $trie[fail[p], c]$  不存在：那么我们继续找到  $trie[fail[fail[p]], c]$ 。重复 1 的判断过程，一直跳 fail 指针直到根结点。
- 如果真的没有，就让 fail 指针指向根结点。

如此即完成了  $fail[u]$  的构建。

以字符串 i he his she hers 举例

### build() 函数

构建 fail 指针，构建自动机

- $tr[u, c]$  可以理解为从状态（结点） $u$  后加一个字符  $c$  到达的状态（结点），即一个状态转移函数  $trans(u, c)$ 。
- $q$  队列，用于 BFS 遍历字典树。
- $fail[u]$  结点  $u$  的 fail 指针。

```
1 void build() {
2     for (int i = 0; i < 26; i++)
3         if (tr[0][i]) //将根结点的子结点一一入队
4             q.push(tr[0][i]);
5     while (q.size()) {
6         int u = q.front(); //取出队首的结点 u
7         q.pop();
8         for (int i = 0; i < 26; i++) { //遍历字符集, 即寻找u的子节点
9             if (tr[u][i]) {
10                 fail[tr[u][i]] = tr[fail[u]][i];
11                 q.push(tr[u][i]);
12             }
13             else
14                 tr[u][i] = tr[fail[u]][i]; //为什么有这个else, 以his的s为例
15         }
16     }
17 }
```

将结点按 BFS 顺序入队，依次求 fail 指针。这里的字典树根结点为 0，我们将根结点的子结点一一入队。

然后开始 BFS：每次取出队首的结点  $u$ 。fail[u] 指针已经求得，我们要求  $u$  的子结点们的 fail 指针。然后遍历字符集（这里是 0-25，对应 a-z）：

- 如果  $tr(u, i)$  存在，我们就将  $tr(u, i)$  的 fail 指针赋值为  $tr(fail[u], i)$ 。这里似乎有一个问题。根据之前的讲解，我们应该用 while 循环，不停的跳 fail 指针，判断是否存在字符  $i$  对应的结点，然后赋值。不过在代码中我们一句话就做完这件事了。
- 否则， $tr(u, i)$  不存在，就让  $tr(u, i)$  指向  $tr(fail[u], i)$  的状态。

### query()匹配函数

```

1  int query(char *t) {
2      int u = 0, res = 0;
3      for (int i = 1; t[i]; i++) {
4          u = tr[u][t[i] - 'a']; // 转移
5          for (int j = u; j && e[j] != -1; j = fail[j]) {
6              res += e[j], e[j] = -1;
7          }
8      }
9      return res;
10 }

```

### 洛谷P3808

<https://www.luogu.org/problem/P3808>

```

1  #include <bits/stdc++.h>
2  using namespace std;
3  const int N = 1e6+6;
4  int tr[N][26], tot;
5  int e[N], fail[N];
6  void insert(char *s) {
7      int u = 0;
8      for (int i = 1; s[i]; i++) {
9          if (!tr[u][s[i] - 'a']) tr[u][s[i] - 'a'] = ++tot;
10         u = tr[u][s[i] - 'a'];
11     }
12     e[u]++;
13 }
14 queue<int> q;
15 void build() {
16     for (int i = 0; i < 26; i++)
17         if (tr[0][i]) q.push(tr[0][i]);
18     while (q.size()) {
19         int u = q.front();
20         q.pop();
21         for (int i = 0; i < 26; i++) {
22             if (tr[u][i])
23                 fail[tr[u][i]] = tr[fail[u]][i], q.push(tr[u][i]);
24             else
25                 tr[u][i] = tr[fail[u]][i];
26         }
27     }
28 }
29 int query(char *t) {
30     int u = 0, res = 0;
31     for (int i = 1; t[i]; i++) {
32         u = tr[u][t[i] - 'a']; // 转移
33         for (int j = u; j && e[j] != -1; j = fail[j]) {
34             res += e[j], e[j] = -1;
35         }
36     }
37 }

```

```

36     }
37     return res;
38 }
39 int main() {
40     int n;
41     scanf("%d", &n);
42     char s1[N];
43     for (int i = 1; i <= n; i++) {
44         scanf("%s", s1 + 1);
45         insert(s1);
46     }
47     char s2[N];
48     scanf("%s", s2 + 1);
49     build();
50     printf("%d\n", query(s2));
51     return 0;
52 }

```

### 洛谷3796

<https://www.luogu.org/problem/P3796>

```

1  #include <bits/stdc++.h>
2  using namespace std;
3  const int N = 156, L = 1e6 + 6;
4  const int SZ = N * 80;
5  int tot, tr[SZ][26];
6  int fail[SZ], idx[SZ], val[SZ];
7  int cnt[N]; // 记录第 i 个字符串的出现次数
8  void init() {
9      memset(fail, 0, sizeof(fail));
10     memset(tr, 0, sizeof(tr));
11     memset(val, 0, sizeof(val));
12     memset(cnt, 0, sizeof(cnt));
13     memset(idx, 0, sizeof(idx));
14     tot = 0;
15 }
16 void insert(char *s, int id) { // id 表示原始字符串的编号
17     int u = 0;
18     for (int i = 1; s[i]; i++) {
19         if (!tr[u][s[i] - 'a']) tr[u][s[i] - 'a'] = ++tot;
20         u = tr[u][s[i] - 'a'];
21     }
22     idx[u] = id;
23 }
24 queue<int> q;
25 void build() {
26     for (int i = 0; i < 26; i++)
27         if (tr[0][i]) q.push(tr[0][i]);
28     while (q.size()) {
29         int u = q.front();
30         q.pop();
31         for (int i = 0; i < 26; i++) {
32             if (tr[u][i])
33                 fail[tr[u][i]] = tr[fail[u]][i], q.push(tr[u][i]);
34             else
35                 tr[u][i] = tr[fail[u]][i];
36         }
37     }
38 }

```

```

36     }
37 }
38 }
39 int query(char *t) { // 返回最大的出现次数
40     int u = 0, res = 0;
41     for (int i = 1; t[i]; i++) {
42         u = tr[u][t[i] - 'a'];
43         for (int j = u; j; j = fail[j]) val[j]++;
44     }
45     for (int i = 0; i <= tot; i++)
46         if (idx[i]) res = max(res, val[i]), cnt[idx[i]] = val[i];
47     return res;
48 }
49 int n;
50 char s[N][100], t[L];
51 int main() {
52     while (~scanf("%d", &n)) {
53         if (n == 0) break;
54         init();
55         for (int i = 1; i <= n; i++) scanf("%s", s[i] + 1), insert(s[i], i);
56         build();
57         scanf("%s", t + 1);
58         int x = query(t);
59         printf("%d\n", x);
60         for (int i = 1; i <= n; i++)
61             if (cnt[i] == x) printf("%s\n", s[i] + 1);
62     }
63     return 0;
64 }

```