

# 网络实验报告

高效路由查找

2016K8009909005 雷正宇

## 实验内容

实现网络路由的前缀树查找

我的具体实现为：多比特前缀树 + 多层优化

## 实验流程

### 测试集的编写

本次试验我使用了 Java。需求是实现一个可以进行路由查找的前缀树结构，它通过配置一个由子网段、前缀长度和端口构成的表来完成整棵树的构建，然后可以通过对一个 ip 地址进行匹配来得到对应的端口号或者返回打印"没有找到"的信息，如下：

```
1 interface Trie {
2     public void put(String ip, int mask, int port);
3     public int get(String ip);
4 }
```

具体实现中没有声明这个接口，直接将 Trie 作为类名。Test 类提供用于配置 Trie 和测试效率的静态方法，这个方法不能用来检测正确性（正确性检测在实验结果与分析部分会提及）：

```
1 public class Test {
2     private static final int testNumber = 398765;
3
4     public static void main(String[] args) {
5         String testFile = "./test-list.txt";
6         String config = "./forwarding-table.txt";
7         Trie trie = new Trie();
8         try {
9             makeTestFile(config, testFile);
10            makeTrieConfig(trie, config);
11            long time = lookUpTest(trie, testFile);
12
13            System.out.println("average time for one lookup: " + time *
14            1000000 / testNumber + " ns");
15        } catch (IOException e) {
16            System.out.println("wtf? It is impossible");
17        }
18    }
19 }
```

```

17     }
18     ...
19 }

```

其中, makeTestFile 用来创建测试信息:

```

1  private static void makeTestFile(String config, String filename) throws
    IOException {
2      PrintWriter pw = new PrintWriter(new FileWriter(filename));
3      BufferedReader br = new BufferedReader(new FileReader(config));
4
5      for (int i = 0; i < testNumber; i++) {
6          String[] strings = br.readLine().split(" ");
7          int ip = Trie.parseIpString(strings[0]);
8          int mask = Integer.parseInt(strings[1]);
9          pw.println(Trie.parseIpNumber(makeRandomIp(ip, mask)));
10     }
11
12     br.close();
13     pw.close();
14 }

```

Test 定义了一个常量 testNumber. 实验中给出的 forwarding-table 的总量为 697882 行, 实验中可以取一个子集用来进行测试, 以获得各方面效率和测试集大小的关系。代码中的 makeRandomIp 是一个根据网段和前缀长度随机生成网段中合法 IP 的函数:

```

1  public class Test {
2      ...
3      private static int makeRandomIp(int ip, int mask) {
4          Random random = new Random();
5          int filled;
6          if (mask == 0) {
7              filled = -1;
8          } else if (mask == 32) {
9              filled = 0;
10         } else {
11             filled = (1 << (31 - mask)) - 1;
12         }
13         return ip | (random.nextInt() & filled);
14     }
15     ...
16 }

```

makeTrieConfig 则用来根据具体文件配置 Trie:

```

1  public class Test {
2      ...

```

```

3     private static void makeTrieConfig(Trie trie, String filename) throws
IOException {
4         BufferedReader br = new BufferedReader(new FileReader(filename));
5         for (int i = 0; i < testNumber; i++) {
6             String[] strings = br.readLine().split(" ");
7             int ip = Trie.parseIpString(strings[0]);
8             int mask = Integer.parseInt(strings[1]);
9             int port = Integer.parseInt(strings[2]);
10            trie.put(ip, mask, port);
11        }
12        br.close();
13    }
14    ...
15 }

```

lookUpTest 测试指定文件中所有的 IP，并把获得的端口号输出到控制台：

```

1     private static long lookUpTest(Trie trie, String testFile) throws
IOException {
2         BufferedReader br = new BufferedReader(new FileReader(testFile));
3         int[] ip = new int[testNumber];
4         int[] port = new int[testNumber];
5         for (int i = 0; i < testNumber; i++) {
6             ip[i] = Trie.parseIpString(br.readLine());
7         }
8
9         long start = System.currentTimeMillis();
10        for (int i = 0; i < testNumber; i++) {
11            port[i] = trie.get(ip[i]);
12        }
13        long end = System.currentTimeMillis();
14
15        for (int i = 0; i < testNumber; i++) {
16            System.out.println(port[i]);
17        }
18        return end - start;
19    }

```

通过控制台打印的内容，对比配置表的端口号，可以初步检查正确性。

## 一个初步的实现

第一个实现是简单的 1bit 前缀树，只考虑结果的正确性，不考虑效率。查找方式用递归算法。以下代码中删去了 setter/getter 和几个构造方法的具体实现：

```

1     class BinTrieNode {
2         private Subnet body;
3         private boolean valid = false;

```

```

4     private int port;
5     private BinTrieNode subNode0 = null;
6     private BinTrieNode subNode1 = null;
7
8     private void addPort(int port) {
9         valid = true;
10        this.port = port;
11    }
12
13    void put(Subnet subnet, int port) {
14        if (body.getMask() == subnet.getMask()) {
15            if (valid) {
16                System.out.println("bad things happened");
17            } else {
18                addPort(port);
19            }
20        } else {
21            int index = 31 - body.getMask();
22            String newIp = Subnet.asString(Subnet.asNumber(body.getIp()) +
(1 << index));
23            if (subnet.getByIndex(index)) {
24                if (subNode1 == null) {
25                    subNode1 = new BinTrieNode(newIp, body.getMask() + 1);
26                }
27                subNode1.put(subnet, port);
28            } else {
29                if (subNode0 == null) {
30                    subNode0 = new BinTrieNode(body.getIp(),
body.getMask() + 1);
31                }
32                subNode0.put(subnet, port);
33            }
34        }
35    }
36 }
37
38 public class BinTrie {
39     private BinTrieNode root;
40
41     public BinTrie() {
42         this.root = new BinTrieNode("0.0.0.0", 0);
43     }
44
45     public void put(String ip, int mask, int port) {
46         root.put(new Subnet(ip, mask), port);
47     }
48
49     public int get(String ip) {
50         int port = 0;

```

```

51     int ipNumber = Subnet.asNumber(ip);
52     boolean found = false;
53     BinTrieNode node = root;
54     while (node != null) {
55         int mask = node.getBody().getMask();
56         if (node.isValid() && node.getBody().matchIp(ip)) {
57             found = true;
58             port = node.getPort();
59         }
60         if ((ipNumber & (1 << (31 - mask))) != 0) {
61             node = node.getSubNode1();
62         } else {
63             node = node.getSubNode0();
64         }
65     }
66     if (!found) {
67         System.out.println(ip + " not found");
68     }
69     return port;
70 }
71 }

```

其中 Subnet 是一个字段为网段和前缀长度的类，提供了用于 IP 字符串和整数相互转化的静态函数。以下代码省略了较多无关紧要的内容，例如各种 setter/getter：

```

1  public class Subnet {
2      private String ip;
3      private int mask;
4
5      public static int asNumber(String net) {
6          Integer[] parts = Arrays.stream(net.split("\\.")).
7              map((s) -> Integer.parseUnsignedInt(s)).
8              toArray(Integer[]::new);
9          return parts[3] + (parts[2] << 8) + (parts[1] << 16) + (parts[0]
10 << 24);
11     }
12
13     public static String asString(int ip) {
14         String[] nets = new String[4];
15         for (int i = 0; i < 4; i++) {
16             nets[3 - i] = Integer.toString((ip >> (i * 8)) & 0xff);
17         }
18         return String.join(".", nets);
19     }
20
21     public boolean getByIndex(int index) {
22         return (asNumber() & (1 << index)) != 0;
23     }
24 }

```

```

24     public boolean matchIp(String target) {
25         int ipa = asNumber();
26         int ipb = asNumber(target);
27         return ipa >> (32 - mask) == ipb >> (32 - mask);
28     }
29 }

```

IP 地址是32位的，Java 中没有 uint32 这样的类型，所以数字类型都是有符号的。int 类型为32位，可以用来存储完整的 IP 地址，但一定要注意用于运算时符号带来的影响。

## 优化策略

### 用 int 而不是 String 来保存数据

在初步实现中，只考虑了功能的正确性，忽视了空间消耗。所以，对于和 IP 以及网段相关的内容，在运用于 Trie 对象之前可以先转换为32位整数类型。

### 用多 bit 而不是 1bit 的前缀树

多 bit 前缀树会稍微牺牲一点空间，但可以在查找上换取成倍的速率。而且，在 JVM 和机器层面而言，当数据访问以 8bit 即一字节的整数倍为单位时，可以得到更好的 Cache 命中率并且在数据索引时省略一些编译器自动补充的移位操作。

### 叶推消除冗余节点

如果一个节点的所有子节点都已经有了表项了（在我的实现中体现为 valid），就可以把这个节点删去，将所有子节点进行层次提升。如果一个节点的孩子节点存在但是无效，则把这个节点的表项下推给孩子节点。

在下面的实现中，没有删除子节点都为 valid 的节点。因为最终我发现 8bit trie 的性能让人满意，而 8bit trie 的一个节点有 256 个子节点，全部提升意味着相当繁重的内存存取。

### 用栈做循环而不是递归查找

递归的好处在于代码清晰简单，但效率的确不如手动用栈实现。不过本例看似是一个递归算法，其实用不着栈，因为是尾递归，优化起来很方便。

## 优化后的版本

根据上述策略，Trie 几乎被重写了。重写后的版本代码量也少了许多（依旧删去了一些无关紧要的部分）：

```

1  public class Trie {
2      public static final int BITS = 8;
3      public static final int FILLED;
4      public static final int POW;
5
6      static {
7          POW = 1 << BITS;
8          FILLED = POW - 1;

```

```

9         if (32 % BITS != 0) {
10             System.out.println("Bits is illegal");
11         }
12     }
13
14     public static int parseIpString(String ipString) {
15         Integer[] parts = Arrays.stream(ipString.split("\\.")).
16             map((s) -> Integer.parseUnsignedInt(s)).
17             toArray(Integer[]::new);
18         return parts[3] | (parts[2] << 8) | (parts[1] << 16) | (parts[0]
19 << 24);
20     }
21
22     public static String parseIpNumber(int ip) {
23         String[] nets = new String[4];
24         for (int i = 0; i < 4; i++) {
25             nets[3 - i] = Integer.toString((ip >> (i * 8)) & 0xff);
26         }
27         return String.join(".", nets);
28     }
29
30     private TrieNode root = new TrieNode(0, 0);
31
32     public static int getBinsByIndex(int i, int index) {
33         return (i >> index) & FILLED;
34     }
35
36     public void put(String ip, int mask, int port) {
37         put(parseIpString(ip), mask, port);
38     }
39
40     public void put(int ip, int mask, int port) {
41         root.put(ip, mask, port);
42     }
43
44     public int get(String ip) {
45         return get(parseIpString(ip));
46     }
47
48     public int get(int ip) {
49         TrieNode find = root.get(ip);
50         if (find == null) {
51             System.out.println(parseIpNumber(ip) + " not found");
52             return 0;
53         } else {
54             return find.getPort();
55         }
56     }

```

```

57
58 class TrieNode {
59     private int ip;
60     private int mask;
61     private boolean valid = false;
62     private int validMask = 0;
63
64     private int port = 0;
65     private TrieNode[] sub = new TrieNode[Trie.POW];
66
67     private boolean match(int ip) {
68         if (this.mask == 0) {
69             return true;
70         } else if (this.valid) {
71             return this.ip >> (32 - this.validMask) == ip >> (32 -
this.validMask);
72         } else {
73             return this.ip >> (32 - this.mask) == ip >> (32 - this.mask);
74         }
75     }
76
77     private TrieNode getSub(int i, boolean autoCreate) {
78         if (i > Trie.POW) {
79             System.out.println("illegal access to sub");
80             return null;
81         }
82         if (i >= Trie.POW || i < 0) {
83             System.out.println("getSub a wrong number: " + i);
84             return null;
85         }
86         if (autoCreate && sub[i] == null) {
87             int newMask = this.mask + Trie.BITS;
88             sub[i] = new TrieNode(this.ip + (i << (32 - newMask)),
newMask);
89         }
90         return sub[i];
91     }
92
93     private TrieNode getSub(int i) {
94         return getSub(i, true);
95     }
96
97     private void addPort(int validMask, int port) {
98         if (validMask <= 0 || validMask > 32) {
99             System.out.println("add port illegally");
100         }
101         valid = true;
102         this.port = port;
103         this.validMask = validMask;

```



```

104     }
105
106     void put(int ip, int mask, int port) {
107         int from = Trie.getBinsByIndex(ip, 32 - this.mask - Trie.BITS);
108         int to = from + (1 << (Trie.BITS - mask + this.mask));
109         if (mask > this.mask) {
110             if (mask >= this.mask + Trie.BITS) {
111                 getSub(from).put(ip, mask, port);
112             } else {
113                 for (int i = from; i < to; i++) {
114                     getSub(i).put(ip, mask, port);
115                 }
116             }
117         } else if (!valid || mask > this.validMask){
118             addPort(mask, port);
119         }
120     }
121
122     TrieNode get(int ip) {
123         if (!match(ip)) {
124             return null;
125         }
126
127         int index = Trie.getBinsByIndex(ip, 32 - this.mask - Trie.BITS);
128         if (sub[index] != null && sub[index].get(ip) != null) {
129             return sub[index].get(ip);
130         } else if (this.valid) {
131             return this;
132         } else {
133             return null;
134         }
135     }
136 }

```

所有的方法都是有弹性的：当 BITS 常量发生改变时，Trie 的结构会发生变化，但行为正确性不变。后面的测试中我将会展示不同 BITS 值下的效率差别。

以上代码没有包括完全消除尾递归的优化，我曾经写了如下去递归版本：

```

1     public int getFaster(int ip) {
2         TrieNode tn = root;
3         while (tn != null) {
4             if (!tn.match(ip)) {
5                 tn = null;
6             }
7
8             int index = getBinsByIndex(ip, 32 - tn.getMask() - BITS);
9             if (tn.getSub(index, false) != null && tn.getSub(index, true)
!= null) {

```

```

10         tn = tn.getSub(index, true).get(ip);
11     } else if (tn.isValid()) {
12         return tn.getPort();
13     } else {
14         tn = null;
15     }
16 }
17 System.out.println(parseIpNumber(ip) + "not found");
18 return 0;
19 }
20
21 public int getFaster(String ip) {
22     return getFaster(parseIpString(ip));
23 }

```

然后发现这个版本几乎没有带来任何效率上的提升，根本不"faster". 所以在最终提交的版本中将它删去了。我的猜想是，编译器已经在尾递归上做了足够的优化（现在似乎很少有不支持尾递归优化的编译器了）。

## 实验结果与分析

这里直接测试优化后的版本，首先是正确性分析：

```

1 Trie trie = new Trie();
2 trie.put("192.168.255.0", 24, 61666);
3 trie.put("192.168.0.0", 16, 12345);
4 trie.put("192.0.0.0", 8, 54321);
5 trie.testIp("192.168.100.0");

```

这个样例用来测试最长匹配规则，结果如下（我似乎难以证明这是用上面的代码跑出来的结果，如果老师有兴趣可以用我的代码进行测试）：

```

/Library/Java/JavaVirtualMachines/jdk-10.0.2.jdk/Con
12345

```

```

Process finished with exit code 0

```

对于更多的数据测试，可以查看在"实验流程#测试集的编写"部分运行结果的端口号，它的顺序和 forwarding-table 中端口号的顺序一致。

接下来分别是优化后版本的 2bit, 4bit, 8bit 的测试结果（测试集大小 testNumber 为398765）：

```
/Library/Java/JavaVirtualMachines/jdk-10.0.2.  
average time for one lookup: 31898 ns  
  
Process finished with exit code 0
```

```
/Library/Java/JavaVirtualMachines/jdk-10.0.2.jdk/Content  
average time for one lookup: 692 ns  
  
Process finished with exit code 0
```

```
/Library/Java/JavaVirtualMachines/jdk-10.0  
average time for one lookup: 250 ns  
  
Process finished with exit code 0  
|
```

可以看出，2bit 的平均查找时间非常恐怖，达到了 30ms 以上，而 4bit 和 8bit 就控制在 1ms 以内了（数据波动不大，4bit 版本有时会达到 800ns 以上，8bit 版本则稳定在 250ns 左右）。

接下来是 8bit 版本在 testNumber 分别为 765, 8765, 98765, 398765 时的平均查找时间：

```
/Library/Java/JavaVirtualMachines/jdk-10.0.2.jdk/Content  
average time for one lookup: 3921 ns  
  
Process finished with exit code 0
```

```
/Library/Java/JavaVirtualMachines/jdk-10.0.2.  
average time for one lookup: 684 ns  
  
Process finished with exit code 0
```

```
/Library/Java/JavaVirtualMachines/jdk-10.0.2.  
average time for one lookup: 415 ns  
  
Process finished with exit code 0
```

```
/Library/Java/JavaVirtualMachines/jdk-10.0.2.jc
average time for one lookup: 240 ns

Process finished with exit code 0
```

测试集数量越大，查找速度越快对于这个结果，我也分析不出具体原因。时间获取函数自身占用的时间应该不到 1ms，不会造成很大影响。这个问题先保留着。

对于内存的使用，我为 JVM 设置的内存上限是 725 M，在 2bit、4bit Trie 版本下，即使将 forwarding-table 中的所有条目全部配置进 Trie 对象也不会发生内存溢出，但是 8bit 前缀树会在配置到第 600 000 行左右的数据时发生堆溢出。具体获得堆的使用情况有点麻烦，需要先使用 jsp 命令获得进程号，然后用 jmap 命令 dump 出堆的二进制信息，这个文件非常庞大，有 5GB。最后使用 jhat 命令进行解析，用浏览器访问<http://localhost:7000/>获得具体内存使用情况。

#### References to this object:

```
trie.TreeNode@0x734c62570 (41 bytes) : ??
trie.TreeNode@0x6f7f42b30 (41 bytes) : ??
trie.TreeNode@0x6e189a180 (41 bytes) : ??
trie.TreeNode@0x747c34b98 (41 bytes) : ??
trie.TreeNode@0x732505870 (41 bytes) : ??
trie.TreeNode@0x6f436e348 (41 bytes) : ??
trie.TreeNode@0x6cd689b0 (41 bytes) : ??
trie.TreeNode@0x79d3d1460 (41 bytes) : ??
trie.TreeNode@0x737c80410 (41 bytes) : ??
trie.TreeNode@0x7a4c936d8 (41 bytes) : ??
trie.TreeNode@0x748c97b88 (41 bytes) : ??
trie.TreeNode@0x7aa393ed8 (41 bytes) : ??
trie.TreeNode@0x74888a460 (41 bytes) : ??
trie.TreeNode@0x6dab3f458 (41 bytes) : ??
trie.TreeNode@0x6f4791eb8 (41 bytes) : ??
trie.TreeNode@0x6fc431088 (41 bytes) : ??
trie.TreeNode@0x7ab35f370 (41 bytes) : ??
trie.TreeNode@0x7276b14c0 (41 bytes) : ??
trie.TreeNode@0x6e12daf78 (41 bytes) : ??
trie.TreeNode@0x7a01fca80 (41 bytes) : ??
trie.TreeNode@0x70fa256e8 (41 bytes) : ??
trie.TreeNode@0x6fa903ae8 (41 bytes) : ??
trie.TreeNode@0x6f9f67598 (41 bytes) : ??
trie.TreeNode@0x6f91e3a88 (41 bytes) : ??
trie.TreeNode@0x6e850ca58 (41 bytes) : ??
trie.TreeNode@0x6e1a6fc98 (41 bytes) : ??
trie.TreeNode@0x744252850 (41 bytes) : ??
```

最终页面呈现如上图。这里罗列的应该是所有的 TrieNode 对象，每个占 41 字节，但很难确定到底有多少个。我的电脑目前启动 jhat 后会变得相当卡顿。但是我可以通过在 TrieNode 类中增加对象计数器来统计对象个数。最终计算出的 testNumber 为 398765 时，2bit, 4bit, 8bit 的 TrieNode 所占用的总空间分别为：

```
total space: 27650 KB
```

```
total space: 35317 KB
```

```
total space: 99409 KB
```

即不到 100MB，而且测试集越大，8bit Trie 占用的额外空间比例就越小（虽然还是很耗费空间）。

