

ট্রাই বা প্রিফিক্স ট্রি ব্যবহার করে ডাটাবেস থেকে খুব সহজে স্ট্রিং খুঁজে বের করা যায়। ধরো একটা ফোনবুকে একটা শহরের সব মানুষের ফোন নম্বর রাখা আছে। শহরে মানুষ আছে হয়তো কয়েক লক্ষ্য, কিন্তু প্রতিটা মানুষের নাম সর্বোচ্চ ২০টা অক্ষর ব্যবহার করে লেখা যায়। আমরা এমন একটা ডাটা স্ট্রাকচার ব্যবহার করে নাম খুঁজবো যে নির্ভর করে শুধু মাত্র নামটিতে কয়টি অক্ষর আছে তার উপর। যেমন “Alice” নামটি খুঁজতে মাত্র ৫টি অপারেশন করা লাগবে ডাটাবেস যত বড়ই হোক না কেন।

এই লেখা পড়ার আগে **লিংকড লিস্ট** এবং রিকার্সন সম্পর্কে ধারণা থাকতে হবে।
ধরো তোমাকে একটা ডিকশনারী দেয়া হলো যেখানে নিচের শব্দগুলো আছে:

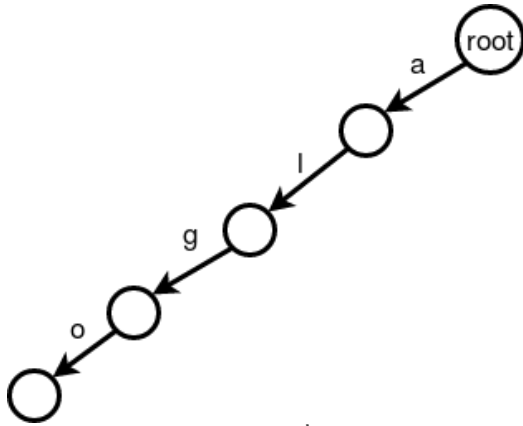
algo
algea
also
tom
to

এখন আমরা এই ডিকশনারিটাকে এমনভাবে মেমরিতে রাখতে চেষ্টা করবো যেন খুব সহজে কোনো একটা শব্দ খুঁজে পাওয়া যায়। একটা উপায় হলো শব্দগুলোকে সর্ট করে রাখা যেটা কাগজের ডিকশনারি গুলোতে রাখা হয়, তাহলে বাইনারি সার্চ করেই আমরা কোনো একটা শব্দ খুঁজে বের করতে পারবো। আরেকটা উপায় হলো প্রিফিক্স ট্রি বা সংক্ষেপে ট্রাই(trie) ব্যবহার করা। trie শব্দটা এসেছে “retrieval” শব্দটা থেকে। সেই হিসাবে এটার উচ্চারণ “ট্রি” হওয়ার কথা কিন্তু গ্রাফ থিওরীতে ট্রি এর আরো ব্যাপক ব্যবহার আছে তাই এটাকে বলা হয় “ট্রাই”। **প্রিফিক্স মানে হলো একটা স্ট্রিং এর শুরু থেকে কয়েকটা ক্যারেকটার নিয়ে নতুন স্ট্রিং তৈরি করা। যেমন blog এর প্রিফিক্স হলো b,bl,blo এবং blog।**

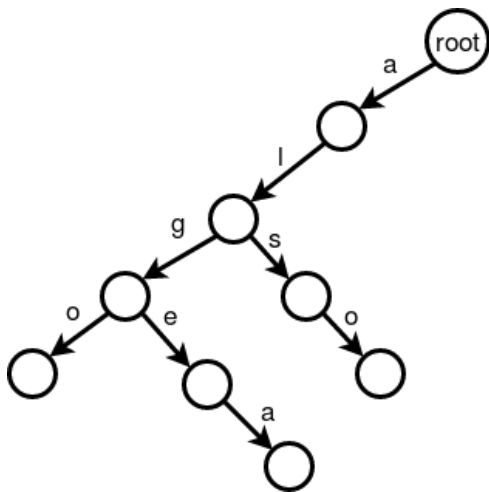
শুরুতে আমাদের একটা রুট নোড ছাড়া কিছুই নেই।



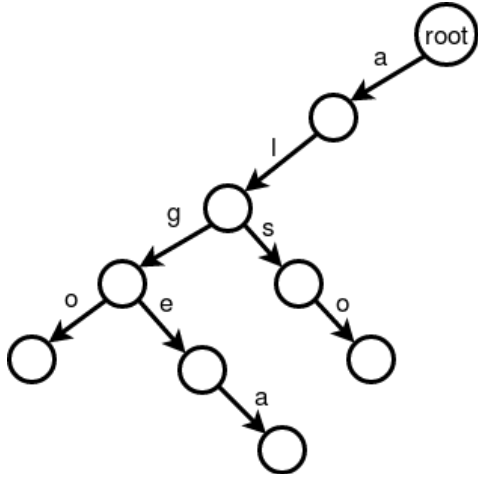
এখন আমরা algo শব্দটাকে যোগ করবো। নিচের ছবিতে দেখো কিভাবে শব্দটা যোগ করা হয়েছে।
রুট নোড থেকে আমরা একটা এজ দিবো যেই এজ এর নাম হবে “a”। তারপর নতুন তৈরি হওয়া নোড থেকে “l” নামের একটা এজ তৈরি করবো। এভাবে “g” আর “o” এজ দুইটাও তৈরি করবো।
লক্ষ্য করো নোডে আমরা কোনো তথ্য রাখছি না, খালি নোড থেকে এজ বের করছি।



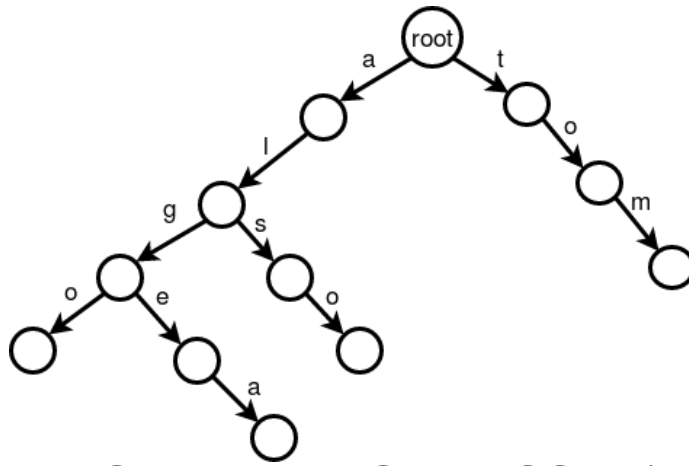
এখন আমরা alga শব্দটা যোগ করতে চাই। রুট থেকে “a” নামের এজ দরকার, সেটা অলরেডি আছে, নতুন করে যোগ করা দরকার নাই। ঠিক সেরকম a থেকে। এবং l থেকে g তেও এজ আছে। তারমানে “alg” অলরেডি ট্রাই তে আছে, আমরা শুধু e আর a যোগ করবো।



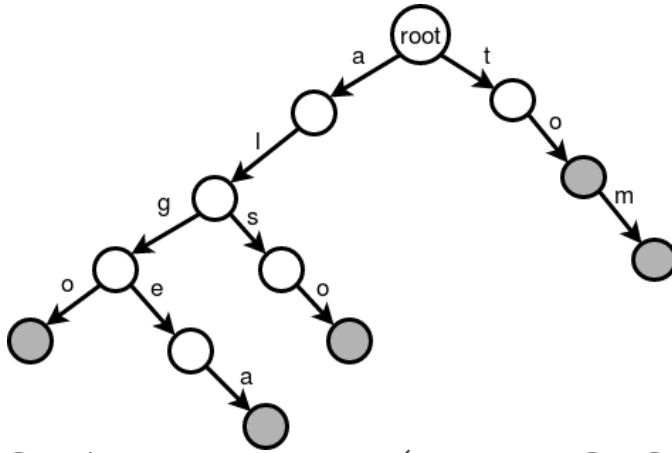
also শব্দটাকে যোগ করবো এবার। রুট থেকে “al” প্রিফিক্স এরইমধ্যে আছে, শুধু “so” যোগ করতে হবে।



এবার “tom” যোগ করি। এবার রুট থেকে নতুন এজ তৈরি করতে হবে কারণ tom এর কোনো প্রিফিক্স আগে যোগ করা হয়নি।



এখন “to” শব্দটা যোগ করবো কিভাবে? “to” পুরোপুরি tom এর প্রিফিক্স তাই নতুন কোনো এজ যোগ করা দরকার নাই। আমরা যে কাজটা করতে পারি সেটা বলে নোডগুলোতে কিছু এন্ড-মার্ক বসানো। যেসব নোডে এসে অন্তত একটা শব্দ কমপ্লিট হয়েছে সেসব নোডে আমরা এন্ডমার্ক বসিয়ে দেই, ছবিতে ধূসর রঙ দিয়ে এন্ডমার্ক বোঝানো হয়েছে। আগের সব শব্দের জন্য এবং সেই সাথে নতুন শব্দ “to” এর জন্য এন্ডমার্ক বসালে ট্রাইটা এরকম দেখাবে:



নিশ্চয়ই বুঝতে পারছেন এন্ডমার্কগুলো কেন বসিয়েছি। মার্ক দেখে সহজেই বুঝা যাচ্ছে কোন কোন শব্দ ট্রাইতে আছে। কোন ক্যারেকটার নিচ্ছি সেই তথ্য থাকবে এজ এ, আর এন্ডমার্কগুলো থাকবে নোড এ।

এভাবে শব্দগুলো রাখার সুবিধা কি? ধরো তোমাকে বলা হলো “alice” শব্দটা ডিকশনারিতে আছে কিনা বলতে। তুমি শুরু থেকে ট্রাই ধরে আগাতে থাকো। প্রথমে দেখো রুট থেকে a নামের এজ আছে নাকি, তারপর চেক করো a থেকে i নামের এজ আছে নাকি। এরপরে i থেকে e নামের এজ খুঁজে পাওয়া যাচ্ছেনা, তাই বলতে পারো alice শব্দটা নেই।

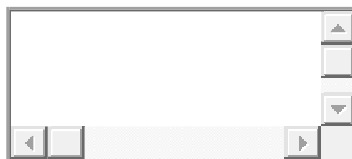
“alg” শব্দটা খুঁজতে দিলে তুমি root->a, a->l এবং l->g এজগুলো সবই খুঁজে পাবে, কিন্তু শেষ পর্যন্ত কোনো ধূসর নোডে যেতে পারবেনা, তারমানে alg ও ডিকশনারিতে নেই। “tom” খুঁজতে গেলে তুমি একটা ধূসর নোডে গিয়ে শেষ করবে তাই শব্দটা ডিকশনারিতে আছে।

ট্রাই ইমপ্লিমেন্ট করার সহজ একটা উপায় হলো লিংকড লিস্ট ব্যবহার করা। লিংকড লিস্ট, পয়েন্টার এসব শুনে ভয়ের কিছুই নেই, তুমি যদি লিংকলিস্ট ব্যবহার করতে অভ্যস্ত নাও হও আশা করি এই ইমপ্লিমেন্টেশনটা দেখে শিখে ফেলতে পারবে। আমাদের প্রতিটা নোডে ২টি জিনিস থাকবে:

১. এন্ড-মার্ক রাখার জন্য একটা ভ্যারিয়েবল।

২. প্রতিটা নোড থেকে a,b,c,...,x,y,z ইত্যাদি নামের এজ তৈরি হতে পারে। আমরা প্রতিটা ক্যারেকটারের জন্য একটা পয়েন্টার রাখবো। পয়েন্টারের সাহায্যে একটা নোড আরেকটার সাথে যোগ হবে। a নামের পয়েন্টার দিয়ে যোগ হলে বুঝতে হবে কারেন্ট নোড থেকে a নামের একটা এজ আছে। শুরুতে সবগুলো পয়েন্টার “নাল” থাকবে।

আমরা প্রথমেই নোডটা তৈরি করে ফেলি:



```
1 struct node {
2     bool endmark;
3     node* next[26 + 1];
4     node()
```

```

5  {
6      endmark = false;
7      for (int i = 0; i < 26; i++)
8          next[i] = NULL;
9  }
10 } *root;
11 int main()
12 {
13     root = new node();
14     return 0;
15 }

```

next[]next[] অ্যারের প্রতিটা এলিমেন্ট আরেকটা নোডকে পয়েন্ট করে। next[0]next[0] দিয়ে নতুন নোডকে পয়েন্ট করা হলে সেই এজ এর নাম “aa”, next[1]next[1] এর জন্য এজ এর নাম “bb”, next[25]next[25] এর জন্য “zz”। শুরুতে সবগুলো পয়েন্টার নাল। লক্ষ্য করো নোডের ভিতর একটা কনস্ট্রাক্টর “node()” তৈরি করেছি। যখনই নতুন নোড তৈরির জন্য new node()new node() কল করবো তখনই ভ্যারিয়েবলগুলোকে শূন্য বা নাল বানিয়ে দিবে। এটা না দিলে গারবেজ ভ্যালু থাকতো। rootroot ভ্যারিয়েবলটা হলো আমাদের রুট নোড, উপরের ছবিগুলোতে লাল রঙ এর নোড। আসলে রুট একটা পয়েন্টার, যখন root=new node();root=new node(); লাইনটা এক্সিকিউট হবে তখনই একটা নতুন নোড তৈরি করে rootroot যে মেমরি অ্যাড্রেসকে পয়েন্ট করে সেখানে অ্যাসাইন করে দেয়া হবে। এটাকে একটু গালভরা ভাষায় বলে instanceinstance তৈরি করা। এবার আমাদের একটা ফাংশন লাগবে নতুন শব্দ ট্রাইতে যোগ করার জন্য:



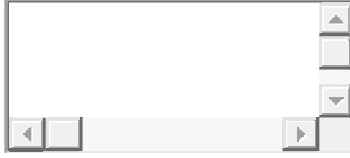
```

1 void insert(char* str, int len)
2 {
3     node* curr = root;
4     for (int i = 0; i < len; i++) {
5         int id = str[i] - 'a';
6         if (curr->next[id] == NULL)
7             curr->next[id] = new node();
8         curr = curr->next[id];
9     }
10    curr->endmark = 1;
11 }

```

রুট ভ্যারিয়েবলটা আমাদের সবসময় দরকার হবে তাই “curr” এর মধ্যে সেটার কপি তৈরি করে কাজ করি। যেহেতু পয়েন্টার নিয়ে কাজ করছি তাই রুট থেকে নতুন এজ তৈরি করা আর “curr” থেকে নতুন এজ তৈরি করা একই কথা। আমরা এখন শুধু a-z নিয়ে কাজ করছি, তাই অ্যাসকি ভ্যালুগুলোকে 0-25 এ কনভার্ট করে নিবো ‘a’ এর অ্যাসকি ভ্যালু বিয়োগ করে। insert করা খুব সহজ, আমরা শুধু চেক করবো কারেন্ট নোড (curr) থেকে বর্তমানে যে ক্যারেকটারে আছে সেই নামের কোনো এজ আছে নাকি, না থাকলে নতুন নোড তৈরি করতে হবে। এরপরে সেই এজ ধরে আমরা পরের নোডে যাবো। সবার শেষ নোডটায় এন্ড-মার্ক true করে দিবে।

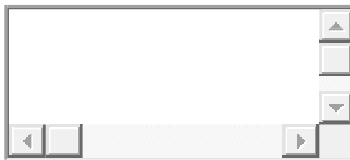
ইনসার্ট করার পর এখন সার্চ করবো। এটা আসলে ঠিক ইনসার্ট এর মতোই। পার্থক্য হলো যে এজটা দরকার সেটা না থাকলে তৈরি করে নিচ্ছিলাম, এখন এজ না থাকলে false রিটার্ন করে দিবে।



```
1 bool search(char* str, int len)
2 {
3     node* curr = root;
4     for (int i = 0; i < len; i++) {
5         int id = str[i] - 'a';
6         if (curr->next[id] == NULL)
7             return false;
8         curr = curr->next[id];
9     }
10    return curr->endmark;
11 }
```

লক্ষ্য করো সবকাজ ইনসার্ট এর মতোই করেছি। সবার শেষে লাস্ট নোডটার এন্ডমার্ক রিটার্ন করে দিয়েছি। এন্ডমার্ক true হলে শব্দটা আছে, false হলে নাই।

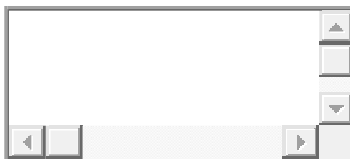
আমাদের মূল কোড শেষ। আমরা এখন যেকোনো শব্দ ট্রাইতে যোগ করতে পারবো, আবার ট্রাই থেকে কোনো শব্দ খুজতে পারবো। অনেক সময় প্রতিটা টেস্টকেস এর জন্য ট্রাই তৈরি করতে গেলে মেমরি লিমিট নিয়ে সমস্যা হয়। তাই নিরাপদ উপায় হলো প্রতি কেস এর পর ব্যবহৃত মেমরি-সেল গুলোকে ডিলিট করে দেয়া। শুধু root ডিলিট করলে হবেনা, প্রতিটা নোড করতে হবে। আমরা সে জন্য একটা রিকার্সিভ ফাংশন লিখতে পারি:



```
1 void del(node* cur)
2 {
3     for (int i = 0; i < 26; i++)
4         if (cur->next[i])
5             del(cur->next[i]);
6     delete (cur);
7 }
```

এই ফাংশনটা প্রতিটা নোডে গিয়ে আগে চাইল্ডগুলোকে ডিলিট করে এসে তারপর নোডটাকে ডিলিট করে দিবে।

সম্পূর্ণ কোডটা এরকম:



```
1 struct node {
2     bool endmark;
3     node* next[26 + 1];
4 }
```

```

4   node()
5   {
6       endmark = false;
7       for (int i = 0; i < 26; i++)
8           next[i] = NULL;
9   }
10 } * root;
11 void insert(char* str, int len)
12 {
13     node* curr = root;
14     for (int i = 0; i < len; i++) {
15         int id = str[i] - 'a';
16         if (curr->next[id] == NULL)
17             curr->next[id] = new node();
18         curr = curr->next[id];
19     }
20     curr->endmark = true;
21 }
22 bool search(char* str, int len)
23 {
24     node* curr = root;
25     for (int i = 0; i < len; i++) {
26         int id = str[i] - 'a';
27         if (curr->next[id] == NULL)
28             return false;
29         curr = curr->next[id];
30     }
31     return curr->endmark;
32 }
33 void del(node* cur)
34 {
35     for (int i = 0; i < 26; i++)
36         if (cur->next[i])
37             del(cur->next[i]);
38     delete (cur);
39 }
40 }
41 int main()
42 {
43     puts("ENTER NUMBER OF WORDS");
44     root = new node();
45     int num_word;
46     cin >> num_word;
47     for (int i = 1; i <= num_word; i++) {
48         char str[50];
49         scanf("%s", str);
50         insert(str, strlen(str));
51     }
52     puts("ENTER NUMBER OF QUERY");
53     int query;
54     cin >> query;
55     for (int i = 1; i <= query; i++) {
56         char str[50];
57         scanf("%s", str);
58         if (search(str, strlen(str)))
59             puts("FOUND");
60         else
61             puts("NOT FOUND");
62     }
63     del(root); //ট্রাইটা ধবংস করে দিলাম
64     return 0;
65 }

```

কমপ্লেক্সিটি: প্রতিটা শব্দ খুঁজতে লুপ চালাতে হচ্ছে শব্দটার লেংথ পর্যন্ত, সার্চিং এর কমপ্লেক্সিটি $O(\text{length})$ । প্রতিটা শব্দ ইনসার্ট করার কমপ্লেক্সিটিও একই। মেমরি কতখানি লাগবে সেটা ডিপেন্ড করে ইমপ্লিমেন্টেশন এবং শব্দগুলোর প্রিফিক্স কতখানি ম্যাচ করে তার উপর। উপরের ইমপ্লিমেন্টেশন দিয়ে প্রায় 10^6 টা ক্যারেকটার ট্রাইতে ইনসার্ট করা যাবে (10^6 টা ওয়ার্ড নয়, ক্যারেকটার বা লেটার)।

তুমি চাইলে ট্রাই লিংকলিস্ট ছাড়া সাধারণ অ্যারে ব্যবহার করে ইমপ্লিমেন্ট করতে পারো, নিজে চেষ্টা করো!

ট্রাই এর কিছু ব্যবহার:

১. একটা ডিকশনারিতে অনেকগুলো শব্দ আছে, কোনো একটা শব্দ আছে নাকি নাই খুঁজে বের করতে হবে। এই প্রবলেমটা আমরা উপরের কোডেই সলভ করেছি।

২. ধরো তোমার ৩ বন্ধুর টেলিফোন নম্বর হলো “৫৬৭৮”, “৪৩২২”, “৫৬৭”। তুমি যখন প্রথম বন্ধুকে ডায়াল করবে তখন ৫৬৭ চাপার সাথে সাথে ৩য় বন্ধুর কাছে ফোন চলে যাবে কারণ ৩য় বন্ধুর নাম্বার প্রথম জনের প্রিফিক্স! অনেকগুলো ফোন নম্বর দেয়া আছে, বলতে হবে এরকম কোনো নম্বর আছে নাকি যেটা অন্য নম্বরের প্রিফিক্স। (UVA 11362)।

৩. একটা ডিকশনারিতে অনেকগুলো শব্দ আছে। এখন কোনো একটা শব্দ কয়বার “prefix” হিসাবে এসেছে সেটা বের করতে হবে। যেমন “al” শব্দটা উপরের ডিকশনারিতে ৩বার প্রিফিক্স হিসাবে এসেছে (algo, algea, also এই সবগুলো শব্দের প্রিফিক্স “al”)। এটা বের করার জন্য প্রতিটা নোডে একটা কাউন্টার ভ্যারিয়েবল রাখতে হবে, কোনো নোডে যতবার যাবে ততবার কাউন্টারের মান বাড়িয়ে দিবে। সার্চ করার সময় প্রিফিক্সটা খুঁজে বের করে কাউন্টারের মান দেখবে।

৪. মোবাইলের ফোনবুকে সার্চ করার সময় তুমি যখন কয়েকটা লেটার লিখো তখন সেই প্রিফিক্স দিয়ে কি কি নাম শুরু হয়েছে সেগুলো সাজেশন বক্সে দেখায়। এটা তুমি ট্রাই দিয়ে ইমপ্লিমেন্ট করতে পারবে?

৪. দুটি স্ট্রিং এর “longest common substring” বের করতে হবে। (subsequence হলে ডিপি দিয়ে সহজে করা যায়, এখানে substring চেয়েছি)।

(হিন্টস: একটা স্ট্রিং এর শেষ থেকে এক বা একাধিক ক্যারেকটার নেয়া হলে সেটাকে স্ট্রিংটার সাফিক্স বলে, যেমন blog এর সাফিক্স g,og,log,blog। আর প্রতিটা substring ই কিন্তু কোনো না কোনো সাফিক্স এর প্রিফিক্স!! তাই সবগুলো সাফিক্সকে ট্রাইতে ইনসার্ট করলে কাজটা সহজ হয়ে যায়!)

৫. (অ্যাডভান্সড) সম্ভবত ২০১১তে ডেফোডিল ইউনিভার্সিটির ন্যাশনাল কনটেস্টে এসেছিলো প্রবলেমটা। একটা ডিকশনারি ইনপুট দেয়া থাকবে। প্রতিবার ডিকশনারির ২টা শব্দ কুয়েরি দিবে, বলতে হবে তাদের মধ্যে common prefix এর দৈর্ঘ্য কত। যেমন algo আর algea এর কমন প্রিফিক্স alg, দৈর্ঘ্য ৩। ট্রাইতে ডিকশনারিতে ইনসার্ট করে প্রতি কুয়েরিতে শব্দদুটি এন্ড-মার্ক থেকে LCA(lowest common ancestor) বের করে প্রবলেমটা সলভ করা যায়।

কিছু প্রবলেম:

<https://leetcode.com/problems/search-suggestions-system/>

UVA 10226

(UVA 11362 Phonebook)

UVA 11488 Hyper prefix sets

POJ 2001 Shortest Prefix

POJ 1056

হ্যাপি কোডিং!