

# Reiknirit Ahos og Corasicks (1975)

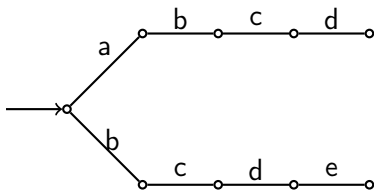
Bergur Snorrason

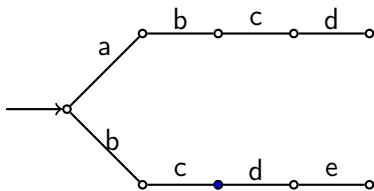
13. apríl 2023

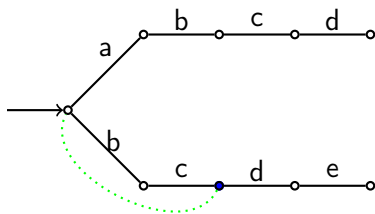
- ▶ Gerum ráð fyrir að við séum með stafróf  $\Sigma$ , streng  $s$  og lista af  $n$  strengjum  $p$ , þar sem  $j$ -ti strengurinn kallast  $p_j$ .
- ▶ Látum  $|s|$  tákna lengd strengsins  $s$  og  $|p| = |p_1| + \dots + |p_n|$ .
- ▶ Við viljum finna alla hlutstrengi  $s$  sem eru í listanum  $p$ .
- ▶ Við getum notað reiknirit Knuths, Morrisar og Pratt's  $n$  sinnum.
- ▶ Þessi aðferð hefur tímaflækjuna  $\mathcal{O}(n \cdot |s| + |p|)$ .
- ▶ Reiknirit Aho og Corasick's bætir þetta.

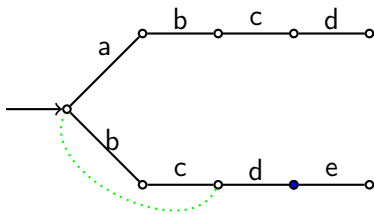
- ▶ Byrjum á að setja all strengina í  $p$  inn í forstrengstré  $T$ .
- ▶ Við viljum síðan búa til stöðuvél úr  $T$ .
- ▶ Hnúturnir í trénu eru stöðurnar, en okkur vantar að finna færslur fyrir hverja stöðu og bókstaf í  $\Sigma$ .
- ▶ Gerum ráð fyrir að við séum í hnút  $v$  í  $T$  og viljum færast fyrir staf  $c$  í  $\Sigma$ .
- ▶ Ef það er leggur í  $T$  úr  $v$  merktur með  $c$  þá ferðumst við eftir honum.
- ▶ Annars þurfum við að fara aftur í hnút  $w$  þannig að strengurinn sem svarar til  $w$  er bakstrengur strengsins sem svarar til  $v$ .
- ▶ Við viljum hafa strenginn  $w$  sem lengstan (með öðrum orðum viljum við fara sem styst aftur).
- ▶ Í hvern hnút bætum við við legg sem svarar til slíkrar færslu.
- ▶ Við köllum þessa leggi *bakstrengsleggi* (e. *suffix links*).
- ▶ Takið eftir að þeir eru í raun óháðir bókstafnum  $c$ .
- ▶ Við látum bakstrengslegg rótarinn benda á sjálfa sig.

- ▶ Hvernig finnum við alla bakstrengsleggina?
- ▶ Við notum kvika bestun.
- ▶ Látum  $f(w, c)$  tákna færslu úr stöðu  $w$  með staf  $c$  og  $g(w)$  tákna bakstrengslegg  $w$ .
- ▶ Gerum einnig ráð fyrir að foreldri  $v$  sé  $p$  og  $f(p, a) = v$ .
- ▶ Við sjáum þá að  $g(v) = f(g(p), a)$ .
- ▶ Með öðrum orðum förum við upp í foreldrið, ferðumst eftir bakstrenglegg þaðan og færum okkur í stöðuvélinni eftir  $a$ .
- ▶ Við höfum þá rakningarformúlu sem við getum notað til að reikna bakstrengshlekkina.

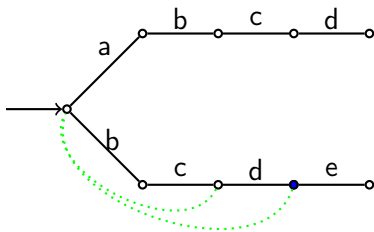


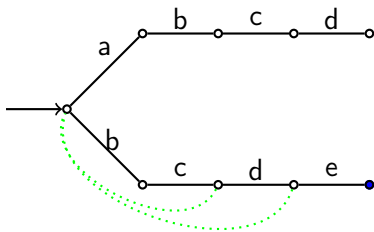


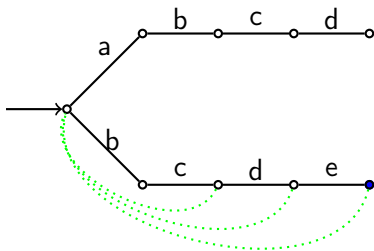


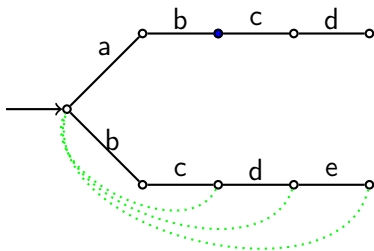


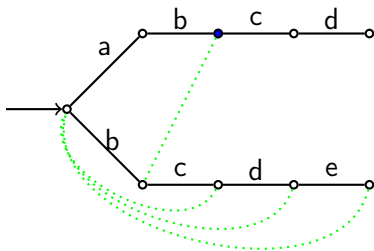


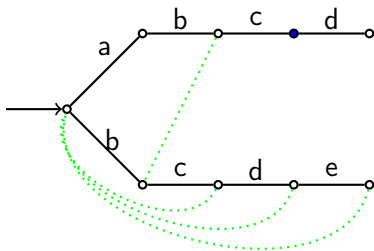


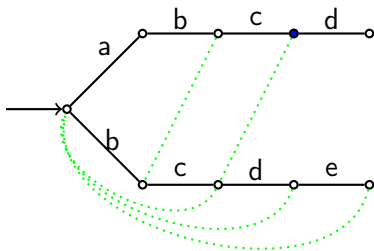


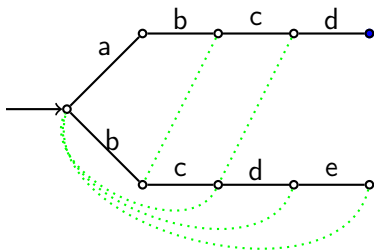




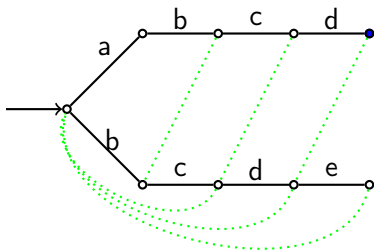






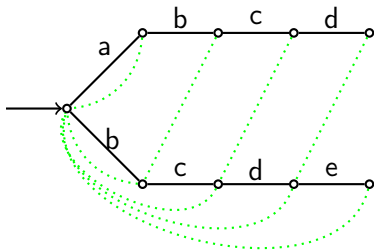




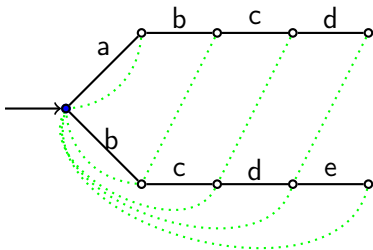


- ▶ Í  $T$  merkjum við lokastöður þar sem strengir enda.
- ▶ Við ferðumst svo í gegnum strenginn  $s$  og hliðrum stöðvélinni miðað við stafina í  $s$ .

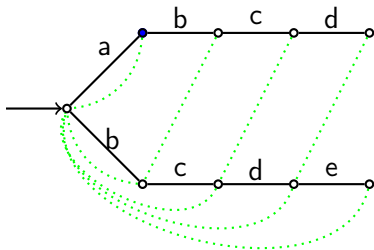
„abcdcdeaaaabcdeabcxab”



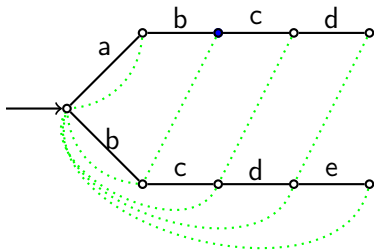
„abcdcdeaaaabcdeabcxab”



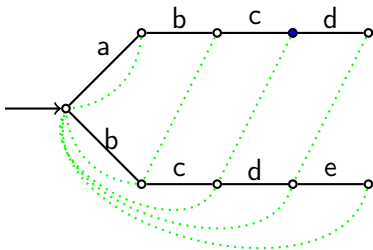
„bcdcdcaaaabcdeabcxab”



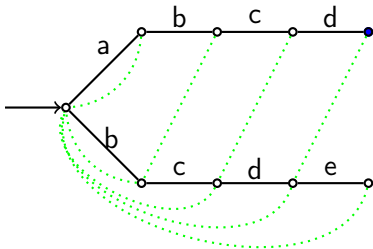
„cdcdeaaaabcdeabcxab”



„dcdeaaaabcdeabcxab”

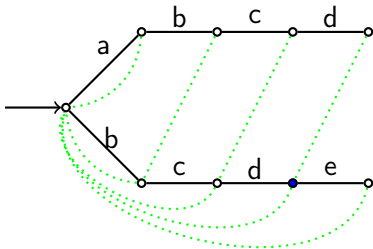


„cdeaaabcdeabcxab”

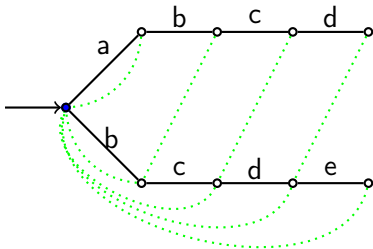




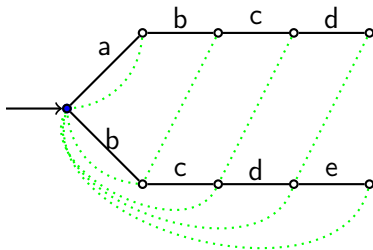
„cdeaaabcdeabcxab”



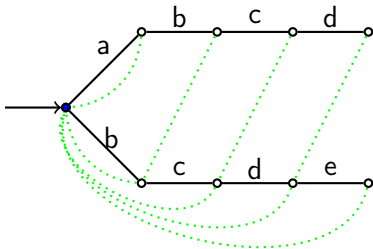
„cdeaaaabcdeabcxab”



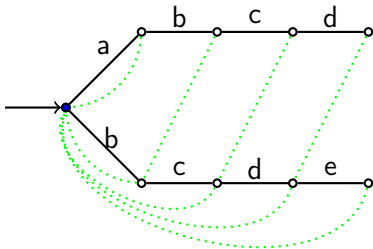
„deaaaabcdeabcxab”



„eaaabcdeabcxab”

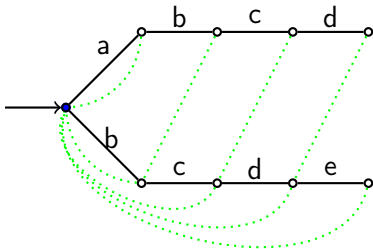


„aabcdeabcxab”

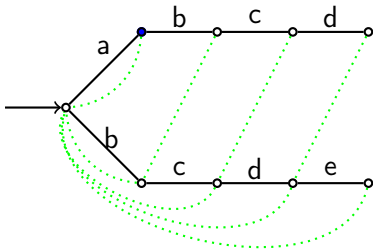




„aabcdeabcxab”

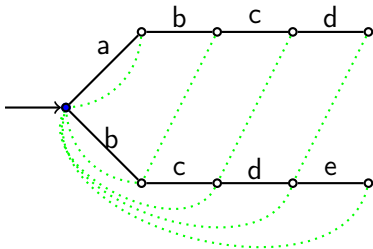


„abcdeabcxab”

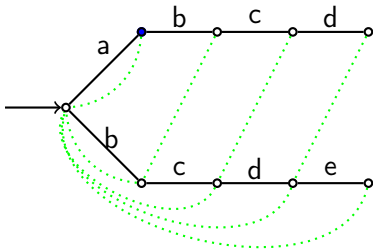




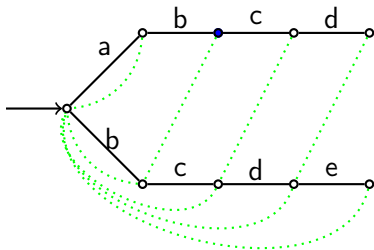
„abcdeabcxab”



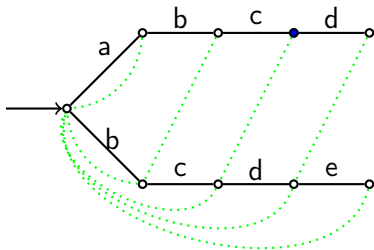
„bcdeabcxab”



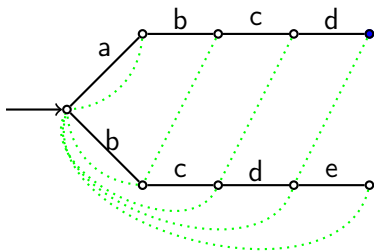
„cdeabcxab”



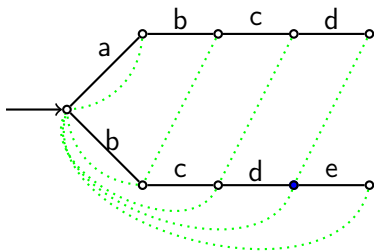
„deabcxab”



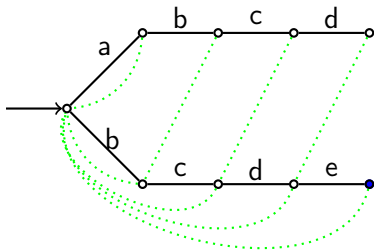
„eabcxab”



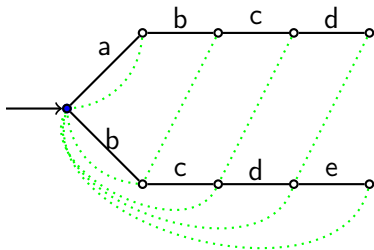
„eabcxab”



„abcxab”

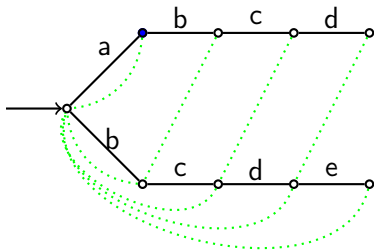


„abcxab”

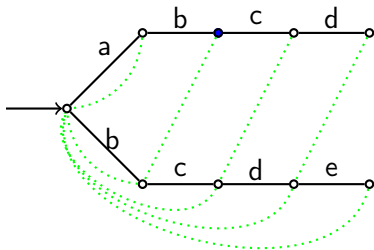




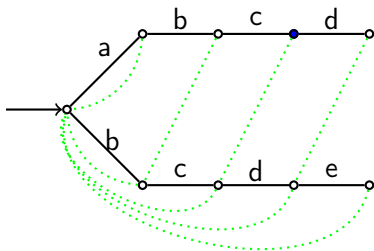
„bcxab”



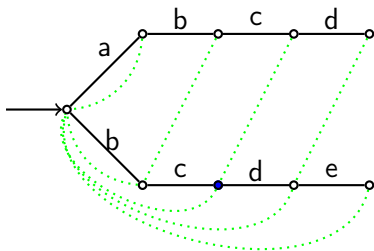
„cxab”



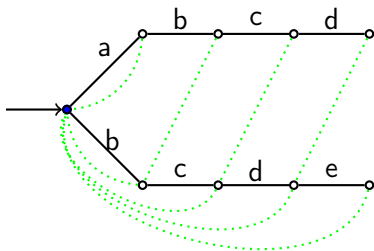
„xab“



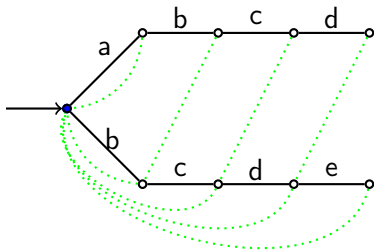
„xab“

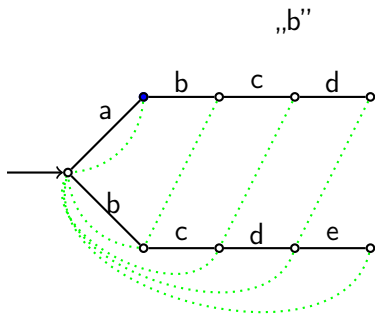


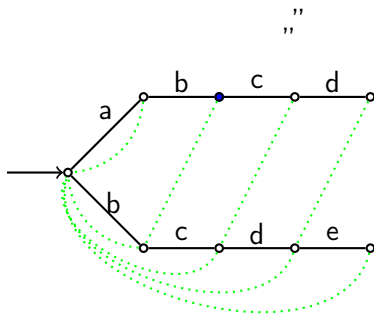
„xab“



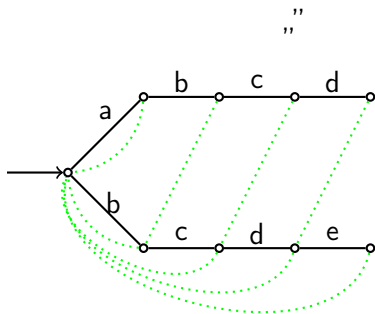
„ab”





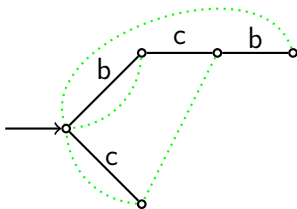




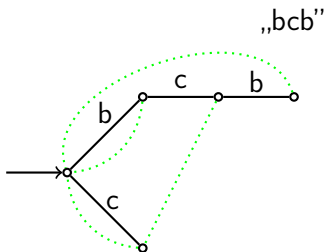


- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?

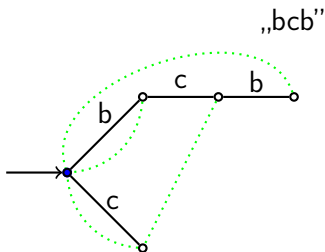
- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?



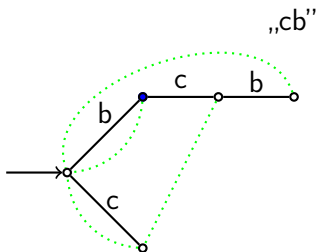
- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?



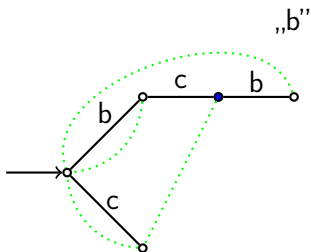
- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?



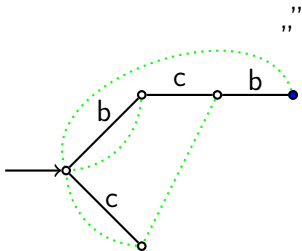
- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?



- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?

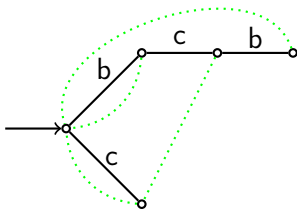


- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?

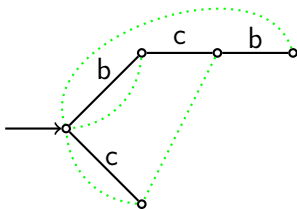




- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?



- ▶ Ljóst er að alltaf þegar við erum í lokastöðu þá erum við með hlutstreng í  $s$  sem er í  $p$ .
- ▶ En eru þetta einu slíku hlutstrengirnir?



- ▶ Nei, við þurfum líka, í hverju skrefi, að athuga hvort við getum komist í lokastöðu ef við ferðumst eftir bakstrengsleggjum.
- ▶ Til að koma í veg fyrir að tímaflækjan verði of slæm þá notum við aftur kvika bestun.
- ▶ Við bætum í raun leggjum inn í tréð, sem við köllum *lokaleggi* (e. *exit links*).

- ▶ Í útfærslunni munum við notast við þrjú hjálparföll.
- ▶ Fyrsta heitir `trie_step(...)` sem er notað til að ferðast um stöðuvélina.
- ▶ Annað heitir `trie_suffix(...)` sem er notað til að finna bakstrengsleggi.
- ▶ Þriðja heitir `trie_exit(...)` sem er notað til að finna lokaleggina.
- ▶ Öll þessi föll eru endurkvæm og notast við minnun.

```

5 #define ALPHABET 128
6 #define MAXN 1000000
7 typedef struct { int v, n; } listnode;
8 typedef struct { int t[ALPHABET], l, e, p, c, d; } trienode;
9 typedef struct { int s, r, l; trienode m[MAXN + 1]; listnode w[MAXN]; } trie;
10 int val(char c) { return c; }
11 int list_node(trie *t, int v, int n)
12 {
13     t->w[t->l].v = v, t->w[t->l].n = n;
14     return t->l++;
15 }
16 int trie_node(trie *t, int p, int c)
17 {
18     for (int i = 0; i < ALPHABET; i++) t->m[t->s].t[i] = -1;
19     t->m[t->s].l = t->m[t->s].e = t->m[t->s].d = -1;
20     t->m[t->s].p = p, t->m[t->s].c = c,
21     return t->s++;
22 }
23 void trie_init(trie *t) { t->s = t->l = 0, t->r = trie_node(t, -1, -1); }
24
25 void trie_insert(trie *t, char *s, int x)
26 {
27     int h;
28     for (h = t->r; *s; h = t->m[h].t[val(*s++)])
29         if (t->m[h].t[val(*s)] == -1)
30             t->m[h].t[val(*s)] = trie_node(t, h, val(*s));
31     t->m[h].l = list_node(t, x, t->m[h].l);
32 }

```

```

35 int trie_suffix(trie *t, int h)
36 {
37     if (t->m[h].d != -1) return t->m[h].d;
38     if (h == t->r || t->m[h].p == t->r) return t->m[h].d = t->r;
39     return t->m[h].d = trie_step(t, trie_suffix(t, t->m[h].p), t->m[h].c);
40 }
41
42 int trie_step(trie *t, int h, int c)
43 {
44     if (t->m[h].t[c] != -1) return t->m[h].t[c];
45     return t->m[h].t[c] = h == t->r ? t->r :
46         trie_step(t, trie_suffix(t, h), c);
47 }
48
49 int trie_exit(trie *t, int h)
50 {
51     if (t->m[h].e != -1) return t->m[h].e;
52     if (h == 0 || t->m[h].l != -1) return t->m[h].e = h;
53     return t->m[h].e = trie_exit(t, trie_suffix(t, h));
54 }

```

```

57 int aho_corasick(char *s, char **p, int m)
58 {
59     trie_init(&t);
60     int h, i, j, k, w, l[m];
61     for (i = 0; i < m; i++) l[i] = strlen(p[i]);
62     for (i = 0; i < m; i++) trie_insert(&t, p[i], i);
63     printf(" ");
64     for (i = 0; i < strlen(s); i++) printf("%0d", i%10); printf("\n");
65     printf("searching in %s\n", s);
66     for (i = 0, j = 0, h = t.r; ; j++)
67     {
68         k = trie_exit(&t, h);
69         while (t.m[k].l != -1)
70         {
71             for (w = t.m[k].l; w != -1; w = t.w[w].n)
72                 printf(" '%s' found at %d\n", p[t.w[w].v], j - l[t.w[w].v]);
73             k = trie_exit(&t, trie_suffix(&t, k));
74         }
75         h = trie_step(&t, h, val(*s));
76         if (*s++ == '\0') break;
77     }
78     return i;
79 }
80
81 int main()
82 {
83     int i, j, n;

```

- ▶ Gerum ráð fyrir að strengirnar í  $p$  komi fyrir  $k$  sinnum í  $s$ .
- ▶ Þá er tímaflækjan  $\mathcal{O}(|s| + |\Sigma| \cdot |p| + k)$ .
- ▶ Ef við höfum bara áhuga á að finna töluna  $k$  getum við breytt `trie_exit(...)` þannig að það reikni fjölda lokastaða á leiðinni að rót eftir bakstrengsleggjum.
- ▶ Þá verður tímaflækjan  $\mathcal{O}(|s| + |\Sigma| \cdot |p|)$ .
- ▶ Takið eftir að ef stafrófið er takmarkað þá er seinni tímaflækjan línuleg.

