

UNIVERZITET U NIŠU
PRIRODNO-MATEMATIČKI FAKULTET
DEPARTMAN ZA RAČUNARSKE NAUKE



Algoritmi za pretragu stringova

MASTER RAD

Student:

Ivan Stošić

Mentor:

Marko Petković

Niš
Septembar, 2019.

Sadržaj

1	Uvod	3
1.1	Uvodne definicije	3
1.2	Složenost algoritama	4
1.3	Implementacije algoritama	5
2	Osnovni algoritmi pretrage	6
2.1	Naivna pretraga	6
2.2	Knuth-Morris-Pratt algoritam	7
2.3	Z-algoritam	10
3	Dinamičko programiranje nad stringovima	13
3.1	LCP matrica	13
3.2	Najduži zajednički podniz	14
3.2.1	Hirschberg-ov algoritam	16
3.3	Najduži palindromski podniz	19
3.4	Levenshtein udaljenost	20
4	Prefiksne strukture podataka	23
4.1	Prefiksno stablo	23
4.1.1	Konstrukcija	24
4.1.2	Primene	24
4.2	Aho-Corasick algoritam	26
5	Sufiksne strukture podataka	32
5.1	Sufiks niz	32
5.1.1	Konstrukcija	33
5.1.2	Brži algoritmi za konstrukciju	34
5.1.3	LCP niz	38
5.1.4	Primene	40
5.2	Sufiks stablo	43
5.3	Sufiks automat	47
5.3.1	Definicija	47
5.3.2	Algoritam za konstrukciju	49
5.3.3	Uopštenja i primene	52
6	Heširanje	54
6.1	Rabin-Karp algoritam	54
6.2	Kolizije	57
6.2.1	Odabir parametara heš funkcije	57
6.2.2	Konstrukcija kolizije	60

7	Palindromi	65
7.1	Manacher-ov algoritam	66
7.2	Palindromsko stablo	67
7.2.1	Algoritam za konstrukciju	68
7.2.2	Primene	70

1 Uvod

1.1 Uvodne definicije

Neformalno, string je niz simbola iz nekog alfabeta. U opštem slučaju, alfabet može biti bilo koji konačan skup. Međutim, za potrebe pojedinih algoritama, potrebno je da alfabet bude totalno uređen skup. Kod implementacije algoritama, dodatno možemo pretpostaviti da je alfabet jednak nekom konačnom skupu uzastopnih celih brojeva, najčešće $\{0, 1, \dots, k-1\}$ ili $1, 2, \dots, k$ za neko $k \in \mathbb{N}$. Tradicionalno, ovaj alfabet se označava grčkim slovom Σ .

Definicija 1.1 *Ako je Σ alfabet a n prirodan broj, Σ^n označava skup svih uređenih n -torki $(s_0, s_1, \dots, s_{n-1})$, gde je $s_i \in \Sigma$ za svako $i \in \{0, \dots, n-1\}$.*

Ovakvu n -torku možemo kraće zapisati sa s , a njenu dužinu (broj elemenata) sa $|s|$. Koristi se i kraći zapis n -torke: $s_0 s_1 \dots s_{n-1}$.

Definicija 1.2 *Ako je Σ alfabet, onda je*

$$\Sigma^+ = \bigcup_{n=1}^{\infty} \Sigma^n$$

skup svih nepraznih reči nad alfabetom Σ .

Ovom skupu možemo dodati i praznu reč, koju označavamo sa e ili ϵ .

Definicija 1.3 *Skup svih reči nad alfabetom Σ je skup $\Sigma^* = \Sigma^+ \cup \{e\}$.*

Za string možemo definisati i njegove podstringove na sledeći način.

Definicija 1.4 *Podstring počev od pozicije l , do pozicije r ne uključujući r , nekog stringa s , gde važi $0 \leq l \leq r \leq |s|$ je string $s_l s_{l+1} \dots s_{r-1}$. Ovaj podstring kraće zapisujemo i $s_{[l,r)}$.*

Podstringove stringa s kod kojih je $l = 0$ nazivamo prefiksima tog stringa, dok podstringove kod kojih je $r = |s|$ nazivamo sufiksima tog stringa. Ukoliko je $l = r$ radi se o praznom podstringu. Ukoliko je podstring različit od celog stringa, onda se radi o pravom podstringu, sufiksu odnosno prefiksu. Izbor notacije sa indeksiranjem od nule i poluotvorenim intervalom olakšava implementaciju većine algoritama sa stringovima.

Definicija 1.5 *Ciklični podstring stringa s dužine n počev od pozicije l do pozicije r je string $s_{[l,r)} = s_l \bmod n s_{(l+1) \bmod n} \dots s_{(r-1) \bmod n}$.*

Ciklični podstring uopštava pojam podstringa. Zaista, ako je $0 \leq l \leq r \leq n$, ciklični podstring jednak je običnom podstringu.

Definicija 1.6 *Ciklični pomerač stringa s dužine n počev od pozicije i je string $s_{[i,i+n)}$.*

Stringovi se mogu i nadovezivati odnosno konkatenerirati. Skup Σ^+ , odnosno Σ^* zajedno sa operacijom konkatencije čini algebarsku strukturu polugrupe, odnosno monoida.

Definicija 1.7 *Ako su s, p stringovi, tada je njihova konkatencija string $sp = s_0 s_1 \dots s_{|s|-1} p_0 p_1 \dots p_{|p|-1}$.*

Definicija 1.8 *Ako je s string dužine n , sa \bar{s} označavamo string $s_{n-1} s_{n-2} \dots s_0$.*

Definicija 1.9 *String s je palindrom ako važi $s = \bar{s}$.*

Ukoliko je skup Σ totalno uređen, definišemo leksikografsko poređenje stringova kao uređenje skupa Σ^* , na sledeći način.

Definicija 1.10 *Za string s kažemo da je leksikografski manji od stringa p ukoliko postoji ceo broj $k \geq 0$ takav da je $k < \min\{|s|, |p|\}$, $s_{[0,k)} = p_{[0,k)}$ i $s_k < p_k$ ili ako je s pravi prefiks stringa p .*

Teorema 1.1 *Leksikografsko poređenje je totalno uređenje skupa Σ^* . \square*

Definicija 1.11 *Za svako $x \in \Sigma$, $\text{Ord}(x)$ je broj elemenata skupa Σ koji su strogo manji od x .*

1.2 Složenost algoritama

Vreme, odnosno broj koraka i količina utrošene memorije tokom izvršenja nekog algoritma zavisi od ulaznih parametara. *Veliko* O notacija nam olakšava opisivanje i izračunavanje ovih funkcionalnih zavisnosti. Neka je u narednim definicijama domen funkcija f, g skup \mathbb{N}_0 a kodomen $\mathbb{R}^+ \cup \{0\}$.

Definicija 1.12 *Skup $O(g)$ definišemo kao skup svih funkcija f za koje važi da postoje konstante c i n_0 takve da je $f(n) \leq cg(n)$ za svako $n \geq n_0$.*

Ovu notaciju koristimo kada želimo da opišemo gornju granicu neke funkcije, do na proizvod sa konstantom. Problem ove notacije je upravo u tome što samo daje gornju granicu ponašanja neke funkcije. Zato se uvodi Θ -notacija.

Definicija 1.13 *Skup $\Theta(g)$ definišemo kao skup svih funkcija f za koje važi da postoje pozitivne konstante c_1, c_2 i n_0 takve da je $c_1g(n) \leq f(n) \leq c_2g(n)$ za svako $n \geq n_0$.*

Za algoritam čiji je ulazni parametar n , što može biti broj elemenata niza, broj vrsta matrice, broj čvorova grafa, itd. kažemo da ima vremensku složenost $\Theta(g(n))$ odnosno $O(g(n))$ ako je f , gde je $f(n)$ broj elementarnih koraka tokom izvršenja algoritma, u skupu $\Theta(g)$ odnosno $O(g)$. Slično definišemo memorijsku složenost preko broja iskorišćenih elementarnih memorijskih lokacija.

1.3 Implementacije algoritama

Za sve implementacije biće korišćen programski jezik C++, kompajler GCC, verzija 9.1.0. Verzija standarda jezika biće C++17. Radi jednostavnosti, kodovi će biti dati bez `main` funkcije, `#include` direktiva i naredbe `using namespace std`; . Pretpostavlja se da su uključene sve standardne biblioteke, što se kod GCC-a može postići sa `#include <bits/stdc++.h>`.

2 Osnovni algoritmi pretrage

Jedan od osnovnih problema pretrage stringova je problem nalaženja svih pojavljivanja jednog stringa u drugom. Formalno, za data dva stringa s, p dužina n i m , redom, naći sve indekse i takve da je $0 \leq i \leq n-m$ i $s_{[i, i+m)} = p$. String s unutar kojeg se traži se često naziva tekstom, dok se p naziva rečju (iako ne mora biti jedna reč) ili *pattern*-om odnosno šablonom. U literaturi na engleskom jeziku se često sreću i pojmovi *haystack* (plast sena) i *needle* (igla).

2.1 Naivna pretraga

Naivna pretraga rešava prethodno opisani problem tako što za svako celo i iz segmenta $[0, n - m]$ upoređuje karakter po karakter stringove $s_{[i, i+m)}$ i p , pri čemu se odmah zaustavlja ukoliko naiđe na dva različita karaktera. Ovo poređenje ima vremensku složenost $O(m)$ pa ceo algoritam ima vremensku složenost $O((n - m)m)$, ova složenost se dostiže npr. za stringove koji se sastoje samo od slova a.

Implementacija naive pretrage

```
1 vector<int> naive_search(const string& s, const string& p) {
2     int n = s.size(), m = p.size();
3     vector<int> result;
4     for (int i=0; i<=n-m; i++) {
5         bool ok = true;
6         for (int j=0; j<m; j++) {
7             if (s[i+j] != p[j]) {
8                 ok = false;
9                 break;
10            }
11        }
12        if (ok)
13            result.push_back(i);
14    }
15    return result;
16 }
```

2.2 Knuth-Morris-Pratt algoritam

KMP je prvi otkriveni algoritam koji rešava problem pretrage stringa u linearnom vremenu, odnosno u složenosti $O(n + m)$.¹ Da bismo razumeli rad algoritma, definišimo sledeće pojmove.

Definicija 2.1 *Sufiks-prefiks nepraznog stringa s je svaki string p različit od s , uključujući i prazan string, koji je istovremeno i sufiks i prefiks stringa s , tj. $s_{[0,|p|)} = s_{[|s|-|p|,|s|)} = p$.*

Označimo sa $g(s)$ dužinu najdužeg sufiks-prefiksa stringa s . Za prazan string s definišemo $g(s) = -1$. Ova funkcija zadovoljava sledeću, važnu osobinu:

Teorema 2.1 *Neka je s neprazan string. Neka je $s' = s_{[0,|s|-1)}$, tada je $g(s) \leq g(s') + 1$.*

Dokaz. Pretpostavimo suprotno, da je $g(s) > g(s') + 1$. Neka je p najduži sufiks prefiks stringa s . Neka je $p' = p_{[0,|p|-1)}$. Tada je p' sufiks-prefiks stringa s' , pa je $g(s') \geq |p'| = |p| - 1 = g(s) - 1$, kontradikcija.

Ukoliko je p najduži sufiks-prefiks stringa s , tada je i svaki drugi sufiks-prefiks stringa s ujedno i sufiks-prefiks stringa p . Ova osobina nam omogućava da opišemo sve sufiks-prefikse nekog stringa s kao lanac, gde je svaki naredni string najduži sufiks-prefiks prethodnog, sve dok se ne dođe do praznog stringa.

Definicija 2.2 *Niz neuspeha stringa s je niz $f_0, f_1, \dots, f_{|s|}$, gde je $f_i = g(s_{[0,i)})$.*

Prvi korak KMP algoritma je nalaženje niza neuspeha za string p . Prvo upisujemo $f_0 = -1$. Zatim, za svako $i > 0$, tražimo najduži sufiks-prefiks stringa $p_{[0,i-1)}$ koji se može proširiti slovom p_{i-1} , odnosno, nalazimo najveći broj r takav da je r dužina nekog sufiks-prefiksa stringa $p_{[0,i-1)}$ i važi $p_r = p_{i-1}$. Ukoliko ne nađemo takav broj r , onda je $f_i = 0$. U suprotnom, pošto je $p_{[0,r)} = p_{[i-1-r,i-1)}$ i $p_r = p_{i-1}$ imamo da je $p_{[0,r+1)} = p_{[i-r-1,i)}$, odnosno, $r + 1$ je dužina jednog sufiks-prefiksa stringa $p_{[0,i)}$. Iz prethodnog razmatranja imamo da je ovo upravo najduži sufiks-prefiks stringa $p_{[0,i)}$, pa upisujemo $f_i = r + 1$. Niz svih sufiks-prefiksa možemo naći primenom opisane osobine lanca i činjenice da smo u i -tom koraku već izračunali dužine najdužih sufiks-prefiksa za sve prefikse stringa p dužine manje od i .

Primer. Posmatrajmo string `atamatata`. Njegov niz neuspeha dat je u sledećoj tabeli. Indeksi u donjem redu su pomereni za jedno mesto da bi se bolje videlo o kom prefiksu je reč.

	a	t	a	m	a	t	a	t	a
-1	0	0	1	0	1	2	3	2	3

Posmatrajmo šta se dešava pri računanju f_i za pretposlednji prefiks, dužine 8. U prethodnom koraku smo imali string `atamata`, kome je najduži sufiks-prefiks string `ata`. Ako pokušamo da dodamo na njega slovo `t`, nećemo dobiti poklapanje jer je na poziciji 3 slovo `m`, zato pokušavamo sa narednim sufiks-prefiksom odnosno $r = f_3 = 1$. Sada uspešno dodajemo slovo `t` i upisujemo $r + 1 = 2$ u f_8 .

Implementacija nalaženja niza neuspeha kod KMP algoritma

```

1  vector<int> kmp_ff(const string& p) {
2      int m = p.size();
3      vector<int> f(m+1);
4      f[0] = -1;
5      for (int i=1; i<=m; i++) {
6          int r = f[i-1];
7          while (r != -1 && p[r] != p[i-1])
8              r = f[r];
9          f[i] = r+1;
10     }
11     return f;
12 }
```

Ocenimo složenost ovog algoritma. Jasno je da složenost srazmerna zbiru broja iteracija ove dve petlje. Posmatrajmo vrednost $2i - r$. Nakon svake iteracije *while* petlje imamo da se r smanjuje, jer je $f_r < r$ pa se $2i - r$ povećava. U jednoj iteraciji *for* petlje imamo da se i i r povećavaju za tačno 1, pa se $2i - r$ ponovo povećava. Kako je početna vrednost $2i - r$ jednaka 3, a važi $2i - r \leq 2m + 1$, ukupan broj iteracija, a samim tim i složenost algoritma je $O(m)$.

Opišimo sada glavni algoritam za traženje stringa. Algoritam za svaki prefiks i stringa s nalazi najduži sufiks koji je ujedno i prefiks stringa p . Ukoliko je dužina tog prefiksa jednaka $|p|$, tada dolazi do poklapanja na poziciji $i - |p|$. Neka je ta dužina jednaka r_i . Za $i = 0$ imamo $r_0 = 0$. Za svako $i > 0$, krećemo od prefiksa stringa p dužine r_{i-1} i proveravamo da li je naredno slovo jednako slovu s_{i-1} . Ako jeste, zaustavljamo se i upisujemo

dužinu nađenog prefiksa u r_i , u suprotnom, sledeći prefiks koji pokušavamo da produžimo je prefiks dužine $f_{r_{i-1}}$. Ovaj postupak ponavljamo sve dok ne dođemo do poklapanja ili do fiktivnog prefiksa -1 , u tom slučaju dužina poklapanja je 0 . Primetimo da je ovaj deo algoritma veoma sličan nalaženju niza neuspeha.

Implementacija glavnog dela KMP algoritma

```

1  vector<int> kmp_main(const string& s, const string& p) {
2      vector<int> f = kmp_ff(p), result;
3      int n = s.size(), m = p.size(), r = 0;
4      for (int i=1; i<=n; i++) {
5          while (r != -1 && p[r] != s[i-1])
6              r = f[r];
7          r++;
8          if (r == m) {
9              result.push_back(i-m);
10             r = f[r];
11         }
12     }
13     return result;
14 }
```

Na potpuno isti način se ocenjuje složenost glavnog dela algoritma, posmatranjem vrednosti izraza $2i - r$. Složenost je $O(n + m)$, zajedno sa prvom fazom, ukupna složenost je takođe $O(n + m)$.

Cela implementacija algoritma se može značajno pojednostaviti na sledeći način: Posmatrajmo string $p\$s$, gde je $\$$ karakter koji se ne javlja u stringovima p, s . Na osnovu njegovog niza neuspeha možemo da zaključimo gde se sve pojavljuje p u s , tačnije, i je pojavljivanje p u s akko je $f_{i+2m+1} = m$.

Pojednostavljena implementacija celog algoritma

```
1 vector<int> kmp_simple(const string& s, const string& p) {
2     string q = p + '\0' + s;
3     int n = s.size(), m = p.size();
4     vector<int> f(n+m+2), result;
5     f[0] = -1;
6     for (int i=1; i<=n+m+1; i++) {
7         int r = f[i-1];
8         while (r != -1 && q[r] != q[i-1])
9             r = f[r];
10        f[i] = ++r;
11        if (r == m && i >= 2*m+1)
12            result.push_back(i-2*m-1);
13    }
14    return result;
15 }
```

2.3 Z-algoritam

Osnovna ideja ovog algoritma je da se izračuna Z-niz, koji se definiše na sledeći način.

Definicija 2.3 Za string s dužine n , Z-niz je niz z_1, \dots, z_{n-1} gde je z_i dužina najdužeg zajedničkog prefiksa za stringove s i $s_{[i,n]}$.

Z-niz se može iskoristiti za pretragu stringova. Ukoliko nađemo Z-niz za string ps , i je pojavljivanje p u s akko je $z_{i+m} \geq m$.

Z-algoritam je efikasan algoritam za nalaženje Z-niza.² Označimo sa q string za koji nalazimo Z-niz. Algoritam za svako i direktno izračuna z_i poklapanjem slova na pozicijama $z_i, i + z_i$. Ključna ideja je da se prethodno izračunate vrednosti Z-niza mogu iskoristiti da se postavi bolja početna vrednost za z_i , umesto da se svaki put kreće od nule. Naime, neka je r najveća nađena vrednost izraza $j + z_j$ za $j < i$, a neka je pritom $l = j$ za koje se dostiže taj maksimum. Drugim rečima, $[l, r)$ je prozor koji odgovara nekom do sada pronađenom podstringu koji je jednak nekom prefiksu stringa, i to onom kod kojeg je r najveće. Ukoliko važi $l \leq i < r$, tada znamo da je $q_{[i,r)} = q_{[i-l,r-l)}$. Takođe, važi $q_{[i-l,i-l+z_{i-l})} = q_{[0,z_{i-l})}$ pa, ako označimo sa $t_i = \min(z_{i-l}, r - i)$, važi $q_{[i,i+t_i)} = q_{[0,t_i)}$ odnosno $z_i \geq t_i$.

Implementacija Z-algoritma za pretragu stringa

```
1 vector<int> z_algorithm(const string& s, const string& p) {
2     int n = s.size(), m = p.size();
3     string q = p + s;
4     vector<int> z(n+m, 0), result;
5     for (int i=1, l=0, r=0; i<n+m; i++) {
6         if (i < r)
7             z[i] = min(z[i-l], r-i);
8         while (i+z[i] < n+m && q[i+z[i]] == q[z[i]])
9             z[i]++;
10        if (i+z[i] > r)
11            l = i, r = i+z[i];
12        if (z[i] >= m)
13            result.push_back(i-m);
14    }
15    return result;
16 }
```

Dokažimo da ovaj algoritam ima složenost $O(n+m)$. Očigledno, kritična je unutrašnja *while* petlja. Dokazaćemo da svaka iteracija *while* petlje odgovara povećanju vrednosti promenljive r za bar 1. Posmatrajmo sledeće slučajeve:

- Pre ulaska u *while* petlju važi $i \geq r$. U ovom slučaju z_i ima početnu vrednost nula, na kraju *while* petlje će važiti $i + z_i \geq r$ pa će r dobiti vrednost $i + z_i$, odnosno r će se povećati za barem z_i , što je veće ili jednako od broja iteracija *while* petlje.
- Pre ulaska u *while* petlju važi $i < r$ i $z_{i-l} < r - i$. Sada z_i dobija početnu vrednost z_{i-l} . Dokažimo da će *while* petlja izvršiti tačno nula iteracija. Pretpostavimo suprotno, da je $q_{i+z_i} = q_{z_i}$, odnosno $q_{i+z_{i-l}} = q_{z_{i-l}}$. Iz definicije prozora $[l, r)$ imamo da je $q_{[0, r-l)} = q_{[l, r)}$, pošto je $l \leq i + z_{i-l} < r$, odnosno, ova pozicija je unutar prozora, važi $q_{i+z_{i-l}} = q_{i-l+z_{i-l}}$, odnosno, po pretpostavci, $q_{z_{i-l}} = q_{i-l+z_{i-l}}$, što znači da z_{i-l} ima pogrešno izračunatu vrednost jer se poklapaju karakteri na kraju odgovarajućih podstringova, što dovodi do kontradikcije.
- Pre ulaska u *while* petlju važi $i < r$ i $z_{i-l} \geq r - i$. Sada z_i dobija početnu vrednost $r - i$. Ako petlja izvrši k iteracija, na kraju će važiti $z_i = k + r - i$ odnosno $i + z_i = k + r$, što znači da će nova vrednost r biti za bar k veća.

Kako je $r \leq n + m$ zaključujemo da je ukupna složenost $O(n + m)$.

3 Dinamičko programiranje nad stringovima

Dinamičko programiranje je tehnika rešavanja optimizacionih problema i problema prebrojavanja gde se glavni problem rešava tako što se identifikuju slični potproblemi manje veličine, koji se zatim rešavaju i čija se rešenja kombinuju u rešenje glavnog problema. Svaki ovako dobijeni potproblem se rešava najviše jedanput, nakon čega se njegovo rešenje pamti u memoriji.

Prvi korak u primeni dinamičkog programiranja na rešavanje nekog problema jeste da se identifikuju *potproblemi* tog problema. Najveća instanca među svim potproblemima jeste *glavni problem*. Najmanje instance su *trivijalni potproblemi*, koji se ne dele dalje na potprobleme i čija rešenja se dobijaju na neki drugi način. Zatim je neophodno, za svaki potproblem naći relaciju između njega i jednog ili više manjih potproblema. Ova relacija odnosno rešenje za potproblem je matematički izraz ili rezultat nekog jednostavnog algoritma u kojem figurišu rešenja manjih potproblema. Kažemo da potproblem A zavisi od potproblema B ako rešenje potproblema B figuriše u izrazu koji je rešenje potproblema A . Da bismo našli rešenje nekog potproblema neophodno je da prethodno nađemo rešenja za sve potprobleme od kojih on zavisi.

3.1 LCP matrica

Jedna od jednostavnijih primena dinamičkog programiranja jeste izračunavanje LCP matrice stringa. $LCP(s, p)$ (od *longest common prefix*) je dužina najdužeg zajedničkog prefiksa stringova s, p .

Definicija 3.1 *LCP matrica za string s dužine n je kvadratna matrica A za koju važi $A_{i,j} = LCP(s_{[i,n]}, s_{[j,n]})$.*

Ukoliko važi $s_i \neq s_j$, onda je $A_{i,j} = 0$. U suprotnom, posmatrajmo stringove $s_{[i+1,n]}, s_{[j+1,n]}$. Važi da je p njihov zajednički prefiks ako i samo ako je $s_i p = s_j p$ zajednički prefiks stringova $s_{[i,n]}, s_{[j,n]}$, odakle dobijamo da je $A_{i,j} = A_{i+1,j+1} + 1$, pod uslovom da su indeksi $i+1, j+1$ validni, inače se radi o praznim podstringovima pa možemo smatrati da je $A_{i+1,j+1} = 0$ odnosno $A_{i,j} = 1$. Ovo nas dovodi do sledećeg, jednostavnog algoritma za računanje LCP matrice.

Implementacija algoritma za nalaženje LCP matrice

```
1 vector<vector<int>> lcp_matrix(const string& s) {  
2     int n = s.size();  
3     vector<vector<int>> a(n, vector<int>(n));  
4     for (int i=n-1; i>=0; i--)  
5         for (int j=n-1; j>=0; j--)  
6             if (s[i] != s[j])  
7                 a[i][j] = 0;  
8             else if (i != n-1 && j != n-1)  
9                 a[i][j] = a[i+1][j+1] + 1;  
10            else  
11                a[i][j] = 1;  
12    return a;  
13 }
```

Vremenska i memorijska složenost ovog algoritma je $O(n^2)$.

3.2 Najduži zajednički podniz

Definicija 3.2 Podniz stringa s dužine n je string $s_{i_0}s_{i_1}\dots s_{i_{k-1}}$ gde važi $0 \leq s_0 < s_1 < \dots < i_{k-1} < n$.

Dužina najdužeg zajedničkog podniza (LCS, od *longest common subsequence*) dva stringa s, p se može koristiti kao mera njihove sličnosti. Preciznije, ukoliko su dužine s, p redom n, m a dužina njihovog LCS-a je k , tada je njihova udaljenost $n + m - 2k$, gde se udaljenost odnosi na minimalan broj izmena potrebnih da se od stringa s dobije string p , gde su dozvoljene operacije brisanje jednog slova i umetanje jednog slova u string.

Neka je za fiksne stringove s, p dužina n, m , redom, $d_{i,j}$ dužina LCS-a za stringove $s_{[0,i)}, p_{[0,j)}$.

- Ako je $i = 0$ ili $j = 0$, tada je $d_{i,j} = 0$.
- U suprotnom, posmatrajmo slova s_{i-1} i p_{j-1} . Ukoliko su ona jednaka, tada se svaki zajednički podniz stringova $s_{[0,i-1)}, p_{[0,j-1)}$ može proširiti za karakter $s_{i-1} = p_{j-1}$ tako da se dobije zajednički podniz stringova $s_{[0,i)}, p_{[0,j)}$. Inače, svaki zajednički podniz stringova $s_{[0,i)}, p_{[0,j)}$ je ili zajednički podniz za $s_{[0,i-1)}, p_{[0,j)}$ ili za $s_{[0,i)}, p_{[0,j-1)}$.

Oдавde dobijamo sledeću rekurentnu vezu: Ako su $i, j > 0$,

$$d_{i,j} = \begin{cases} \max(d_{i-1,j}, d_{i,j-1}) & s_{i-1} \neq p_{j-1} \\ \max(d_{i-1,j}, d_{i,j-1}, d_{i-1,j-1} + 1) & s_{i-1} = p_{j-1} \end{cases} \quad (3.2.1)$$

Jasno je da je $d_{n,m}$ dužina LCS-a za cele stringove s, p .

Implementacija nalaženja dužine LCS-a

```

1  int lcs_len(const string& s, const string& p) {
2      int n = s.size(), m = p.size();
3      vector<vector<int>> d(n+1, vector<int>(m+1));
4      for (int i=0; i<=n; i++)
5          for (int j=0; j<=m; j++)
6              if (i == 0 || j == 0)
7                  d[i][j] = 0;
8              else if (s[i-1] != p[j-1])
9                  d[i][j] = max(d[i-1][j], d[i][j-1]);
10             else
11                 d[i][j] = max({d[i-1][j], d[i][j-1], d[i-1][j-1]+1});
12     return d[n][m];
13 }
```

Ako posmatramo d kao matricu, primećujemo da se vrednosti u i -toj vrsti mogu izračunati samo na osnovu stringova s, p i vrednosti u i -toj i vrsti $i-1$. Ovo nam omogućava da implementiramo algoritam tako da se u svakom trenutku pamte samo poslednje dve vrste.

Implementacija nalaženja dužine LCS-a sa $O(m)$ memorije

```

1  int lcs_len_mem(const string& s, const string& p) {
2      int n = s.size(), m = p.size();
3      vector<int> di(m+1, 0), dim1(m+1);
4      for (int i=1; i<=n; i++) {
5          swap(di, dim1);
6          for (int j=1; j<=m; j++)
7              if (s[i-1] != p[j-1])
8                  di[j] = max(dim1[j], di[j-1]);
9              else
10                 di[j] = max({dim1[j], di[j-1], dim1[j-1]+1});
11     }
12     return di[m];
13 }
```


Ukoliko je potrebno naći ceo podniz a ne samo njegovu dužinu, rekonstrukciju radimo kretanjem unazad kroz matricu d , uvek idući ka odgovarajućem polju koje ima najveću vrednost, odnosno, ako smo u polju i, j idemo ka polju od kojeg je $d_{i,j}$ "uzelo" vrednost.

Implementacija nalaženja LCS-a

```

1  string lcs_string(const string& s, const string& p) {
2      int n = s.size(), m = p.size();
3      vector<vector<int>> d(n+1, vector<int>(m+1));
4      for (int i=0; i<=n; i++)
5          for (int j=0; j<=m; j++)
6              if (i == 0 || j == 0)
7                  d[i][j] = 0;
8              else if (s[i-1] != p[j-1])
9                  d[i][j] = max(d[i-1][j], d[i][j-1]);
10             else
11                 d[i][j] = max({d[i-1][j], d[i][j-1], d[i-1][j-1]+1});
12
13     int i = n, j = m;
14     string q;
15     while (i > 0 && j > 0) {
16         if (s[i-1] == p[j-1]) {
17             if (d[i][j] == d[i-1][j-1] + 1)
18                 q += s[i-1], i--, j--;
19             else if (d[i][j] == d[i-1][j])
20                 i--;
21             else
22                 j--;
23         } else {
24             if (d[i][j] == d[i-1][j])
25                 i--;
26             else
27                 j--;
28         }
29     }
30     reverse(q.begin(), q.end());
31     return q;
32 }

```

3.2.1 Hirschberg-ov algoritam

Iako se algoritam koji samo nalazi dužinu LCS-a dva stringa suštinski ne razlikuje od onog koji nalazi ceo taj podniz, prvi se može jednostavno realizovati tako da mu je memorijska složenost $O(m)$. Hirschberg-ov algoritam nalazi ceo

LCS u memorijskoj složenosti $O(n+m)$ bez žrtvovanja vremenske složenosti.³ Ideja algoritma je da string s predstavimo kao $s = s_1s_2$ gde s_1, s_2 imaju približno jednake dužine, a da zatim nađemo predstavljanje stringa $p = p_1p_2$ takvo da je $|LCS(s_1, p_1)| + |LCS(s_2, p_2)| = |LCS(s, p)|$, odnosno, ako posmatramo LCS za s, p slova iz s_1 su uparena tačno sa slovima iz p_1 i slova iz s_2 su uparena tačno sa slovima iz p_2 . Znajući particije ovih stringova, rekursivno nalazimo $LCS(s_1, p_1)$ i $LCS(s_2, p_2)$ a zatim konkatenujemo rezultate.

Opišimo prvo pomoćni algoritam koji za stringove s, p dužina n, m nalazi, za svako $j \in \{0, 1, \dots, m\}$ vrednost $|LCS(s, p_{[0,j]})|$. Primetimo da je ovaj algoritam identičan algoritmu koji nalazi dužinu LCS-a u $O(m)$ memorije, osim što vraća ceo vektor a ne samo njegov poslednji element.

Pomoćna funkcija Hirschberg-ovog algoritma

```

1  vector<int> lcs_vector(const string& s, const string& p) {
2      int n = s.size(), m = p.size();
3      vector<int> di(m+1, 0), dim1(m+1);
4      for (int i=1; i<=n; i++) {
5          swap(di, dim1);
6          for (int j=1; j<=m; j++)
7              if (s[i-1] != p[j-1])
8                  di[j] = max(dim1[j], di[j-1]);
9              else
10                 di[j] = max({dim1[j], di[j-1], dim1[j-1]+1});
11         }
12     return di;
13 }
```

Nađimo ovaj vektor za parove stringova s_1, p i $\overline{s_2}, \overline{p}$, neka su to vektori v_1, v_2 . Ako je $i \in \{0, 1, \dots, m\}$, tada je $v_1(i) + v_2(m-i)$ dužina LCS-a koji odgovara predstavljanju $p = p_1p_2$ sa $p_1 = p_{[0,i]}, p_2 = p_{[i,m]}$, pa maksimiziranjem prethodnog izraza po i nalazimo traženu particiju za p .

Implementacija Hirschberg-ovog algoritma

```
1 string hirschberg(const string& s, const string& p) {
2     if (s.size() > p.size())
3         return hirschberg(p, s);
4     int n = s.size(), m = p.size();
5     if (n == 0)
6         return string();
7     if (n == 1) {
8         if (p.find(s[0]) != string::npos)
9             return string(1, s[0]);
10        else
11            return string();
12    }
13    string s1 = s.substr(0, n/2), s2 = s.substr(n/2);
14    string p_rev = p, s2_rev = s2;
15    reverse(p_rev.begin(), p_rev.end());
16    reverse(s2_rev.begin(), s2_rev.end());
17    vector<int> v1 = lcs_vector(s1, p);
18    vector<int> v2 = lcs_vector(s2_rev, p_rev);
19    int i_best = 0;
20    for (int i=1; i<=m; i++)
21        if (v1[i] + v2[m-i] > v1[i_best] + v2[m-i_best])
22            i_best = i;
23    return hirschberg(s1, p.substr(0, i_best))
24        + hirschberg(s2, p.substr(i_best));
25 }
```

Uzećemo da je s duži string. Ukoliko nije, rekursivno zovemo istu funkciju gde parametri menjaju mesto. Koristimo matematičku indukciju da dokažemo da je memorijska složenost algoritma $O(n + m)$. Indukciju radimo po zbiru $n + m$. Dokazaćemo da postoji realan broj c_2 takav da algoritam koristi ne više od $c_2(n + m)$ bajtova memorije. Za $n + m \leq 1$ tvrđenje očigledno važi jer nema rekursivnih poziva. U suprotnom, telo funkcije koristi ne više od $c_1(n + m)$ bajtova dodatne memorije, dok je kod rekursivnih poziva zbir dužina stringova ne više od $\frac{n}{2} + m \leq \frac{3}{4}(n + m)$ (ovo važi jer je $m \leq n$), odnosno, ako uzmemo da je $c_2 = 4c_1$, važi da algoritam za dužine n, m koristi ne više od $c_2 \cdot \frac{3}{4}(n + m) + c_1(n + m) = 4c_1(n + m) = c_2(n + m)$ memorije, čime završavamo indukcijski korak.

Procenimo sada vremensku složenost algoritma. Označimo sa $H = nm$. Algoritam za računanje nizova v_1, v_2 koristi $O(H)$ vremena, dok rekursivni pozivi zajedno imaju veličinu $\frac{n}{2} \cdot m = \frac{H}{2}$, pa vremenska složenost izražena preko H zadovoljava relaciju $T(H) = O(H) + T(\frac{H}{2})$, pa je na osnovu Master

teoreme vremenska složenost upravo $O(H)$ odnosno $O(nm)$.

3.3 Najduži palindromski podniz

Definicija 3.3 Za dati string s dužine n , najduži palindromski podniz je palindrom najveće dužine koji se javlja kao podniz stringa s .

Problem nalaženja najdužeg palindromskog podniza se efikasno može rešiti pomoću dinamičkog programiranja. Nađimo dužinu najdužeg palindromskog podniza za svaki podstring stringa s , preciznije, neka je $d_{l,r}$ za $0 \leq l < r \leq n$ dužina najdužeg palindromskog podniza za string $s_{[l,r]}$. Posmatrajmo sledeće slučajeve:

- $r - l = 1$. Tada string $s_{[l,r]}$ sadrži jedno slovo pa je $d_{l,r} = 1$.
- $r - l = 2$. Tada string $s_{[l,r]}$ sadrži dva slova, ukoliko su ona jednaka $d_{l,r} = 2$, inače je $d_{l,r} = 1$.
- $r - l > 2, s_l \neq s_{r-1}$. Svaki palindromski podniz stringa $s_{[l,r]}$ je sigurno sadržan u celosti ili u $s_{[l,r-1]}$ ili u $s_{[l+1,r]}$, pa je $d_{l,r} = \max(d_{l+1,r}, d_{l,r-1})$.
- $r - l > 2, s_l = s_{r-1}$. Svaki palindromski podniz stringa $s_{[l,r]}$ je sigurno sadržan u celosti ili u $s_{[l,r-1]}$ ili u $s_{[l+1,r]}$ ili počinje na poziciji l , završava se na poziciji $r - 1$, a ostatak, koji je takođe palindrom, je u celosti sadržan u $s_{[l+1,r-1]}$ pa je $d_{l,r} = \max(d_{l+1,r}, d_{l,r-1}, d_{l+1,r-1} + 2)$.

Rešenje glavnog problema, odnosno dužina najdužeg palindromskog podniza je po definiciji $d_{0,n}$. Sâmo rešenje se može rekonstruisati sličnim postupkom kao kod nalaženja LCS-a.

Implementacija algoritma za nalaženje najdužeg palindromskog podniza

```
1 string lpalsubseq(const string& s) {
2     int n = s.size();
3     vector<vector<int>> d(n, vector<int>(n+1));
4     for (int l=n-1; l>=0; l--) {
5         d[l][l+1] = 1;
6         if (l+2 <= n)
7             d[l][l+2] = 1 + (s[l] == s[l+1]);
8         for (int r=l+3; r<=n; r++) {
9             d[l][r] = max(d[l+1][r], d[l][r-1]);
10            if (s[l] == s[r-1])
11                d[l][r] = max(d[l][r], d[l+1][r-1]+2);
12        }
13    }
14    int l = 0, r = n;
15    string prefix, middle;
16    while (r-l > 2) {
17        if (s[l] == s[r-1] && d[l+1][r-1]+2 == d[l][r])
18            prefix += s[l], l++, r--;
19        else if (d[l+1][r] == d[l][r])
20            l++;
21        else
22            r--;
23    }
24    if (r-l == 2 && s[l] == s[l+1])
25        middle = s.substr(l, 2);
26    else
27        middle = string(1, s[l]);
28    string suffix = prefix;
29    reverse(suffix.begin(), suffix.end());
30    return prefix + middle + suffix;
31 }
```

Vremenska i memorijska složenost algoritma je $O(n^2)$. Ukoliko je potrebna samo dužina, dovoljno je vratiti vrednost $d_{0,n}$ nakon kraja prve spoljne *for* petlje.

3.4 Levenshtein udaljenost

Definišemo pojam udaljenosti između stringova na drugačiji način, tačnije, dodavanjem nove operacije – izmene slova.

Definicija 3.4 *Levenshtein udaljenost između stringova s, p je minimalan broj operacija potreban da se string s prevede u string p , gde su operacije brisanje slova, umetanje slova i menjanje jednog slova u drugo.*

Ovaj problem rešavamo algoritmom dinamičkog programiranja koji je veoma sličan prethodno opisanom algoritmu za nalaženje LCS-a. Kako je skup operacija na stringovima invertibilan, odnosno, za svaku operaciju postoji inverzna operacija iste težine, svedeno je da li operaciju radimo na stringu s ili na stringu p . Ovo znači da je dovoljno da nađemo string q takav da se minimizuje ukupan broj operacija da se oba stringa dovedu do stringa q . Takođe primetimo da u tom slučaju možemo u potpunosti da zanemarimo operaciju umetanja slova, jer svako umetnuto slovo u , na primer, string s , odgovara nekom slovu u u stringu q . Ako je to slovo bilo umetnuto i u p , onda je ono suvišno i može se eliminisati, pri čemu se smanjuje broj operacija. U suprotnom, ono je ili nastalo izmenom nekog slova u p ili od nekog slova koje je na početku bilo u p . U oba slučaja možemo jednostavno obrisati to slovo u p i ne dodavati ga u s , pri čemu se ne povećava broj operacija. Posmatrajmo dva stringa s, p dužina n, m . Tražimo string q takav da se minimizuje ukupan broj poteza da se i s i p prevedu u q , gde su operacije brisanje slova i izmena slova. Neka je $d_{i,j}$ minimalan broj operacija potreban da se stringovi $s_{[0,i)}$ i $p_{[0,j)}$ dovedu do istog stringa, odnosno, to je njihova Levenshtein udaljenost. Nakon svih operacija njima moraju da se poklapaju poslednja dva slova. Ukoliko se u početku njima poklapaju poslednja dva slova, možemo samo da rešimo problem za $s_{[0,i-1)}$ i $p_{[0,j-1)}$. U suprotnom, možemo da izjednačimo ta dva slova u jednoj operaciji, pri čemu ponovo problem svodimo na stringove $s_{[0,i-1)}$ i $p_{[0,j-1)}$ ili možemo jedno od tih slova obrisati, pri čemu svodimo problem na $s_{[0,i-1)}$ i $p_{[0,j)}$ ako brišemo iz stringa s , ili na $s_{[0,i)}$ i $p_{[0,j-1)}$ ako brišemo iz stringa p . Jasno je da, ako je $i = 0$ ili $j = 0$, onda je Levenshtein udaljenost j odnosno i . Dolazimo do sledeće rekurentne veze:

$$d_{i,j} = \begin{cases} i + j & i = 0 \vee j = 0 \\ \min(d_{i-1,j} + 1, d_{i,j-1} + 1, d_{i-1,j-1} + \delta(s_{i-1} \neq p_{j-1})) & i, j > 0 \end{cases} \quad (3.4.1)$$

Implementacija algoritma za nalaženje Levenshtein udaljenosti

```
1 int levenshtein(const string& s, const string& p) {  
2     int n = s.size(), m = p.size();  
3     vector<vector<int>> d(n+1, vector<int>(m+1));  
4     for (int i=0; i<=n; i++)  
5         for (int j=0; j<=m; j++)  
6             if (i == 0 || j == 0)  
7                 d[i][j] = i+j;  
8             else  
9                 d[i][j] = min({d[i-1][j]+1, d[i][j-1]+1,  
10                    d[i-1][j-1] + (s[i-1] != p[j-1])});  
11     return d[n][m];  
12 }
```

Kao i kod algoritma za nalaženje LCS-a, mogu se primeniti ideje za smanjenje memorijske složenosti, uključujući i Hirschberg-ov algoritam. U svakom slučaju, vremenska složenost je $O(nm)$.

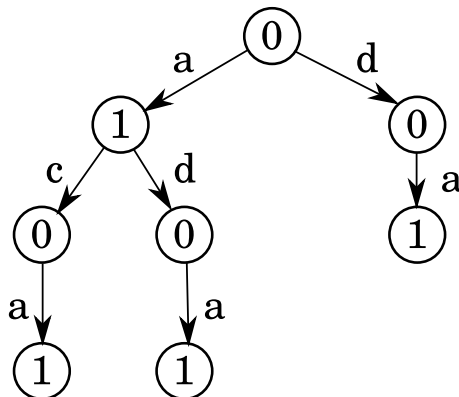
4 Prefiksne strukture podataka

4.1 Prefiksno stablo

Prefiksno stablo, odnosno *trie* je struktura podataka za čuvanje i pretragu kolekcije stringova, i ima oblik n -arnog korenskog stabla, odnosno stabla gde svaki čvor ima stepen najviše $|\Sigma|$. Kod prefiksnog stabla grane imaju labela koje odgovaraju nekom slovu alfabeta Σ , dok čvorovi čuvaju nenegativan ceo broj c_i koji ćemo zvati *višestrukost*.

Definicija 4.1 *Prefiksno stablo za kolekciju stringova $S, s \in \Sigma^*$ je korensko stablo sa minimalnim brojem čvorova, kod kojeg za svaki string $s \in S$ postoji put od korena do čvora koji označavamo sa $\delta(s)$ tako da labela običenih grana obrazuju string s , a u čvoru $\delta(s)$ je višestrukost jednaka broju pojavljivanja stringa s u kolekciji S , i kod kojeg izlazne grane svakog čvora imaju međusobno različite labela.*

Na osnovu definicije sledi da različiti putevi u stablu koji počinju od istog čvora odgovaraju različitim stringovima. Posmatrajmo skup svih prefiksa svih stringova u S , neka je to skup P . Kako postojanje puta od korena sa labelom $s \in S$ implicira postojanje takvog puta za sve prefikse stringa s , važi da će za svaki string $p \in P$ postojati jedinstven čvor u stablu. Pokazuje se da je moguće konstruisati jedinstveno stablo sa tačno $|P|$ čvorova koje zadovoljava definiciju 4.1.



Prefiksno stablo za kolekciju $\{aca, ada, a, da\}$.

4.1.1 Konstrukcija

Kao što je prethodno opisano, svaki čvor u sebi čuva višestrukost, i takođe će čuvati izlazne grane. Grane će biti čuvane kao kolekcija parova (k, v) , gde će k biti slovo, a v pokazivač na čvor do kojeg se dolazi tom granom. Ovo je najjednostavnije uraditi pomoću rečnika, što nam omogućava brzo umetanje nove grane a takođe i brzu pretragu grane sa određenim slovom.

Struktura čvora prefiksnog stabla

```
1 struct trie_node {  
2     map<char, trie_node*> next;  
3     int c;  
4     trie_node() : c(0) {}  
5 };
```

Algoritam za konstrukciju prefiksnog stabla dodaje redom jedan po jedan string iz kolekcije. Pri dodavanju jednog stringa, krećemo se po postojećim granama dokle god je to moguće, a ukoliko nije kreiramo novi čvor. Na kraju se u čvoru u kom se završi kretanje dodaje 1 na vrednost c .

Algoritam za dodavanje jednog stringa u prefiksno stablo

```
1 void trie_insert(trie_node* root, const string& s) {  
2     trie_node* t = root;  
3     for (char x : s)  
4         if (t->next.count(x))  
5             t = t->next[x];  
6         else  
7             t = t->next[x] = new trie_node;  
8     t->c++;  
9 }
```

Ukoliko pretpostavimo da operacije na mapi imaju konstantnu složenost, što je opravdano konstantnom veličinom alfabeta, dolazimo do toga da je vremenska složenost algoritma linearna po veličini stringa koji se dodaje.

4.1.2 Primene

Prefiksno stablo ima veliki broj primena. Svaka od primena može zahtevati drugačiji izgled strukture jednog čvora, ali svima je zajedničko to da čvor

čuva svoje izlazne grane i skelet algoritma za dodavanje stringa je isti kao što je prikazano u kodu gore.

Najosnovnija primena prefiksnog stabla je za implementaciju rečnika, odnosno, ono nam omogućava da brzo proverimo da li se neki string javlja u kolekciji ili ne.

Algoritam za brojanje pojavljivanja stringa u prefiksnom stablu

```
1 int trie_count(trie_node* root, const string& s) {
2     trie_node* t = root;
3     for (char x : s)
4         if (t->next.count(x))
5             t = t->next[x];
6     else
7         return 0;
8     return t->c;
9 }
```

Ukoliko pri dodavanju jednog stringa svim usput obišenim čvorovima inkrementiramo višestrukost c , prethodni algoritam će vratiti broj stringova u kolekciji koji imaju s kao svoj prefiks, odnosno, moguće je dodati sve prefikse jednog stringa odjednom, bez povećanja vremenske ili memorijske složenosti.

Moguće je i pronaći leksikografski najmanji string u kolekciji koji je veći ili jednak zadatom, odnosno, naći *lower bound* za dati string. Za to nam je potrebna pomoćna funkcija koja za dati čvor nalazi leksikografski najmanji put od tog čvora do nekog čvora sa pozitivnom višestrukošću.

Algoritam za nalaženje najmanjeg stringa počev od zadanog čvora

```
1 string trie_smallest(trie_node* t) {
2     string s;
3     while (!t->c) {
4         s += t->next.begin()->first;
5         t = t->next.begin()->second;
6     }
7     return s;
8 }
```

Primetimo da $c = 0$ implicira da čvor ima bar jednu izlaznu granu.

Algoritam za nalaženje lower bound-a za zadati string

```
1 string trie_lb(trie_node* root, const string& s) {
2     int n = s.size();
3     trie_node* t = root;
4     vector<trie_node*> path = {t};
5     for (int i=0; i<n; i++) {
6         char x = s[i];
7         if (t->next.count(x)) {
8             t = t->next[x];
9             path.push_back(t);
10        } else {
11            for (int j=i; j>=0; j--) {
12                if (auto it = path[j]->next.upper_bound(s[j]);
13                    it != path[j]->next.end())
14                {
15                    return s.substr(0, j) + it->first +
16                        trie_smallest(it->second);
17                }
18            }
19            return "";
20        }
21    }
22    return s + trie_smallest(t);
23 }
```

Algoritam je grabljive prirode. Posmatrajmo najduži zajednički prefiks stringa s za koji tražimo *lower bound* i nekog stringa iz kolekcije. Ukoliko je taj prefiks jednak celom stringu s , potrebno je da nađemo najmanji string q dostižan iz čvora $\delta(s)$, i da vratimo sq . U suprotnom, ako je dužina tog zajedničkog prefiksa i , biramo najveću poziciju $j \leq i$ takvu da u čvoru $\delta(s_{[0,j]})$ postoji izlazna grana čija je labela slovo strogo veće od s_j , zatim vraćamo string $s_{[0,j]}yq$, gde je y to slovo a q najmanji string dostižan iz čvora $\delta(s_{[0,j]}y)$. Ukoliko takav indeks ne postoji, ne postoji ni leksikografski veći ili jednak string, a algoritam vraća prazan string. Ukupna memorijska i vremenska složenost celog algoritma je $O(n + m)$. gde je $n = |s|$, a m je dužina rezultujućeg stringa.

4.2 Aho-Corasick algoritam

Aho-Corasick algoritam omogućava pretragu više stringova odjednom, odnosno, za dati skup nepraznih stringova P i string s , on pronalazi sva pojavljivanja

svih stringova iz P u s . Kako je ukupan broj pojavljivanja svih stringova u najgorem slučaju srazmeran $|s| \cdot |P|$, tako je i ovo donja granica na složenost bilo kog algoritma koji pronalazi sva pojavljivanja.

Aho-Corasick algoritam generalizuje KMP algoritam. Osnova za ovaj algoritam je upravo prefiksno stablo i prvi korak jeste njegova konstrukcija, uz malu modifikaciju, koja se odnosi na postojanje dodatnih polja, koje se zovu *sufiks veza* i *rečnička veza*. Takođe, umesto višestrukosti, za čvor $\delta(p)$ pamtimo *id*, odnosno redni broj stringa $p \in P$, ukoliko se ne radi o stringu iz P već o nekom prefiksu, upisujemo $id = -1$.

Struktura čvora kod Aho-Corasick algoritma

```

1 struct aho_node {
2     map<char, aho_node*> next;
3     int id;
4     aho_node* link;
5     aho_node* dict;
6     aho_node() : id(-1), link(nullptr), dict(nullptr) {}
7 };

```

Dodavanje stringa u prefiksno stablo kod Aho-Corasick algoritma

```

1 void aho_insert(aho_node* root, const string& s, int id) {
2     aho_node* t = root;
3     for (char x : s)
4         if (t->next.count(x))
5             t = t->next[x];
6         else
7             t = t->next[x] = new aho_node;
8     t->id = id;
9 }

```

Definicija 4.2 Za prefiksno stablo izgrađeno od skupa stringova P , sufiks veza čvora $\delta(s)$ pokazuje na čvor $\delta(q)$ takav da je q najduži string koji je prefiks nekog stringa u P , a koji je ujedno pravi sufiks stringa s . Za koren se sufiks veza ne definiše.

Primetimo da za jednoelementni skup P ove sufiks veze tačno odgovaraju nizu neuspeha kod KMP algoritma. Sufiks veza nekog čvora uvek pokazuje na čvor strogo manje dubine. Zbog ovoga, čvorove je neophodno obraditi u redosledu sortiranom po dubini, za šta se koristi pretraga u širinu. Prirodno,

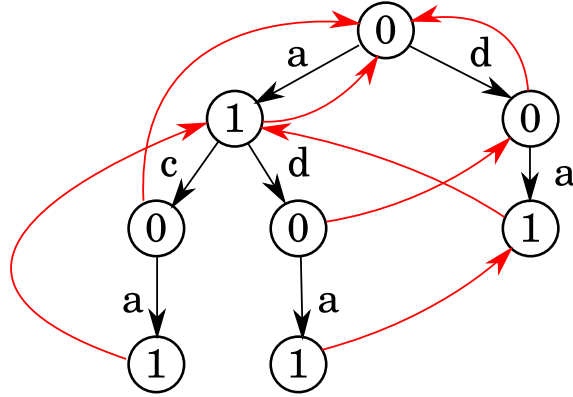
algoritam za nalaženje sufiks veza podseća na prvu fazu KMP algoritma. Za svaki čvor p , krećemo od sufiks veze njegovog roditelja t , i tražimo čvor l koji ima izlaznu granu sa istom labelom kao grana koja spaja t, p . Pretragu vršimo praćenjem sufiks veza. Ukoliko nađemo takav čvor, sufiks veza čvora p će pokazivati na odgovarajuće dete čvora l , u suprotnom, sufiks veza pokazuje na koren.

Nalaženje sufiks veza kod Aho-Corasick algoritma

```

1  vector<aho_node*> aho_find_links(aho_node* root) {
2      vector<aho_node*> q = {root};
3      size_t qs = 0;
4      while (qs != q.size()) {
5          aho_node* t = q[qs++];
6          for (auto [c, p] : t->next) {
7              aho_node* l = t->link;
8              while (l && !l->next.count(c))
9                  l = l->link;
10             if (l)
11                 p->link = l->next[c];
12             else
13                 p->link = root;
14             q.push_back(p);
15         }
16     }
17     return q;
18 }
```

Funkcija vraća niz svih čvorova stabla, sortiran po dubini. Ovaj niz će biti koristan u kasnijoj obradi. Neka je L ukupna dužina svih stringova u P . Vremenska složenost nalaženja svih sufiks veza je $O(L)$.⁵ Dokaz je sličan dokazu složenosti kod KMP algoritma.



Sufiks veze za prefiksno stablo za kolekciju {aca, ada, a, da}.

Prefiksno stablo podseća na deterministički konačni automat, jer ima jedan polazni čvor i svaki čvor ima izlazne grane označene simbolima iz alfabeta Σ . Međutim, neki čvorovi, npr. listovi stabla nemaju izlazne grane za svaki simbol $x \in \Sigma$. Pomoću izračunatih sufiks veza možemo dopuniti stablo do pravog konačnog automata. Neka je u proizvoljan čvor stabla, a $x \in \Sigma$ proizvoljan simbol. Počev od čvora u , uključujući i u , krećemo se po sufiks vezama i tražimo prvi čvor koji ima izlaznu granu sa labelom x . Ukoliko takav čvor ne postoji, grana automata sa labelom x iz u ide ka korenu stabla. U suprotnom, neka je to čvor v . Tada ta grana ide ka čvoru na koji pokazuje grana sa labelom x iz čvora v .

Smisao ovakve definicije automata je sledeći. Ukoliko se kroz automat propusti string s , ako se automat nalazi u stanju $\delta(p)$, tada je p najduži prefiks nekog stringa iz P koji je ujedno sufiks stringa s . Tada je ovakva definicija automata valjana, odnosno, ovako definisan automat zaista raspoznaje sve različite prefikse stringova iz P . Zbog ovakve definicije, nije nužno eksplicitno pamtititi sve prelaze automata – dovoljno je pri svakom prelazu potražiti pomoću sufiks veza odgovarajući čvor. Slično kao kod KMP-a, ukoliko se kroz ovaj automat propusti string s , ukupna vremenska složenost za kretanje kroz automat će biti $O(|s|)$.⁵

Pošto svaka sufiks veza ide od čvora veće ka čvoru manje dubine, graf koji obrazuju ove veze je acikličan, a pošto svaki čvor osim tačno jednog ima izlaznu granu, radi se o obrnutom korenskom stablu. Nadalje ćemo ovo stablo zvati *stablo sufiks veza*. I originalni Aho-Corasick algoritam i njegove modifikacije se oslanjaju na preprocesiranje ovog stabla.

Originalni Aho-Corasick algoritam nalazi, za skup stringova P i svaki prefiks $s_{[0,i]}$ stringa s sve stringove $p \in P$ koji su sufixi stringa $s_{[0,i]}$, u vremenskoj složenosti $O(L + |s| + k)$, gde je k broj poklapanja. Posmatrajmo šta se dešava nakon tačno i prelaza, odnosno, kada smo stigli do prefiksa $s_{[0,i]}$. Neka se automat nalazi u stanju $\delta(q)$. Posmatrajmo put od čvora q do korena kroz stablo sufiks veza. Svaki čvor na ovom putu koji ima pozitivnu višestrukost odgovara pojavljivanju nekog stringa $p \in P$ kao sufix stringa $s_{[0,i]}$. Dakle, naš cilj treba da bude da efikasno otkrijemo sve ove čvorove, poželjno u linearnoj složenosti po njihovom broju. Zato uvodimo novi tip veze – rečničku vezu.

Definicija 4.3 *Rečnička veza za čvor u stabla sufiks veza je prvi čvor dostižan iz u kod kojeg je višestrukost pozitivna. Ukoliko takav čvor ne postoji, rečnička veza se ne definiše.*

Rečničke veze nam omogućavaju da preskočimo sve usputne čvorove koji ne odgovaraju celim stringovima iz P , odnosno, da direktno obiđemo sve čvorove sa pozitivnom višestrukošću. Ove veze se jednostavno konstruišu ukoliko je već izračunat niz čvorova sortiran po dubini.

Nalaženje rečničkih veza kod Aho-Corasick algoritma

```

1 void aho_find_dict(const vector<aho_node*>& q) {
2     int L = q.size();
3     q[0]->dict = nullptr;
4     for (int i=1; i<L; i++)
5         if (q[i]->id != -1)
6             q[i]->dict = q[i];
7         else
8             q[i]->dict = q[i]->link->dict;
9 }

```

Niz q je prethodno izračunati niz čvorova, jasno je da je q_0 koren i da je on jedini čvor koji nema sufiks vezu. Za sve ostale čvorove važi da, ako oni imaju pozitivnu višestrukost, onda rečnička veza pokazuje upravo na njih same, a inače će pokazivati na isti čvor na koji pokazuje njihova sufiks veza.

Konačno, opišimo rad celog algoritma. Prvo, konstruišemo prefiksno stablo ubacivanjem svih stringova iz P . Zatim nalazimo sufiks veze i rečničke veze. Zatim, obrađujemo string s , u svakom trenutku pamtimo pokazivač na trenutnu poziciju u automatu. Nakon svakog dodatog slova, pomoću rečničkih veza obiđemo sva sufiksna poklapanja sa stringovima iz P .

Glavna funkcija Aho-Corasick algoritma

```
1 vector<pair<int, int>> aho_main(  
2     const vector<string>& P,  
3     const string& s  
4 ) {  
5     aho_node* root = new aho_node;  
6     for (int i=0; i<(int)P.size(); i++)  
7         aho_insert(root, P[i], i);  
8     vector<aho_node*> q = aho_find_links(root);  
9     aho_find_dict(q);  
10    aho_node* curr = root;  
11    vector<pair<int, int>> result;  
12    for (int i=1; i<=(int)s.size(); i++) {  
13        while (curr && !curr->next.count(s[i-1]))  
14            curr = curr->link;  
15        curr = curr ? curr->next[s[i-1]] : root;  
16        for (auto l = curr->dict; l; l = l->link->dict)  
17            result.emplace_back(l->id, i-(int)P[l->id].size());  
18    }  
19    for (auto l : q)  
20        delete l;  
21    return result;  
22 }
```

Na osnovu svega do sad opisanog, složenost algoritma je linearna po zbiru veličina ulaza i izlaza. Jedan od nedostataka ovakvog pristupa je što izlaz može biti dosta veliki, na primer, ako je $P = \{a, a^2, \dots, a^k\}$ a $s = a^n$, ukupna veličina ulaza je $\Theta(k^2 + n)$ dok je veličina izlaza $\Theta(kn)$, što je za, na primer $k = \Theta(\sqrt{n})$ superlinearna funkcija dužine ulaza. Ukoliko je potrebno eksplicitno enumerisati sva pojavljivanja prethodno opisani algoritam je optimalan. U suprotnom, moguće varijacije su da treba da se prijavi broj pojavljivanja ili prvo pojavljivanje svakog stringa $p \in P$ u s . Ovo je moguće postići u linearnom vremenu po veličini ulaza modifikacijom glavnog algoritma. Naime, za svaki čvor prefiksnog stabla zapamtimo sve trenutke i kada smo, nakon obrade prefiksa $s_{[0,i]}$ završili baš u tom čvoru. Broj pojavljivanja nekog stringa $p \in P$ je onda ukupan broj trenutaka kada smo bili u podstablu stabla sufiks veza sa korenem u čvoru $\delta(p)$, dok je prvo pojavljivanje jednako minimumu svih trenutaka pojavljivanja u tom podstablu. Obe ove vrednosti se mogu efikasno izračunati postprocesiranjem stabla nakon obilaska celog stringa s , primenom dinamičkog programiranja.

5 Sufiksne strukture podataka

5.1 Sufiks niz

Sufiks niz je struktura podataka koja omogućava brzo traženje pojavljivanja stringa unutar stringa za koji se konstruiše sufiks niz, tačnije, vremenska složenost pretrage je sublinearna funkcija dužine stringa unutar kojeg se vrši pretraga. Pored ovoga, pomoću sufiks niza se mogu brzo vršiti poređenja podstringova unutar samog stringa.⁷

Sufiks niz za string s je niz sortiranih nepraznih sufiksa tog stringa. Formalno,

Definicija 5.1 *Sufiks niz za string s dužine $|s| = n$ je niz p koji se sastoji od n različitih celih brojeva iz skupa $\{0, \dots, n-1\}$ takav da je niz sufiksa čije su početne pozicije p_0, p_1, \dots, p_{n-1} leksikografski rastući niz.*

Za svaki string postoji jedinstven sufiks niz, zato što je leksikografsko uređenje totalno a ne postoje dva jednaka sufiksa. Primera radi, nađimo sufiks niz za string **banana**. Označimo sa u_i string $s_{[i,n]}$. Svi sufiksi ovog stringa su:

u_0	banana
u_1	anana
u_2	nana
u_3	ana
u_4	na
u_5	a

Sortiranjem dobijamo niz sufiksa:

u_5	a
u_3	ana
u_1	anana
u_0	banana
u_4	na
u_2	nana

Sortirani niz sufiksa je $u_5, u_3, u_1, u_0, u_4, u_2$, pa je sufiks niz $p = (5, 3, 1, 0, 4, 2)$.

5.1.1 Konstrukcija

Sufiks niz se može konstruisati prostim sortiranjem svih sufiksa u vremenskoj složenosti $O(n^2 \log n)$, ukoliko je dostupan algoritam za sortiranje opšte namene koji radi u složenosti $O(n \log n)$. Primer takvog algoritam je *mergesort*. Radi uštede memorijskog prostora sortiraćemo samo niz celih brojeva $0, 1, \dots, n-1$, dok ćemo kao funkciju za poređenje koristiti funkciju koja je "svesna" stringa s i koja za data dva sufiksa određuje koji je manji.

Implementacija klase za poređenje sufiksa stringa

```
1 struct suffix_cmp {
2     const string& s;
3     suffix_cmp(const string& s) : s(s) {}
4     bool operator() (int u, int v) const {
5         if (u == v)
6             return false;
7         return lexicographical_compare(
8             s.begin()+u, s.end(),
9             s.begin()+v, s.end());
10    }
11};
```

Vremenska složenost poređenja dva sufiksa je $O(n)$, a kako *mergesort* sortira niz sa $O(n \log n)$ poziva funkcije za poređenje, ukupna vremenska složenost je $O(n^2 \log n)$.

Glavni algoritam za nalaženje sufiks niza, složenosti $O(n^2 \log n)$

```
1 vector<int> sarray_slow(const string& s) {
2     int n = s.size();
3     vector<int> p(n);
4     iota(p.begin(), p.end(), 0);
5     sort(p.begin(), p.end(), suffix_cmp(s));
6     return p;
7 }
```

Napomena. Funkcija *iota* puni zadati opseg uzastopnim vrednostima počev od trećeg parametra, u ovom slučaju 0. Koristimo je da niz p inicijalizujemo vrednostima $p_i = i$. Prema standardu jezika, počev od C++11, funkcija *sort* zove funkciju za poređenje $O(n \log n)$ puta.¹⁶ Treći parametar je funkcija ili drugi objekat koji ima implementiran *operator()* koji može

da poredi elemente niza.

Označimo sa k dužinu najdužeg stringa koji se javlja više od jednom u stringu s . Nije teško pokazati da, ako se pažljivo implementira, funkcija poređenja dva sufiksa radi u vremenskoj složenosti $O(k)$, pa se složenost konstrukcije sufiks niza može i bolje proceniti sa $O(kn \log n)$.

5.1.2 Brži algoritmi za konstrukciju

Algoritmi dati u nastavku efikasno rešavaju problem nalaženja sortiranog niza svih cikličnih pomeraja stringa. Prvo ćemo pokazati kako se problem nalaženja sufiks niza svodi na sortiranje svih cikličnih pomeraja stringa. Neka je $s \in \Sigma^+$. Proširimo alfabet Σ dodavanjem novog simbola, koji ćemo označiti sa $\$$ koji je po uređenju manji od svih simbola iz Σ . Tada je $s' = s\$ \in (\Sigma \cup \{\$\})^+$. Neka je $n = |s|, n' = n + 1 = |s'|$. Posmatrajmo sve ciklične pomeraje stringa s' . Pošto se karakter $\$$ javlja samo jednom u stringu s' , svi ovi stringovi su različiti pa postoji jedinstven leksikografski poredak. Leksikografski najmanji pomeraj biće $n' - 1$, jer je taj pomeraj jedini koji počinje karakterom $\$$ koji je po definiciji manji od svih ostalih.

Teorema 5.1 *Ako su i, j proizvoljni ciklični pomeraji stringa s' različiti od $n' - 1$, tada važi $s'_{[i, i+n')} < s'_{[j, j+n')}$ akko je $s_{[i, n)} < s_{[j, n)}$. \square*

Odavde sledi da ako je niz q_0, q_1, \dots, q_n sortiran niz cikličnih pomeraja stringa s' tada je niz $p_i = q_{i+1}, i \in \{0, \dots, n-1\}$ sufiks niz stringa s .

Svođenje konstrukcije sufiks niza na sortiranje cikličnih pomeraja

```

1  template<class T>
2  vector<int> sarray_scs(const string& s, T func) {
3      vector<int> v = func(s + '\0');
4      v.erase(v.begin());
5      return v;
6  }
```

Konačno, opišimo algoritam za sortiranje cikličnih pomeraja stringa s dužine n . Za svako $k \in \{0, 1, \dots, \lceil \log_2(n) \rceil\}$ nađimo poredak svih cikličnih podstringova dužine 2^k . Ovaj poredak ćemo opisati nizom celih brojeva $C^{(k)}$, gde stringu $s_{[i, i+2^k)}$ pridružujemo broj $C_i^{(k)}$ na takav način da, ako su podstringovima $s_{[i, i+2^k)}$ i $s_{[j, j+2^k)}$ pridruženi isti brojevi, tada su oni jednaki, a ako

je jednom pridružen manji broj, tada je taj podstring leksikografski manji. Za $k = 0$, možemo jednostavno postaviti $C_i^{(0)} = \text{Ord}(s_i)$.

Neka je $k > 0$. Naš cilj je da odredimo niz $C^{(k)}$ u vremenskoj složenosti $O(n \log n)$. Posmatrajmo podstringove $s_{[i, i+2^k)}$ i $s_{[j, j+2^k)}$ i posmatrajmo uređene parove $u_i = (C_i^{(k-1)}, C_{(i+2^{k-1}) \bmod n}^{(k-1)})$ i $u_j = (C_j^{(k-1)}, C_{(j+2^{k-1}) \bmod n}^{(k-1)})$. Tada je poredak ovih podstringova jednak poretку parova u_i, u_j . Ako za svaki ciklični podstring odredimo ovaj uređeni par, zatim sve dobijene parove sortiramo, a zatim ovim parovima dodelimo ordinalne vrednosti, dobićemo upravo niz $C^{(k)}$. Ovo se može jednostavno izvesti bilo kojim algoritmom za sortiranje koji radi u složenosti $O(n \log n)$.

Za kraj, neka je $k = \lceil \log_2(n) \rceil$. Odavde je $2^k \geq n$. U kontekstu konstrukcije sufiks niza, nijedna dva ciklična pomeraja dužine 2^k neće biti jednaka, pa će niz $C^{(k)}$ sadržati različite brojeve. Odavde sortiran niz cikličnih pomeraja možemo naći kao inverznu permutaciju niza $C^{(k)}$. Važno je napomenuti da će poredak podstringova dužine 2^k biti jednak poretку podstringova dužine n , zato što dužina najdužeg zajedničkog prefiksa za bilo koja dva različita podstringa (bilo koje dužine) iznosi $n - 1$, upravo jer postoji karakter $\$$ koji se javlja samo jednom u stringu.

Kako ovaj algoritam ima $O(\log n)$ faza, a svaka faza radi u složenosti $O(n \log n)$, ukupna vremenska složenost je $O(n \log^2 n)$. Memorijska složenost je $O(n \log n)$, ali se može smanjiti na $O(n)$ pamćenjem samo nizova $C^{(k)}$ prethodne i trenutne faze.

Algoritam za sortiranje cikličnih pomeraja, složenosti $O(n \log^2(n))$

```

1  vector<int> scs_fast(const string& s) {
2      int n = s.size();
3      vector<int> c(s.begin(), s.end());
4      for (int h=1; h<n; h*=2) {
5          vector<array<int, 3>> t(n);
6          for (int i=0; i<n; i++)
7              t[i] = {c[i], c[(i+h)%n], i};
8          sort(t.begin(), t.end());
9          vector<int> cnew(n);
10         cnew[t[0][2]] = 0;
11         int numc = 1;
12         for (int i=1; i<n; i++)
13             if (t[i][0] == t[i-1][0] && t[i][1] == t[i-1][1])
14                 cnew[t[i][2]] = numc-1;
15             else
16                 cnew[t[i][2]] = numc++;
17         swap(c, cnew);
18     }
19     vector<pair<int, int>> g(n);
20     for (int i=0; i<n; i++)
21         g[i] = {c[i], i};
22     sort(g.begin(), g.end());
23     vector<int> p(n);
24     for (int i=0; i<n; i++)
25         p[i] = g[i].second;
26     return p;
27 }

```

Prethodno opisani algoritam za nalaženje cikličnih pomeraja je moguće modifikovati tako da radi u složenosti $O(n \log n)$. Naime, ukoliko bismo sortirali parove u_i u linearnom vremenu, dobili bismo upravo tu vremensku složenost. Svaki par se sastoji iz dva broja iz skupa $\{0, 1, \dots, n-1\}$. Možemo primeniti ideju iz algoritma *radix sort*. Naime, *radix sort* se oslanja na *counting sort*, koji ima jednostavnu implementaciju i može da sortira bilo koji niz od n elemenata čiji su ključevi za poređenje brojevi iz skupa $\{0, 1, \dots, k-1\}$ u vremenskoj složenosti $O(n+k)$. Pritom, moguće je implementirati *counting sort* kao stabilan algoritam sortiranja, odnosno algoritam koji ekvivalentnim elementima ne menja relativni poredak. Algoritam *radix sort* za uređene parove bi prvo pomoću *counting sort*-a sortirao sve parove po drugom elementu, a zatim po prvim, vodeći računa da se ne naruši prethodno ustanovljen poredak korišćenjem stabilne varijante *counting sort*-a.

Moguće je pojednostaviti prethodni algoritam, odnosno, svesti ga na samo jedno pozivanje *counting sort*-a. Naime, pored nizova $C^{(k)}$ čuvaćemo eksplicitno i permutaciju $p^{(k)}$ koja odgovara poretku podstringova dužine 2^k . Posmatrajmo šta se dešava ukoliko pomoću *counting sort*-a sortiramo vrednosti $(j - 2^{k-1}) \bmod n$, gde je ključ $C_{(j-2^{k-1}) \bmod n}^{(k-1)}$, uzete redom za svako j iz niza $p^{(k-1)}$. Stabilan *counting sort* će urediti ove indekse $j' = (j - 2^{k-1}) \bmod n$ po vrednosti $C_{j'}^{(k-1)}$, dok će, ukoliko više njih ima istu vrednost, očuvati prethodno ustanovljeni poredak. Kako je ovaj prethodni poredak indukovao vrednostima $C_j^{(k-1)} = C_{(j'+2^{k-1}) \bmod n}^{(k-1)}$, dobijamo upravo leksikografski redosled parova $(C_{j'}^{(k-1)}, C_{(j'+2^{k-1}) \bmod n}^{(k-1)})$. Sada na osnovu ovog sortirano niz računamo nove vrednosti $C^{(k)}$, dok je $p^{(k)}$ dobijen upravo pomenutim sortiranjem.

Konačno rešenje, odnosno sortirani niz cikličnih pomeraja je upravo niz $p^{(k)}$ za $k = \lceil \log_2(n) \rceil$.

Algoritam za sortiranje cikličnih pomeraja, složenosti $O(n \log n)$

```

1  vector<int> scs_faster(const string& s) {
2      int n = s.size(), k = 256, sz = 0;
3      vector<int> p(n), c(s.begin(), s.end());
4      vector<vector<int>> g(max(n, k));
5      for (int i=0; i<n; i++)
6          g[c[i]].push_back(i);
7      for (auto& gr : g) {
8          for (int i : gr)
9              p[sz++] = i;
10         gr.clear();
11     }
12     for (int h=1; h<n; h*=2) {
13         vector<int> pnew(n), cnew(n);
14         for (int j : p) {
15             int jp = (j+n-h)%n;
16             g[c[jp]].push_back(jp);
17         }
18         sz = 0;
19         for (auto& gr : g) {
20             for (int i : gr)
21                 pnew[sz++] = i;
22             gr.clear();
23         }
24         cnew[pnew[0]] = 0;
25         int numc = 1;
26         for (int i=1; i<n; i++) {
27             int s0 = pnew[i-1], s1 = pnew[i];
28             if (c[s1] == c[s0] && c[(s1+h)%n] == c[(s0+h)%n])
29                 cnew[s1] = numc-1;
30             else
31                 cnew[s1] = numc++;
32         }
33         swap(c, cnew);
34         swap(p, pnew);
35     }
36     return p;
37 }

```

5.1.3 LCP niz

Niz najdužih zajedničkih prefiksa (*longest common prefix*) je niz koji se često zajedno koristi sa sufiks nizom. Za string s dužine n , čiji je sufiks niz p_0, p_1, \dots, p_{n-1} , LCP niz se sastoji od nenegativnih celih brojeva q_0, q_1, \dots, q_{n-2} ,

gde q_i označava dužinu najdužeg zajedničkog prefiksa stringova $s_{[p_i, n)}$ i $s_{[p_{i+1}, n)}$.

Najduži zajednički prefiks zadovoljava jednu veoma važnu osobinu.

Teorema 5.2 *Neka je s_1, s_2, \dots, s_n leksikografski sortiran niz stringova. Neka je $q_i = LCP(s_i, s_{i+1})$. Ako je $i < j$, tada je $LCP(s_i, s_j) = \min\{q_i, q_{i+1}, \dots, q_{j-1}\}$. \square*

Kako se LCP niz konstruiše nad sortiranim nizom sufiksa jednog stringa, on se uz odgovarajuću strukturu podataka za nalaženje minimuma u podnizu može koristiti za određivanje najdužeg zajedničkog prefiksa bilo koja dva sufiksa.

Pre nego što opišemo algoritam za konstrukciju LCP niza, dokažimo sledeću teoremu.

Teorema 5.3 *Neka je s string dužine n , čiji je sufiks niz p_0, \dots, p_{n-1} , LCP niz q_0, \dots, q_{n-2} . Definišemo $r = p^{-1}$ odnosno, $p_{r_i} = i$. Neka je $s_{[i, n)}$ njegov sufiks takav da je $i > 0$ i $p_{n-1} \neq i, i-1$. Tada je $q_{r_i} \geq q_{r_{i-1}} - 1$.*

Dokaz: Uvedimo oznake $i' = i - 1$, $j = p_{r_i+1}$, $j' = p_{r_{i'}+1}$. Drugim rečima, sufiks i' je onaj koji prethodi i , odnosno ima dužinu za 1 veću. j je sufiks koji se nalazi odmah posle i u sufiks nizu. Slično, j' je sufiks koji se nalazi odmah posle i' u sufiks nizu. Ukoliko je $q_{r_{i'}} \leq 1$, onda tvrđenje očigledno važi, jer je $q_{r_i} \geq 0$. Neka je $q_{r_{i'}} \geq 2$. Pošto je j' posle i' u sufiks nizu važi $s_{[j', n)} > s_{[i', n)}$ a pošto je $q_{r_{i'}} \geq 1$ važi $s_{i'} = s_{j'}$. Ako odbacimo prvi karakter sufiksa i' i j' ponovo dobijamo sufikse, $i' + 1, j' + 1$ i važi $s_{[j'+1, n)} > s_{[i'+1, n)}$, odnosno $s_{[j'+1, n)} > s_{[i, n)}$ ili ekvivalentno $r_{j'+1} > r_i$. Dužina najdužeg zajedničkog prefiksa za $i' + 1, j' + 1$ je $q_{r_{i'}} - 1$. Na osnovu teoreme 5.2 imamo da je $q_{r_{i'}} - 1 = LCP(s_{[i, n)}, s_{[j'+1, n)}) = \min\{q_{r_i}, q_{r_i+1}, \dots, q_{r_{j'+1}-1}\}$ odnosno $q_{r_{i'}} - 1 \leq q_{r_i}$ \square .

Algoritam za konstrukciju LCP niza⁸ radi na sledeći način. Kao ulaz se prosleđuju string i izračunati sufiks niz tog stringa. Prvo se računa inverz sufiksnog niza. Zatim se sufiksi obrađuju redom, opadajuće po dužini, odnosno, uzimaju se redom sufiksi čije su početne pozicije $i = 0, 1, \dots, n-1$ tim redom. U svakom trenutku održavamo promenljivu k koja je manja ili jednaka od trenutne vrednosti q_{r_i} koju tražimo. U početku je $k = 0$. Ukoliko je $r_i = n-1$, q_{r_i} se i ne definiše i samo postavljamo $k = 0$. U suprotnom, samo povećavamo k za po 1 sve dok se odgovarajući karakteri sufiksa i i $j = p_{r_i+1}$ poklapaju. Kada dođemo do kraja nekog od ovih sufiksa ili se karakteri ne poklope, prekidamo, upisujemo $q_{r_i} := k$ i zatim, na osnovu

teoreme 5.3 smemo da postavimo $k := \max\{k - 1, 0\}$.

Algoritam za nalaženje LCP niza

```
1 vector<int> lcp_array(const string& s, const vector<int>& p)
2 {
3     int n = s.size(), k = 0;
4     vector<int> q(n-1), r(n);
5     for (int i=0; i<n; i++)
6         r[p[i]] = i;
7     for (int i=0; i<n; i++) {
8         if (r[i] != n-1) {
9             int j = p[r[i] + 1];
10            while (i+k < n && j+k < n && s[i+k] == s[j+k])
11                k++;
12            q[r[i]] = k;
13            k = max(0, k-1);
14        } else {
15            k = 0;
16        }
17    }
18    return q;
19 }
```

Vremenska složenost ovog algoritma je $O(n)$ za računanje inverza i petlju koja redom obrađuje sufikse, plus linearna po ukupnom broju uvećavanja promenljive k . Kako u tačno jednoj iteraciji (kad je $r_i = n - 1$) k smanjujemo direktno na 0, a u svim ostalim iteracijama k smanjujemo za najviše 1, ukupno je smanjujemo za najviše $2n - 2$. Kako je početna vrednost 0 a krajnja ne više od $n - 1$, to znači da je k povećavamo najviše $3n - 3$, odnosno $O(n)$ puta, pa ceo algoritam radi u složenosti $O(n)$.

5.1.4 Primene

Traženje jednog stringa u drugom

Ukoliko je potrebno tražiti veliki broj stringova, recimo njih m unutar jednog istog stringa s dužine n , pritom, za svaki od tih stringova je neophodno odmah naći odgovor pre nego što se počne sa obradom sledećeg, bez preprocesiranja stringa s biće nam potrebno bar $O(nm)$ vremena. Ukoliko su stringovi koji se traže mnogo manje dužine od n , recimo, svaki je dužine d , onda bi algoritam složenosti $O(md \log n)$ radio mnogo brže.

Pomoću sufiks niza realizacija algoritma koji traži string t dužine d u stringu s je krajnje jednostavna. Binarnom pretragom po sufiks nizu tražimo prvu poziciju takvu da je odgovarajući sufiks leksikografski veći ili jednak traženom stringu. Kako se poređenje bilo kog sufiksa i stringa p može izvršiti u vremenu $O(d)$, složenost pretrage je $O(d \log n)$.

Sledeća implementacija nalazi *lower bound* za string t , odnosno, najmanji broj i takav da je $s_{[p_i, n)} \leq t$, ili $i = n$ ukoliko takav broj ne postoji.

Traženje lower bound-a pomoću sufiks niza

```

1  int sarray_lb(
2      const string& s,
3      const vector<int>& p,
4      const string& t
5  ) {
6      int n = s.size(), l = 0, r = n-1, i = n;
7      while (l <= r) {
8          int mid = (l+r) >> 1, j = p[mid];
9          if (lexicographical_compare(
10             s.begin()+j, s.end(),
11             t.begin(), t.end()))
12             {
13                 l = mid + 1;
14             } else {
15                 i = mid;
16                 r = mid - 1;
17             }
18     }
19     return i;
20 }
```

Prethodni algoritam se može iskoristiti za traženje pojavljivanja stringa t u stringu s , tako što za nađeni broj i uporedimo stringove $s_{[p_i, p_i+d)}$ i t , naravno, ukoliko je $p_i + d \leq n$.

Traženje podstringa pomoću sufiks niza

```
1 int sarray_find(  
2     const string& s,  
3     const vector<int>& p,  
4     const string& t  
5 ) {  
6     int i = sarray_lb(s, p, t), n = s.size(), d = t.size();  
7     if (i == n)  
8         return -1;  
9     int j = p[i];  
10    if (j+d <= n && equal(t.begin(), t.end(), s.begin()+j))  
11        return j;  
12    return -1;  
13 }
```

Leksikografski najmanji ciklični pomerač

Pomoću algoritma za sortiranje cikličnih pomerača možemo lako utvrditi koji ciklični pomerač je najmanji, tako što direktno primenimo tu funkciju na string bez dodavanja novog, minimalnog karaktera. Štaviše, jednostavno možemo naći i minimalnu periodu stringa s dužine n , odnosno, najmanji prirodan broj l takav da je $s_i = s_{(i+l) \bmod n}$ za svako i . To je upravo broj različitih elemenata niza $C^{(k)}$ za $k = \lceil \log_2(n) \rceil$.

Leksikografsko poređenje podstringova

Ukoliko kod algoritma za sortiranje cikličnih pomerača upamtimo sve nizove $C^{(k)}$, možemo vrlo jednostavno vršiti leksikografsko poređenje podstringova, uključujući i ciklične podstringove. Ukoliko treba da uporedimo dva stringa različitih dužina, recimo $l_1 < l_2$, prvo uporedimo kraći od njih sa prefiksom dužeg dužine l_1 . Ukoliko dobijemo da su ti stringovi različiti, prekidamo, inače, po definiciji, leksikografski manji string je kraći. Zato pretpostavimo da upoređujemo podstringove iste dužine, neka su to $s_{[i, i+l)}$, $s_{[j, j+l)}$. Ukoliko je $l = 2^k$ za neko k , možemo prosto uporediti vrednosti $C_i^{(k)}$ i $C_j^{(k)}$. U suprotnom, nađimo najveće k takvo da je $2^k < l$. Pritom, jasno je da važi $2^k > \frac{l}{2}$. Ideja je da prvo uporedimo prefikse ovih stringova dužine 2^k isto kao u prethodnom slučaju. Ukoliko dobijemo da su prefiksi jednaki, uporedimo njihove sufikse dužine 2^k – rezultat ovog poređenja biće rezultat poređenja celih stringova.

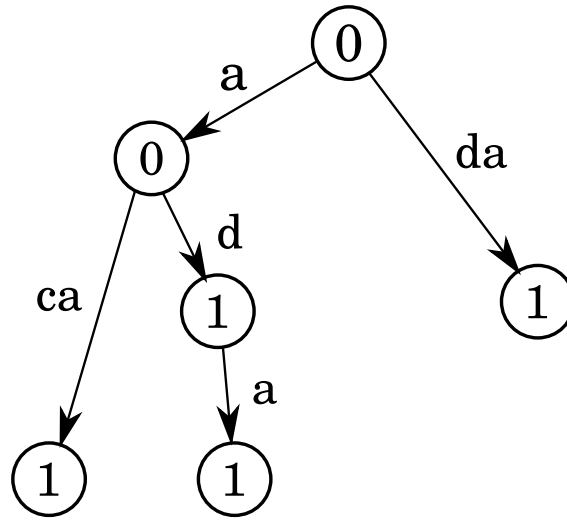
Jasno je da ovaj algoritam radi u složenosti $O(1)$ ali zato koristi $O(n \log n)$

memorije za nizove $C^{(k)}$. Moguće je dostići istu vremensku složenost za poređenje sa samo $O(n)$ dodatne memorije korišćenjem LCP niza. Ideja je da, pri poređenju podstringova $s_{[i,i+l)}, s_{[j,j+l)}$ nađemo $d = LCP(s_{[i,i+l)}, s_{[j,j+l)})$, Ukoliko je $d < l$, samo uporedimo karaktere s_{i+d}, s_{j+d} , inače, stringovi su jednaki. Na osnovu teoreme 5.2 problem se svodi na nalaženje minimuma u podsegmentu niza. Ovaj problem se može efikasno rešiti u vremenskoj složenosti $O(1)$ sa $O(n)$ preprocesiranja i $O(n)$ utrošene memorije.⁶

5.2 Sufiks stablo

Sufiks stablo je struktura podataka koja uopštava sufiks niz i usko je povezana sa njim. Da bismo definisali sufiksno stablo, definišimo prvo kompresovano prefiksno stablo.

Definicija 5.2 *Kompresovano prefiksno stablo za skup nepraznih stringova P se dobija tako što se u prefiksnom stablu obrišu svi čvorovi višestrukosti 0 sa tačno jednim detetom, osim korena, a zatim se dodaju grane koje odgovaraju obrisanim putevima u stablu i imaju labele koje odgovaraju obrisanim putevima.*



Kompresovano prefiksno stablo za skup $\{aca, ad, ada, da\}$.

U primeru sa slike, čvor $\delta(ad)$ ne brišemo iako on ima jedno dete. Na-
jzad, definišimo sufiks stablo.

Definicija 5.3 *Sufiks stablo za string s je kompresovano prefiksno stablo svih nepraznih sufiksa stringa s .*

Iako na prvi pogled deluje da bi ovakvo stablo zauzimalo previše memorije i iz tog razloga bilo nepraktično, pokazuje se suprotno.

Teorema 5.4 *Kompresovano prefiksno stablo od n stringova sadrži najviše $2n$ čvorova.*

Dokaz. Indukcijom po n . Za $n = 1$ tvđenje očigledno važi. Posmatrajmo stablo za skup stringova $P' = P - p$. Neka je q najduži prefiks stringa p koji se javlja u P' . Ukoliko ne postoji čvor $\delta(q)$, put od korena sa labelom q će odgovarati unutrašnjosti neke grane, pa ovu granu delimo na dve umetanjem novog čvora koji će onda biti $\delta(q)$. Zatim, ukoliko je $|q| < |p|$ dodajemo novu granu iz čvora $\delta(q)$ ka novom čvoru i upisujemo joj labelu $p_{[|q|,|p|)}$. Broj čvorova se povećao za najviše 2, što dokazuje tvđenje. \square

Naizgled, čak iako stablo ima mali broj čvorova, ukupna dužina labela svih grana može biti $\Theta(n^2)$, što bi impliciralo veliku memorijsku složenost. Kod sufiks stabla, sve labele grana su podstringovi stringa s pa umesto stringa labelu možemo predstaviti kao par l, r koji bi označavao da je labela te grane jednaka stringu $s_{[l,r)}$.

Postoje algoritmi koji konstruišu sufiks stablo direktno iz datog stringa u linearnom vremenu.⁹ U nastavku će biti prezentovan jednostavan algoritam koji na osnovu izračunatih sufiks i LCP nizova konstruiše sufiks stablo.

Struktura čvora sufiks stabla

```

1 struct stree_node {
2     int l, r, id;
3     stree_node* p;
4     map<char, stree_node*> next;
5     stree_node(int l, int r, int id, stree_node* p)
6         : l(l), r(r), id(id), p(p) {}
7 };

```

Svaki čvor sufiks stabla pamti labelu grane koja spaja taj čvor i njegovog roditelja, odnosno podstring stringa s , pokazivač na svog roditelja kao i pokazivače na svoju decu. Umesto višestrukosti pamti se redni broj id takav da je taj čvor jednak $\delta(s_{[id,n)})$ ukoliko postoji, inače -1 , što odgovara čvorovima višestrukosti 0. Za koren stabla se ne definiše roditelj a brojeve

l, r postavljamo na 0. Iz definicije sufiks stabla zaključujemo da sve grane koji izlaze iz jednog čvora imaju labele kojima se razlikuje prvo slovo, što opravdava izbor mape za čuvanje izlaznih grana. Algoritam radi na sledeći način. U leksikografskom redosledu se dodaju sufiksi stringa s . Ovaj redosled je određen sufiks nizom p . Prvi string se dodaje kao zasebna grana. Za svaki naredni string, odnosno i -ti za $i > 0$, važi da je najduži zajednički prefiks njega i bilo kog drugog stringa upravo jednak q_{i-1} (Teorema 5.2). Penjemo se od kraja prethodno ubačenog sufiksa do dubine q_{i-1} u stablu, mereno po broju slova u labelama. Ukoliko se pozicija na ovoj visini nalazi na sredini grane, delimo tu granu na dva dela. Konačno, ostaje nam da dodamo deo završni deo sufiksa dužine $(n - p_i) - q_{i-1}$, naravno, samo ako je on neprazan.

Implementacija algoritma za konstrukciju sufiks stabla

```

1  using node = stree_node;
2  node* suffix_tree(
3      const string& s,
4      const vector<int>& p,
5      const vector<int>& q
6  ) {
7      int n = s.size();
8      auto root = new node(0, 0, 0, 0);
9      auto curr = root->next[s[p[0]]]
10         = new node(p[0], n, p[0], root);
11      for (int i=1; i<n; i++) {
12          int t = q[i-1], d = n - p[i-1];
13          while (d && d - t >= (curr->r - curr->l)) {
14              d -= curr->r - curr->l;
15              curr = curr->p;
16          }
17          if (d > t) {
18              int m = curr->r - (d - t);
19              node* mid = new node(curr->l, m, -1, curr->p);
20              mid->next[s[m]] = curr;
21              curr->p->next[s[curr->l]] = mid;
22              curr->l = m;
23              curr->p = mid;
24              curr = mid;
25          }
26          if (p[i] + t < n)
27              curr = curr->next[s[p[i] + t]]
28                  = new node(p[i] + t, n, p[i], curr);
29      }
30      return root;
31  }

```

Vremenska složenost ovog algoritma je $O(n)$. Ovo se može dokazati tako što posmatramo vrednost $2i - y$, gde je y dubina čvora $curr$ u stablu, merena po broju grana. Svaka iteracija i spoljne *for* i unutrašnje *while* petlje uvećava vrednost ovog izraza za bar 1, pa je njihov ukupan broj manji od $2n$.

Sufiks stablo omogućava da se nađu prvo, poslednje, broj pojavljivanja nekog stringa t unutar stringa s u linearnom vremenu procedurom spuštanja slično kao kod prefiksnog stabla. Kod traženja svih pojavljivanja stringa korisna je sledeća teorema:

Teorema 5.5 *Kod sufiksnog stabla ne postoji podstablo koje sadrži više čvorova višestrukosti 0 nego čvorova višestrukosti 1.*

Dokaz. Pretpostavimo suprotno, da podstablo ima k čvorova i da čvorova višestrukosti 0 ima više od $\frac{k}{2}$. Na osnovu definicije sufiksnog stabla, svaki takav čvor ima bar dvoje dece, pa je ukupan broj dece svih čvorova u podstablu veći od k , što je nemoguće, jer je ukupan broj dece svih čvorova tačno $k - 1$.

Ovo nam omogućava da, za neki string t , od pozicije $\delta(t)$, ukoliko postoji, izvršimo pretragu u dubinu ili širinu da bismo našli sve čvorove višestrukosti 1 u podstablu pozicije $\delta(t)$. Ako neki ovakav čvor odgovara sufiksu id , znamo da je $t = s_{[id, id+|t|]}$ odnosno, id je pozicija jednog pojavljivanja stringa t u stringu s . Kako je složenost pretrage linearna po veličini podstabla, i kako je ukupan broj čvorova podstabla ne više od duplo veći od broja čvorova višestrukosti 1, vremenska složenost algoritma je $O(|t| + k)$, gde je k broj pojavljivanja stringa t u s .

Od konstruisanog sufiksnog stabla se može dobiti sufiks niz traženjem praznog stringa, ili ekvivalentno, puštanjem pretrage u dubinu od korena i pamćenjem id vrednosti za sve čvorove za koje je $id \neq -1$.

Nalaženje svih pojavljivanja stringa t u stringu s

```
1  using node = stree_node;
2  vector<int> stree_find_all(
3      node* root,
4      const string& s,
5      const string& t
6  ) {
7      vector<int> result;
8      node* curr = root;
9      int pos = 0;
10     for (char x : t)
11         if (pos == curr->r) {
12             if (!curr->next.count(x))
13                 return result;
14             curr = curr->next[x];
15             pos = curr->l + 1;
16         } else if (x != s[pos])
17             return result;
18         else
19             pos++;
20     vector<node*> q = {curr};
21     size_t qs = 0;
22     while (qs != q.size()) {
23         node* tmp = q[qs++];
24         if (tmp->id != -1)
25             result.push_back(tmp->id);
26         for (auto [x, ch] : tmp->next)
27             q.push_back(ch);
28     }
29     return result;
30 }
```

Napomena. Prikazani kod ne vraća sve pozicije u rastućem, već u proizvoljnom redosledu.

5.3 Sufiks automat

5.3.1 Definicija

Sufiks automat je struktura podataka koja se gradi od datog stringa, a iz koje se mogu izvući različite korisne informacije o samom stringu i može se koristiti za brzu pretragu podstringova, slično sufiks nizu i sufiks stablu.

Glavna prednost sufiks automata je ta što ima veoma jednostavan algoritam za konstrukciju koji radi u linearnom vremenu.

Definicija 5.4 *Sufiks automat je parcijalni konačni automat sa minimalnim brojem čvorova koji prepoznaje sve sufikse stringa s .*

Parcijalni konačni automat koji raspoznaje skup stringova $S \subseteq \Sigma^*$ je usmeren graf (V, E) , gde svaka grana ima labelu koja je slovo iz alfabeta Σ , zajedno sa specijalnim čvorom $t_0 \in V$, i skupom čvorova $T \subseteq V$, takav da je $s \in S$ akko postoji put od čvora v_0 do nekog čvora $t \in T$ čiji je niz labela s , i nijedan čvor nema dve izlazne grane sa istom labelom.

Pri određivanju oblika sufiks automata, od velike koristi je *endpos* funkcija. Ispostaviće se da upravo ova funkcija određuje stanja, odnosno čvorove automata. Na dalje, smatrajmo da je s fiksiran string dužine n za koji konstruišemo sufiks automat.

Definicija 5.5 *Za string p , $endpos(p)$ je skup celih brojeva takav da je $i \in endpos(p)$ akko je $|p| \leq i \leq |s|$ i $s_{[i-|p|, i]} = p$.*

Drugim rečima, $endpos(p)$ je skup svih pojavljivanja stringa p u s , gde za poziciju uzimamo desni kraj.

Definicija 5.6 *Za stringove p_1, p_2 , koji su podstringovi stringa s kažemo da su $endpos$ -ekvivalentni ukoliko važi $endpos(p_1) = endpos(p_2)$.*

Ukoliko su p_1, p_2 $endpos$ -ekvivalentni, tada je jedan od njih sufiks drugog. Svaka klasa ekvivalencije se sastoji od nekoliko stringova koji se mogu poređati u niz takav da je svaki naredni string sufiks prethodnog i ima dužinu za jedan manju. Samim tim, za predstavnika klase možemo uzeti najduži string te klase. Obrat ovog tvrđenja ne važi, ali važi da ako je p_1 sufiks stringa p_2 , da je $endpos(p_2) \subseteq endpos(p_1)$. Takođe, ako p_1 nije sufiks od p_2 i p_2 nije sufiks od p_1 , tada je $endpos(p_1) \cap endpos(p_2) = \emptyset$.¹⁰

Sufiks automat kao čvorove ima sve klase $endpos$ -ekvivalencije takve da je njima pridružen $endpos$ skup neprazan. Sve izlazne grane nekog čvora se određuju na sledeći način. Neka je u čvor automata odnosno jedna klasa ekvivalencije i neka je p bilo koji string klase u . Za svako $x \in \Sigma$, ako je px podstring od s , onda postoji grana od u do klase koja sadrži px . Primetimo da rezultat ne zavisi od izbora stringa p iz klase u . Zaista, $endpos(px)$ je skup svih brojeva $i + 1$ takvih da je $s_i = x, i \in endpos(p)$, pa grana ide ka klasi koja odgovara vrednosti $endpos(px)$. Ako posmatramo najkraći string

klase ekvivalencije u , osim klase koja sadrži prazan string, uklanjanjem prvog slova tog stringa dobijamo drugačiju klasu ekvivalencije v , odnosno klasu koja sadrži prethodni *endpos* skup kao svoj strogi podskup. Kažemo da postoji sufiks veza od čvora u do čvora v u automatu. Svi čvorovi imaju jedinstvenu izlaznu sufiks vezu, osim čvora koji odgovara praznom stringu, i veza uvek ide ka klasama koje sadrže kraće stringove, pa je dobijeni graf sufiks veza stablo.

5.3.2 Algoritam za konstrukciju

Opišimo algoritam koji konstruiše sufiks automat. Pored prethodno opisanih izlaznih grana, svaki čvor grafa će pamtit i dužinu najdužeg stringa svoje klase ekvivalencije, kao i sufiks vezu.

Struktura čvora sufiks automata

```

1 struct sautomaton_node {
2     int len;
3     sautomaton_node* link;
4     map<char, sautomaton_node*> next;
5 };

```

Algoritam radi u iteracijama, dodajući jedno po jedno slovo stringa s . Nakon k -te iteracije dobijeni graf odgovara sufiks automatu za string $s_{[0,k)}$. Algoritam počinje inicijalizacijom početnog čvora koji odgovara klasi koja sadrži prazan podstring. Dakle, *len* se postavlja na 0, *link* na *null*-pokazivač, a skup prelaza je prazan. Takođe, u svakom trenutku pamtimo i pokazivač na čvor koji odgovara klasi koja sadrži ceo string do tog trenutka, to je čvor *last*.

Posmatrajmo promene koje se dese na grafu nakon dodavanja karaktera s_k . Sigurno će se javiti nova klasa koja će sadržati ceo string $s_{[0,k+1)}$ i možda još neke njegove sufikse. Tačnije, ta klasa će sadržati sve sufikse stringa $s_{[0,k+1)}$ koji se ne javljaju već u stringu $s_{[0,k)}$. Taj novi čvor nazovimo *curr*. Njegova dužina je $k + 1$, odnosno za jedan veća nego kod čvora *last*. Koji sve čvorovi imaju grane koje idu ka čvoru *curr*? To su neki od čvorova koji odgovaraju stanjima koja sadrže sufikse stringa $s_{[0,k)}$. Podelimo sve ove sufikse u dve grupe, na one koji se javljaju u stringu $s_{[0,k)}$ na takav način da se neko od tih pojavljivanja može produžiti slovom s_k unutar $s_{[0,k)}$, i one kod kojih to ne važi. Jasno je da će na ovaj način ovi sufiksi biti podeljeni

po dužini u odnosu na neku granicu. Svim stanjima koja odgovaraju dužim sufiksima treba dodati granu sa labelom s_k ka čvoru *curr*, dok kraći sufiksi već imaju takvu granu, i te grane pokazuju na već postojeće čvorove, pa ne treba raditi ništa. Sve ove čvorove možemo obići tako što krenemo od čvora *last* i krećemo se duž sufiks veza.

Neka je p čvor na kojem smo se zaustavili. Podsetimo se da sufiks veza iz nekog čvora ide ka čvoru koji sadrži najduži sufiks stringa tog čvora koji se javlja na više mesta od njega samog. Za čvor *curr* potrebno je naći string koji je sufiks stringa $s_{[0,k+1)}$ a javlja se još negde u stringu $s_{[0,k+1)}$, odnosno, koji se javlja negde u $s_{[0,k)}$. Drugim rečima, potrebno je pronaći najduži sufiks stringa $s_{[0,k)}$ koji se javlja na nekoj poziciji koja se može produžiti slovom s_k , a to je upravo čvor p . Ako se slovo s_k ne javlja nigde u stringu $s_{[0,k)}$, onda će se p zaustaviti na *null* pokazivaču, odnosno, biće obišeni svi sufiksi, uključujući i prazan string. Tada sufiks veza ide ka korenu automata. U suprotnom, zaustavili smo se zato što p ima izlaznu granu sa labelom s_k . Neka ta grana ide ka čvoru q . Sufiks veza čvora *curr* treba da ide ka čvoru čija je dužina $len(p) + 1$. Jasno je da je $len(q) > len(p)$. Ukoliko se ove dužine razlikuju za tačno 1, to znači da su stringovi iz čvora q tačno oni koji su u čvoru p a mogu se produžiti za slovo s_k . Ovo znači da sufiks veza iz čvora *curr* treba da pokazuje upravo na q .

U suprotnom, klasa čvora q se mora podeliti na dve klase, jednu u kojoj su svi stringovi dužine $len(p) + 1$ i kraći, i sve ostale. Za ovu kraću klasu kreiramo novi čvor koji nazivamo *clone*, njegove osobine su iste kao za čvor q , osim što je $len(clone) = len(p) + 1$. Sufiks veze čvora *curr*, kao i čvora q idu ka čvoru *clone*. Međutim, potrebno je uraditi i sledeće. Sve grane koje su polazile iz čvora p ili njegovih sufiks-predaka a preko kojih se dolazilo do čvora q sad treba da pokazuju na čvor *clone*. Razlog za to je što čvor q sada sadrži samo stringove dužine bar $len(p) + 2$, što znači da svi produžeci slovom s_k iz p i njegovih sufiks-predaka zapravo treba da idu ka čvoru *clone*. Ovime smo popravili sve sufiks veze i grane i dobili novi automat za string $s_{[0,k+1)}$.

Glavni algoritam za konstrukciju sufiks automata

```
1 using node = sautomaton_node;
2 node* sautomaton(const string& s) {
3     node* root, *curr;
4     root = curr = new node{0, nullptr, {}};
5     for (char x : s)
6         curr = sautomaton_extend(root, curr, x);
7     return root;
8 }
```

Algoritam za proširenje automata jednim slovom

```
1 using node = sautomaton_node;
2 node* sautomaton_extend(node* root, node* last, char x) {
3     node* curr = new node{last->len+1, nullptr, {}};
4     node* p = last;
5     for (; p && !p->next.count(x); p = p->link)
6         p->next[x] = curr;
7     if (!p) {
8         curr->link = root;
9     } else {
10        node* q = p->next[x];
11        if (p->len + 1 == q->len) {
12            curr->link = q;
13        } else {
14            node* clone = new node(*q);
15            clone->len = p->len + 1;
16            for (; p && p->next[x] == q; p = p->link)
17                p->next[x] = clone;
18            curr->link = q->link = clone;
19        }
20    }
21    return curr;
22 }
```

Pošto svaka ekstenzija kreira najviše dva nova čvora, sufiks automat ne može imati više od $2n + 1$ čvorova. Ova granica se može poboljšati na $2n - 1$ ako primetimo da prve dve ekstenzije ne mogu da kreiraju klonove.

Može se pokazati da je ukupan broj prelaza automata ne više od $3n - 4$ za $n \geq 3$. Vremenska složenost celog algoritma je $O(n)$, uz pretpostavku da je veličina alfabeta fiksna.¹⁰

Skup stanja T automata koji raspoznaje sufikse stringa s se može dobiti tako što nakon poslednje iteracije, počev čvora *curr* krećući se po sufiks vezama obiđemo sve čvorove sve do korena, uključujući i njega.

5.3.3 Uopštenja i primene

Ako u čvorovima pamtimo dodatne informacije koje se tiču sufiks automata, možemo rešiti naizgled nevezane probleme. Zajedničko skoro svim primenama je da koriste dinamičko programiranje na stablu sufiks veza, ili na acikličnom grafu prelaza. Za dinamičko programiranje na acikličnom grafu nam je neophodan topološki redosled čvorova tog grafa, koji se može jednostavno dobiti sortiranjem čvorova po vrednosti *len*, što se može uraditi i u linearnoj složenosti *counting sort*-om. Korisno je i to da isti ovaj redosled odgovara i obrnutom topološkom redosledu stabla sufiks veza.

Broj različitih podstringova

Kako svaki podstring odgovara putu od korena stabla do nekog čvora, broj različitih podstringova se može dobiti i kao broj različitih puteva u acikličnom grafu. Neka je u čvor grafa a $N(u)$ skup njegovih suseda. Ako označimo sa $d(u)$ broj puteva koji počinju u čvoru u , onda važi sledeća rekurentna veza:

$$d(u) = 1 + \sum_{v \in N(u)} d(v) \quad (5.3.1)$$

Rešenje je vrednost $d(\text{root})$. Vrednosti računamo u obrnutom topološkom redosledu.

Broj pojavljivanja datog stringa

Za string p nalazimo čvor $\delta(p)$. Ukoliko taj čvor ne postoji, string p se ne javlja u s . U suprotnom, potrebno je naći veličinu skupa $\text{endpos}(p)$, odnosno veličinu endpos skupa za čvor $u = \delta(p)$. Ove veličine se mogu izračunati nakon konstrukcije automata za sve čvorove pomoću dinamičkog programiranja na sledeći način.

Posmatrajmo prefiks $s_{[0,k]}$. Ukoliko je $k \in \text{endpos}(p)$, postoji put pomoću sufiks veza od čvora $\delta(s_{[0,k]})$ do čvora u . U suprotnom, takav put ne postoji.

Dakle, nas zanima koliko različitih čvorova koji odgovaraju prefiksima stringa s se nalazi u podstablu stabla sufiks veza u čvoru u . Čvorovi koji odgovaraju prefiksima stringa s su tačno oni čvorovi koji nisu klonovi, odnosno koji su kreirani na početku svake ekstenzije. Podsetimo se da prilikom kloniranja originalnom čvoru ostaju svi duži sufiksi, pa string $\delta(s_{[0,k]})$ ne može da pređe u klonirani čvor. Najzad, neka je $cl(v) = 1$ ako je v klon nekog čvora, a 0 inače. Broj nekloniranih čvorova $d(u)$ u podstablu čvora u zadovoljava sledeću rekurentnu vezu:

$$d(u) = (1 - cl(u)) + \sum_{v, u=suff(v)} d(v) \quad (5.3.2)$$

Ovo procesiranje automata ima vremensku složenost $O(n)$. Broj pojavljivanja stringa p se onda može naći u linearnom vremenu po veličini stringa p .

Sva pojavljivanja datog stringa

Iskoristićemo prethodno ustanovljenu činjenicu da nas zanimaju samo čvorovi u podstablu čvora u koji nisu klonovi. Ukoliko pustimo pretragu u dubinu iz čvora u i očitamo sve vrednosti len , dobićemo upravo *endpos* skup za čvor u . Umanjenjem svih elemenata rezultata za $|p|$ dobijamo traženi skup pojavljivanja. Jednostavno se pokazuje da u svakom podstablu stabla sufiks veza nema više od polovine klonova. Naime, ako je čvor x klon, čvor čiji je on klon ima sufiks vezu upravo ka njemu, pa će biti u njegovom podstablu. Klonovi ne mogu imati svoje klonove, a ostali čvorovi ne mogu imati više od jednog klona, odakle sledi da bar polovina svih čvorova bilo kog podstabla čine čvorovi koji nisu klonovi. Odavde sledi da pretraga u dubinu ima vremensku složenost $O(k)$, gde je k broj pojavljivanja stringa p u s , pa je vremenska složenost celog algoritma $O(|p| + k)$. Za potrebe pretrage u dubinu po stablu sufiks veza moramo da zapamtimo za svaki čvor sve sufiks veze koje ulaze u njega, što možemo lako uraditi po završetku konstrukcije automata.

6 Heširanje

Za razliku od kriptografskih heš funkcija, heš funkcije koje se koriste u pretraži stringova zadovoljavaju određene relacije, što omogućava dinamičko održavanje vrednosti heš funkcije i nakon izmena stringa, njihovog spajanja ili sečenja. U opštem slučaju, heš funkcija je funkcija koja slika skup stringova Σ^* u skup celih brojeva iz nekog opsega, najčešće, za neki prirodan broj M , taj opseg je $\{0, 1, \dots, M - 1\}$.

6.1 Rabin-Karp algoritam

Pomoću Rabin-Karp algoritma mogu se pronaći sva pojavljivanja jednog stringa u drugom u linearnoj vremenskoj složenosti, ako se dozvoli mala verovatnoća greške, odnosno *false positive*-a, što znači da algoritam prepozna je podstring stringa s kao string p iako to nije slučaj.

Kod ovog algoritma, heš funkcija se definiše na sledeći način. Neka je $|s| = n$.

$$h(s) = \left(\sum_{i=0}^{n-1} f(s_i) q^{n-i-1} \right) \mod M \quad (6.1.1)$$

Ovde je q ceo broj, $f : \Sigma \rightarrow \{1, \dots, M - 1\}$ je funkcija koja mapira karaktere u brojeve. U praksi, najjednostavnije je uzeti da je $f(c)$ ASCII vrednost karaktera c . Druga opcija je da se svim slovima iz Σ nasumično dodele različite vrednosti $f(c)$.

Teorema 6.1 *Prethodno definisana heš funkcija zadovoljava jednakost:*

$$h(sp) \equiv q^{|p|} h(s) + h(p) \mod M \quad (6.1.2)$$

Dokaz. Neka je $n = |s|, m = |p|$.

$$\begin{aligned}
h(sp) &\equiv \sum_{i=0}^{n+m-1} f((sp)_i) q^{n+m-i-1} \pmod{M} \\
&\equiv \sum_{i=0}^{n-1} f(s_i) q^{n+m-i-1} + \sum_{j=n}^{n+m-1} f(p_{j-n}) q^{n+m-j-1} \pmod{M} \\
&\equiv q^m \sum_{i=0}^{n-1} f(s_i) q^{n-i-1} + \sum_{j=0}^{m-1} f(p_j) q^{m-j-1} \pmod{M} \\
&\equiv q^m h(s) + h(p) \pmod{M}
\end{aligned}$$

Stepen q^k po modulu se može brzo izračunati u $O(\log k)$ pomoću binarnog stepenovanja.

Binarno stepenovanje

```

1 int modpow(int x, int k, int m) {
2     if (k == 0)
3         return 1 % m;
4     int y = modpow(x, k >> 1, m);
5     y = 1ll * y * y % m;
6     if (k & 1)
7         y = 1ll * y * x % m;
8     return y;
9 }

```

Alternativa je da računamo sve stepene q^0, q^1, \dots, q^k , zapamtimo ih u niz i po potrebi uzimamo iz tog niza. Za k možemo uzeti dužinu najdužeg stringa koji želimo da obrađujemo. Za broj M se često uzima veliki prost broj. Za broj q je najbolje uzeti primitivni koren po modulu M .

Računanje svih stepena broja x od 0 do k

```
1 vector<int> modpowvec(int x, int k, int m) {  
2     vector<int> a(k+1);  
3     a[0] = 1;  
4     for (int i=1; i<=k; i++)  
5         a[i] = 1ll * a[i-1] * x % m;  
6     return a;  
7 }
```

Da bismo mogli efikasno da tražimo heš vrednosti podstringova stringa s , dovoljno je da izračunamo za svaki prefiks stringa s vrednost $h(s_{[0,i]})$. To se može uraditi korišćenjem relacije 6.1.2.

Računanje heš vrednosti za sve prefikse stringa

```
1 vector<int> prefixhash(const string& s, int q, int m) {  
2     int n = s.size();  
3     vector<int> h(n+1);  
4     h[0] = 0;  
5     for (int i=1; i<=n; i++)  
6         h[i] = (h[i-1] * 1ll * q + s[i-1]) % m;  
7     return h;  
8 }
```

Za izračunavanje $h(s_{[l,r]})$ možemo iskoristiti relaciju 6.1.2 primenjenu na stringove $s_{[0,l]}$ i $s_{[l,r]}$.

$$h(s_{[l,r]}) \equiv (h(s_{[0,r]}) - q^{r-l}h(s_{[0,l]})) \mod M \quad (6.1.3)$$

Rabin-Karp algoritam nalazi sva pojavljivanja stringa p u stringu s tako što izračuna prefiksne heš vrednosti za string s . Za string p je dovoljno izračunati samo $h(p)$. Zatim se pomoću relacije 6.1.3 nalaze heš vrednosti za sve podstringove stringa s odgovarajuće dužine, i ove vrednosti se upoređuju sa $h(p)$. U slučaju poklapanja, postoje dve opcije. Jedna je da se zadovoljimo time da postoji mala verovatnoća greške i da prihvatimo pronađenu poziciju kao pojavljivanje stringa p . Druga je da za svako potencijalno pojavljivanje proverimo da li se zaista radi o pojavljivanju stringa p u složenosti $O(m)$ po jednom pojavljivanju. U najgorem slučaju, drugi algoritam ima složenost istu kao naivna pretraga, odnosno $O(nm)$.

Rabin-Karp algoritam bez provere poklapanja

```
1 vector<int> rabin_karp(  
2     const string& s, const string& p,  
3     int q, int M  
4 ) {  
5     int n = s.size(), m = p.size();  
6     int hp = prefixhash(p, q, M)[m], qm = modpow(q, m, M);  
7     vector<int> h = prefixhash(s, q, M), r;  
8     for (int i=0; i<=n-m; i++)  
9         if ((h[i+m] - h[i]*111*qm - hp) % M == 0)  
10             r.push_back(i);  
11     return r;  
12 }
```

Pored pretrage podstringova, heš funkciju iz Rabin-Karp algoritma možemo koristiti i za konstrukciju strukture podataka koja čuva string, održava heš vrednost za sve svoje podstringove i omogućava izmene stringa na bilo kojoj poziciji. Takva struktura bi bila balansirano binarno stablo, najbolje *splay* stablo ili *treap*, koja bi u svakom čvoru čuvala heš vrednost i broj karaktera u svom podstablu. Sve operacije (sečenje, spajanje stabala, provera heš vrednosti dela stringa, izmena slova) se mogu implementirati da rade u vremenskoj složenosti $O(\log n)$, gde je n ukupna dužina stringa.

6.2 Kolizije

6.2.1 Odabir parametara heš funkcije

Posmatrajmo interakciju dva agenta, nazovimo ih Ana i Branko. Ana dizajnira heš funkciju, odnosno bira vrednosti q, M dok Branko nalazi različite stringove s, p . Brankov cilj je da se ovi stringovi heširaju u istu vrednost, dok je Anin cilj da se to izbegne.

U svakom slučaju, ukoliko Branko zna vrednosti q, M koje je Ana izabrala, u najgorem slučaju primenom grube sile može se naći par različitih stringova s, p koji imaju istu heš vrednost. Međutim, ukoliko to nije slučaj, odnosno ako Ana bira parametre slučajno bez Brankovog znanja, pokazuje se da je verovatnoća kolizije relativno mala.

Teorema 6.2 *Ako je M prost broj, q se uzima slučajno i uniformno iz skupa*

$\{0, \dots, M-1\}$ i ako su stringovi s, p dužine n , tada je verovatnoća kolizije ne više od $\frac{n-1}{p}$.

Dokaz. Posmatrajmo jednu koliziju, i neka je Q slučajna promenljiva iz koje q uzima vrednost. Sve jednakosti su po modulu M .

$$\begin{aligned} h(s) &= h(p) \\ \sum_{i=0}^{n-1} f(s_i)Q^{n-1-i} &\equiv \sum_{i=0}^{n-1} f(p_i)Q^{n-1-i} \\ \sum_{i=0}^{n-1} (f(s_i) - f(p_i))Q^{n-1-i} &\equiv 0 \end{aligned}$$

Leva strana je nenula polinom pa promenljivoj Q stepena najviše $n-1$. U polju \mathbb{F}_M nenula polinom stepena do $n-1$ ne može imati više od $n-1$ nula, pa je verovatnoća da je Q nula polinoma ne više od $\frac{n-1}{p}$. \square

Ovaj argument je validan samo ukoliko je M prost broj. Posmatrajmo slučaj kada je M stepen dvojke, na primer $M = 2^{64}$. Ovakav odabir je prirodan jer aritmetičke operacije $+$, $-$, \times sa *unsigned* tipovima rade po modulu 2^k , gde je k broj bitova tog tipa. Postoje različiti stringovi koji imaju istu heš vrednost bez obzira na to koji broj q je izabran.

Neka je $u(s)$ funkcija koja menja sva pojavljivanja slova a slovom b i obratno, na primer $u(aab) = bba$. Thue-Morse niz¹¹ je definisan na sledeći način: Krećemo od stringa $t_0 = a$. Za svako $i \geq 0$ uzimamo da je $t_{i+1} = t_i u(t_i)$. Prvih nekoliko elemenata su $t_1 = ab$, $t_2 = abba$, $t_3 = abbabaab$. Kako je t_i prefiks stringa t_{i+1} , možemo reći i da se radi o jednoj beskonačnoj sekvenci slova iz $\{a, b\}$.

Teorema 6.3 Za parno k važi $t_k = t_{k-1} \overline{t_{k-1}}$. Za neparno k važi da je $u(t_k) = \overline{t_k}$.

Dokaz. Za $k = 1$ tvrđenje očigledno važi. Dokažimo tvrđenje indukcijom po k . Ako je k parno, onda je $k-1$ neparno pa po induktivnoj pretpostavci važi $t_k = t_{k-1} u(t_{k-1}) = t_{k-1} \overline{t_{k-1}}$. Štaviše, t_k je palindrom za parno k , jer je $\overline{t_k} = \overline{t_{k-1} \overline{t_{k-1}}} = \overline{\overline{t_{k-1}}(t_{k-1})} = t_{k-1} \overline{t_{k-1}} = t_k$. Iskoristili smo osobine okretanja stringa $\overline{ab} = (\overline{b})(\overline{a})$ i $\overline{\overline{a}} = a$.

Za k neparno imamo da je $u(t_k) = u(t_{k-1}u(t_{k-1})) = u(t_{k-1})u(u(t_{k-1})) = u(t_{k-1})t_{k-1}$. Sa druge strane imamo $\overline{t_k} = \overline{t_{k-1}u(t_{k-1})} = (\overline{u(t_{k-1})})(\overline{t_{k-1}}) = u(\overline{t_{k-1}})\overline{t_{k-1}} = u(t_{k-1})t_{k-1}$. Poslednja jednakost važi jer je $k-1$ parno, pa je t_{k-1} palindrom. \square

Teorema 6.4 Za $M = 2^{64}$, proizvoljno f , neparno q i $s = t_{11}$, $h(s) = h(u(s))$.

Dokaz. Posmatrajmo sledeći niz polinoma. Neka je $p_0(x) = 1$, a za $n \geq 0$ je $p_{n+1} = p_n(x)(1 - x^{2^n})$. Prvih nekoliko elemenata su:

$$\begin{aligned} p_1(x) &= 1 - x \\ p_2(x) &= 1 - x - x^2 + x^3 \\ p_3(x) &= 1 - x - x^2 + x^3 - x^4 + x^5 + x^6 - x^7 \end{aligned}$$

Primetimo da znakovi ispred x^i čine upravo Thue-Morse niz. Dokažimo prvo da je $p_k(x) \equiv 0$ po modulu $2^{\frac{k(k-1)}{2}}$, za svako neparno x i $k \geq 1$. Neka je $q_k(x) = (1 - x^{2^k})$. Polinom $p_k(x)$ se po definiciji može faktorisati kao $q_0(x)q_1(x) \dots q_{k-1}(x)$.

Dokažimo da 2^k deli $q_k(x)$ za svako neparno x , indukcijom po k . Za $k = 0$ imamo $q_0(x) = 1 - x$, pa tvrdjenje važi jer je x neparno. Za $k > 0$ imamo da važi $q_k(x) = 1 - x^{2^k} = (1 - x^{2^{k-1}})(1 + x^{2^{k-1}}) = q_{k-1}(x)(1 + x^{2^{k-1}})$. Broj $1 + x^{2^{k-1}}$ je paran jer je x neparno, a po induktivnoj pretpostavci važi $2^{k-1} | q_{k-1}(x)$, pa $2 \cdot 2^{k-1} = 2^k$ deli njihov proizvod.

Oдавde sledi da $2^0 \cdot 2^1 \cdot \dots \cdot 2^{k-1}$ deli $p_k(x)$, odnosno da je $p_k(x) \equiv 0$ po modulu $2^{\frac{k(k-1)}{2}}$. Konačno, pokažimo da je $h(s) = h(u(s))$. Neka je $c = f(a) - f(b)$. Dužina stringova je $n = 2048$. Posmatrajmo vrednost $h(s) - h(u(s))$.

$$\begin{aligned}
h(s) - h(u(s)) &\equiv \sum_{i=0}^{n-1} f(s_i)q^{n-i-1} - \sum_{i=0}^{n-1} f(u(s_i))q^{n-i-1} \\
&\equiv \sum_{i=0}^{n-1} (f(s_i) - f(u(s_i)))q^{n-i-1} \\
&\equiv \sum_{i=0}^{n-1} (f(s_{n-1-i}) - f(u(s_{n-1-i})))q^i \\
&\equiv \sum_{i=0}^{n-1} (f(\bar{s}_i) - f(u(\bar{s}_i)))q^i
\end{aligned}$$

Sada koristimo teoremu 6.3 za $k = 11$.

$$\begin{aligned}
h(s) - h(u(s)) &\equiv \sum_{i=0}^{n-1} (f(\bar{s}_i) - f(u(\bar{s}_i)))q^i \\
&\equiv \sum_{i=0}^{n-1} (f(u(s_i)) - f(s_i))q^i
\end{aligned}$$

Član $f(u(s_i)) - f(s_i)$ ima vrednost $-c$ ako je $s_i = a$, inače ima vrednost c . Odavde dobijamo da je $h(s) - h(u(s)) = -c \cdot p_{11}(q)$, pošto je q neparno, $p_{11}(q) \equiv 0$ po modulu $2^{\frac{11+12}{2}} = 2^{66}$, a samim tim je 0 i po modulu M , pa je $h(s) - h(u(s)) \equiv 0$ tj. $h(s) = h(u(s))$. \square

Napomena: Ukoliko se dozvoli da je q parno, važiće $q^{64} \equiv 0 \pmod{M}$, pa će heš funkcija za npr. stringove a^{65} i ba^{64} dati istu vrednost.

6.2.2 Konstrukcija kolizije

Opišimo sada algoritme pomoću kojih se mogu naći različiti stringovi s, p iste dužine n takve da je $h(s) = h(p)$, za neke fiksne vrednosti parametara q, M . Pretpostavićemo da je M prost broj, osim ako nije drugačije naglašeno.

Definicija 6.1 *Multiplikativni red broja q po modulu M je najmanji prirodan broj r takav da je $q^r \equiv 1 \pmod{M}$.*

Iz teorije brojeva je poznata činjenica da multiplikativni red broja postoji akko su q, M uzajamno prosti, i važi da r deli $\varphi(M)$, gde je φ Ojlerova funkcija.

Ukoliko je multiplikativni red broja q po modulu M mali, moguće je jednostavno konstruisati koliziju – uzmimo stringove $a^r b^r$ i $b^r a^r$:

$$\begin{aligned} h(a^r b^r) &\equiv q^r \cdot h(a^r) + h(b^r) \\ &\equiv h(a^r) + h(b^r) \\ &\equiv q^r \cdot h(b^r) + h(a^r) \\ &\equiv h(b^r a^r) \end{aligned}$$

Ukoliko broj M nije ogroman, može se naći kolizija prostim generisanjem nasumičnih stringova. Verovatnoća nalaženja kolizije je oko $\frac{1}{2}$ za srazmerno mali broj generisanih stringova, njih $\Theta(\sqrt{M})$, o čemu govori čuveni paradoks rođendana.¹²

Implementacija jednostavnog algoritma za nalaženje kolizije

```

1 pair<string, string> birthday1(int q, int M, int n) {
2     map<int, string> d;
3     mt19937 eng(q^M);
4     while (1) {
5         string s(n, 0);
6         for (char& x : s)
7             x = uniform_int_distribution<char>('a', 'z')(eng);
8         int h = prefixhash(s, q, M).back();
9         if (d.count(h) && d[h] != s) {
10             return {d[h], s};
11         } else {
12             d[h] = s;
13         }
14     }
15 }
```

Ukoliko je broj M reda veličine 2^{64} , prethodna tehnika neće dati zadovoljavajuće rezultate. Naredna tehnika ima dosta bolje performanse, i može

se primeniti u slučaju da je f jednostavna linearna funkcija, na primer ako f vraća ASCII vrednost karaktera. Tačnije, dovoljan uslov je da postoje različiti karakteri $x, y, z \in \Sigma$ takvi da je $f(x) - f(y) = f(y) - f(z)$. Bez gubljenja opštosti, uzmimo da je $a, b, c \in \Sigma$ i $f(b) - f(a) = f(c) - f(b)$. Algoritam kao ulazne parametre ima q, M , kao i dužinu stringa n . Algoritam konstruiše string koji ima istu heš vrednost kao string b^n , tako što rešava sistem jednačina

$$\sum_{i=0}^{n-1} \alpha_i \cdot q^{n-1-i} \equiv 0 \pmod{M} \quad (6.2.1)$$

gde su $\alpha_i \in \{-1, 0, 1\}$ nepoznate. Iz ovog rešenja se jednostavno konstruiše string, ako je $\alpha_i = -1, 0, 1$ redom, onda je $s_i = a, b, c$. Algoritam održava kolekciju klastera. U početku postoji n klastera, svaki od njih odgovara jednom vektoru parametara α , naime, kod i -tog klastera važi $\alpha_j = 1$ ako je $i = j$, a $\alpha_j = 0$ inače. Nakon ovoga, algoritam spaja klastere koji imaju bliske heš vrednosti.

$$H(\alpha) = \left(\sum_{i=0}^{n-1} \alpha_i \cdot q^{n-1-i} \right) \pmod{M} \quad (6.2.2)$$

Cilj je da dobijemo da jedan klaster ima vrednost 0. Najjednostavniji način da spajamo klastere jeste da ih sortiramo po $H(\alpha)$, i da zatim spajamo dva po dva tako što njihove parametre α , a samim tim i H -vrednosti oduzmemo. Štaviše, radi efikasnosti, nećemo eksplicitno čuvati sve parametre α već ćemo klastere čuvati u obliku binarnog stabla. Tek po završetku rada iz konkretnog stabla izvlačimo vrednosti α_i .

Struktura čvora heš klastera

```

1 struct hash_cluster {
2     hash_cluster* l, *r;
3     long h;
4     int id;
5 };

```

Glavna funkcija kao parametre prima samo brojeve q, M, n i, ukoliko uspe, vraća par stringova s, p dužine n koji imaju istu heš vrednost. Za

razliku od prethodnih, ova implementacija radi sa 64-bitnim brojevima (tip `long`), dok pri množenju po modulu koristi tip `__int128_t` koji postoji u GCC kompajleru.

Glavna funkcija za nalaženje heš kolizije

```

1 pair<string, string> hash_merge(long q, long M, int n) {
2     vector<hash_cluster*> c(n);
3     long qi = 1;
4     for (int i=n-1; i>=0; i--) {
5         c[i] = new hash_cluster {nullptr, nullptr, qi, i};
6         qi = (__int128_t)qi * q % M;
7     }
8     auto cmp = [](hash_cluster* a, hash_cluster* b) {
9         return a->h < b->h;
10    };
11    while (c.size()) {
12        sort(c.begin(), c.end(), cmp);
13        if (c[0]->h == 0) {
14            string s(n, 'b'), p = s;
15            hash_collect(c[0], p, 1);
16            return {s, p};
17        }
18        if (c.size() == 1)
19            break;
20        int j = 0;
21        vector<hash_cluster*> d;
22        if (c.size() % 2) {
23            d.push_back(c[0]);
24            j++;
25        }
26        for (; j != (int)c.size(); j+=2) {
27            auto t = new hash_cluster
28                {c[j+1], c[j], c[j+1]->h - c[j]->h, 0};
29            d.push_back(t);
30        }
31        swap(c, d);
32    }
33    return {"", ""};
34 }

```

Oslobađanje memorije zauzete operatorom `new` je izostavljeno zbog preglednosti. Nakon konstruisanog stabla koje ima heš vrednost 0, od njega konstruišemo string `s` koristeći sledeću rekurzivnu funkciju.

Funkcija za rekonstrukciju stringa koji odgovara heš klasteru

```
1 void hash_collect(hash_cluster* root, string& s, int sgn) {  
2     if (root->l) {  
3         hash_collect(root->l, s, sgn);  
4         hash_collect(root->r, s, -sgn);  
5     } else {  
6         s[root->id] = 'b' + sgn;  
7     }  
8 }
```

Funkcija je testirana sa $M = 2^{61} - 1$, što je prost broj, $n = 4000$ i za svako q počev od $q_0 = 314159265$ do $q_0 + 999$. Funkcija je našla koliziju u 801 od 1000 scenarija, i ukupno vreme celokupnog izvršenja je bilo oko 0.54 sekunde, odnosno, u proseku oko 0.54ms po scenariju.

Nije jednostavno proceniti vrednost broja n kod koje je verovatnoća uspeha približno $\frac{1}{2}$. Jedna procena¹⁵ je da se radi o vrednosti $n = 2^{\sqrt{2\log_2(M)+1}}$. Za $M = 2^{61} - 1$, ova procena daje vrednost $n \approx 4227$, što se uklapa u empirijski dobijeni rezultat. Vremenska složenost algoritma je $O(n \log n)$, jer dominira sortiranje u prvoj iteraciji, u svakoj narednoj se broj elemenata prepolovi. Sudeći po ovoj proceni, ovaj algoritam je izvodljiv i za vrednosti M reda veličine 2^{256} odnosno 10^{77} , gde je $n \approx 10^7$. Ovaj algoritam se može primeniti i kada M nije prost broj, a i kada q, M nisu uzajamno prosti.

7 Palindromi

U ovom odeljku pokazaćemo da skup svih palindromskih podstringova stringa s ima zanimljivu strukturu, i daćemo jednostavne i efikasne algoritme koji nalaze tu strukturu.

Teorema 7.1 *String s dužine n ima ne više od n različitih nepraznih palindromskih podstringova.*

Dokaz. Indukcijom po dužini stringa n . Za $n = 1$ tvđenje očigledno važi, jer string ima samo jedan palindromski podstring – sebe samog. Neka je $n > 1$. Posmatrajmo sve palindromske sufikse stringa s , neka najduži od njih ima dužinu k . Svaki palindromski sufiks dužine $l < k$ stringa s se javlja i u podstringu $s_{[0, n-1]}$, zato što je $s_{[n-k, n]}$ palindrom, pa je $s_{[n-l, n]} = \overline{s_{[n-k, n-k+l]}} = s_{[n-k, n-k+l]}$, što je podstring jer važi $0 \leq n-k \leq n-k+l \leq n-1$. Dakle, skup palindromskih podstringova za string s ima najviše jedan element više nego taj skup za string $s_{[0, n-1]}$, i to može biti samo podstring $s_{[n-k, n]}$. \square

Jednostavan algoritam koji nalazi sve palindrome je sledeći. Fiksirajmo centar palindroma. Centar može biti jedno slovo ili pozicija između dva uzastopna slova, što odgovara redom palindromima neparne i parne dužine. Sve dok se prvo i poslednje slovo poklapaju, povećavamo granice. Oslanjamo se na ideju da, ako je s palindrom dužine bar 2, tada je i $s_{[1, n-1]}$ takođe palindrom. Algoritam zato vraća samo maksimalne palindrome, tj. palindrome koji se ne mogu proširiti i sa leve i sa desne strane istovremeno.

Algoritam za nalaženje svih maksimalnih palindromskih podstringova

```
1 vector<pair<int, int>> palin_n2(const string& s) {
2     int n = s.size();
3     vector<pair<int, int>> result;
4     for (int l=0, r=0; r<n; (l<r?l:r)++) {
5         int x=l, y=r;
6         while (x>=0 && y<n && s[x] == s[y])
7             x--, y++;
8         if (x+1 < y)
9             result.emplace_back(x+1, y);
10    }
11    return result;
12 }
```

Glavni nedostatak algoritma je njegova vremenska složenost, koja u najgorem slučaju iznosi $\Theta(n^2)$, i postiže se za string $s = a^n$.

7.1 Manacher-ov algoritam

Pre nego što opišemo sam algoritam, da bismo ga pojednostavili svedimo problem nalaženja maksimalnih palindroma na problem nalaženja maksimalnih palindroma neparne dužine. Ovo radimo tako što od stringa s dužine n pravimo string p dužine $2n + 1$, $p = \$s_0\$s_1\$ \dots \$s_{n-1}\$$, gde je $\$$ karakter koji se ne javlja u stringu s . String p ne sadrži parne palindrome, jer se srednja dva karaktera sigurno ne poklapaju – jedan od njih će uvek biti $\$$, a drugi neće. Sa druge strane, palindrom sa centrom u karakteru $\$$, osim na krajevima, odgovara palindromu parne dužine u stringu s , dok palindrom sa centrom u karakteru različitom od $\$$ odgovara palindromu neparne dužine u stringu s . Konkretno, palindrom $p_{[l,r]}$ odgovara palindromu $s_{[\frac{l}{2}, \frac{r-1}{2}]}$. Uzimamo u razmatranje samo neparne palindrome u p koji počinju i završavaju se sa $\$$ i imaju dužinu bar 3.

Manacher-ov algoritam koristi ideju vrlo sličnu Z-algoritmu. Za svaku poziciju stringa p dužine m algoritam nalazi radijus q_i , koji znači da maksimalni palindrom sa centrom u i ima dužinu $2q_i + 1$. Algoritam redom računa vrednosti q_i i održava prozor sa najvećim desnim krajem koji odgovara nekom maksimalnom palindromu. Ovaj prozor se onda može iskoristiti da se nađe donja granica za vrednost q_i . Ukoliko je $s_{[l,r]}$ palindrom (obratiti pažnju da sada obuhvatamo i desni kraj intervala), onda za $i > \frac{r+l}{2}$ važi $q_i \geq q'_i = \min(q_{r+l-i}, r - i)$ – pošto je $s_{[l,r]}$ palindrom, onda je podstring $s_{[i-q'_i, i+q'_i]}$ ceo sadržan u $s_{[l,r]}$ i jednak podstringu $\overline{s_{[r+l-i-q'_i, r+l-i+q'_i]}}$, koji je palindrom, odakle sledi da je $q_i \geq q'_i$.

Manacher-ov algoritam

```
1 vector<pair<int, int>> manacher(const string& s) {
2     int n = s.size(), m = 2*n+1;
3     string p = "$";
4     for (int i=0; i<n; i++)
5         p += s[i], p += '$';
6     vector<int> q(m, 0);
7     vector<pair<int, int>> result;
8     for (int i=1, l=0, r=0; i<m-1; i++) {
9         if (i < r)
10            q[i] = min(q[r+l-i], r-i);
11        while (q[i] <= i && i+q[i] < m && p[i-q[i]] == p[i+q[i]])
12            q[i]++;
13        q[i]--;
14        if (i+q[i] > r) {
15            r = i + q[i];
16            l = i - q[i];
17        }
18        if (q[i] > 0)
19            result.emplace_back((i-q[i])/2, (i+q[i])/2);
20    }
21    return result;
22 }
```

Dokaz da ovaj algoritam ima vremensku složenost $O(n)$ je sličan dokazu kod Z-algoritma. Dokažimo da se vrednost r poveća za bar jedan u svakoj iteraciji *while* petlje. U slučaju da je $i + q'_i < r$, *while* petlja će izvršiti tačno 0 iteracija pošto je $q_i = q'_i$. U suprotnom bismo dobili, na osnovu refleksije unutar palindroma $s_{[l,r]}$ da se palindrom sa centrom u $r + l - i$ takođe može proširiti za bar jedno slovo, što je nemoguće jer je njegova vrednost tačno izračunata. Drugi slučaj je kad je $i + q'_i \geq r$, tada će, u slučaju da se palindrom sa centrom u i proširi sa k iteracija *while* petlje, za isto toliko povećati vrednost r .

7.2 Palindromsko stablo

Palindromsko stablo, poznato i pod imenom *eertree*, je struktura podataka koja za dati string s nalazi sve palindrome u stringu i omogućuje njihovo efikasno pretraživanje i obradu. Svaki čvor palindromskog stabla je jedan jedinstven neprazan palindrom. Pored ovih, postoje dva specijalna čvora koji odgovaraju palindromu dužine 0, tj. stringu ϵ , kao i fiktivnom palindromu

η dužine -1 . Ovaj palindrom zadovoljava sledeću osobinu: ako je $x \in \Sigma$, onda je $x\eta x = x$. Svaki čvor ima izlazne grane sa labelama iz Σ . Ako grana polazi iz čvora p i ima labelu x , onda ta grana ide ka čvoru xpx . Pored ovih grana, čvorovi imaju i sufiks grane – sufiks grana iz čvora p ide ka najdužem palindromu koji je pravi sufiks palindroma p , ukoliko takav ne postoji, sufiks grana ide ka ϵ . Kod ovog čvora sufiks grana ide ka čvoru η , dok se za čvor η sufiks grana ne definiše. Očigledno, sufiks grane formiraju strukturu obrnutog korenskog stabla sa korenom u η , dok izlazne grane formiraju dva korenska stabla, jedno sa korenom u η , a drugo sa korenom u ϵ , jer svi ostali čvorovi imaju jedinstvenu ulaznu granu, i to je grana sa labelom koja odgovara njihovom prvom tj. poslednjem slovu.

7.2.1 Algoritam za konstrukciju

Na osnovu teoreme 7.1 imamo da je broj čvorova palindromskog stabla ne više od $n + 2$, pa je memorijska složenost cele strukture $O(n)$. Pored navedenih grana, svaki čvor će eksplicitno pamtit i dužinu palindroma koji predstavlja.

Struktura čvora palindromskog stabla

```

1 struct eertree_node {
2     int len;
3     eertree_node* link;
4     map<char, eertree_node*> next;
5 };

```

Glavni algoritam će, pored traženog stabla, naći i čvor koji odgovara najdužem sufiksu koji je palindrom za svaki prefiks stringa s , odnosno, čvor p_i odgovara najdužem sufiksu koji je palindrom stringa $s_{[0,i]}$. Sve ove vrednosti, kao i samo palindromsko stablo računamo redom za sve prefikse stringa s . Iz dokaza teoreme 7.1 znamo da je ovaj skup palindroma jednak skupu svih palindroma koji se javljaju u s . Inicijalizujemo algoritam kreiranjem dva čvora ϵ, η i postavljanjem $p_0 = \epsilon$. Zatim, svaki put kad dodajemo novo slovo stringa s , prvo ispitujemo da li se prethodno najduži sufiks palindrom može proširiti novim slovom c . Ukoliko može, onda je novi palindrom upravo taj proširen slovom c , inače, ispitujemo dalje krećući se po sufiks vezama. Kada nađemo prvi sufiks t koji se može proširiti sa obe strane karakterom c , zaustavljamo se. Ukoliko čvor ctc već postoji, ne radimo ništa, inače mu moramo izračunati sufiks vezu. Ukoliko je $t = \eta$, novi palindrom ima samo

jedno slovo, odnosno c se javlja prvi put u stringu i njegova sufiks veza ide ka čvoru ϵ , inače, krećemo od sufiks veze čvora t i tražimo najduži sufiks koji se može proširiti slovom c . Takav će sigurno postojati jer se slovo c ne javlja prvi put u stringu, pa će bar η moći da se produži. Ako je q takav palindrom, onda je sufiks veza za upravo cqc .

Algoritam za konstrukciju palindromskog stabla

```

1  using node = eertree_node;
2  vector<node*> eertree(const string& s) {
3      int n = s.size();
4      node* eta = new node {-1, nullptr, {}};
5      node* eps = new node {0, eta, {}};
6      vector<node*> p(n+1);
7      p[0] = eps;
8      for (int i=1; i<=n; i++) {
9          char c = s[i-1];
10         node* t = p[i-1];
11         while (i-t->len < 2 || s[i-t->len-2] != c)
12             t = t->link;
13         if (t->next.count(c)) {
14             t = t->next[c];
15         } else {
16             node* q = t->link;
17             t = t->next[c] = new node {t->len+2, nullptr, {}};
18             if (!q) {
19                 t->link = eps;
20             } else {
21                 while (s[i-q->len-2] != c)
22                     q = q->link;
23                 t->link = q->next[c];
24             }
25         }
26         p[i] = t;
27     }
28     return p;
29 }

```

Vremenska složenost algoritma je $O(n)$. To se može dokazati ako posmatramo vrednost $5i - \text{len}(t) - \text{len}(\text{link}(t))$. U svakoj iteraciji neke od *while* ili *for* petlji se ova vrednost povećava. Za *while* petlje je to očigledno. Za *for* petlju primetimo da, ukoliko je najduži sufiks palindrom nekog stringa imao dužinu l , ako se doda jedno slovo na kraj stringa, ta dužina može postati najviše $l+2$, i isto važi za dužinu drugog najdužeg sufiks palindroma (tj. sufiks veze), pa se $\text{len}(t) - \text{len}(\text{link}(t))$ povećava za najviše 4, a $5i$ se povećava za

tačno 5. Kako je početna vrednost izraza 1, a na kraju ne prelazi $5n$, ukupan broj izvršenja svih petlji je $5n$, pa algoritam ima vremensku složenost $O(n)$.

7.2.2 Primene

Najduži palindromski podstring

Palindromsko stablo rešava ovaj problem u istoj vremenskoj složenosti kao Manacher-ov algoritam, odnosno u $O(n)$. Dovoljno je za sve čvorove stabla uzeti vrednost $len(u)$ i vratiti onaj koji ima najveću. Pozicija palindroma u stringu se može naći ispitivanjem za koje i je $p_i = u$. U tom slučaju se palindrom nalazi na poziciji $s_{[i-len(u), i]}$.

Ukupan broj palindroma

Nađimo za svaki čvor stabla t sufiks veza dubinu $d(t)$, odnosno broj čvorova na putu od tog čvora do nekog čvora sa $len = 1$. Ovo će nam za palindrom t dati broj sufiksa stringa t koji su palindromi. Rešenje za ceo string s je $\sum_{i=1}^n d(p_i)$.

Prethodni problemi se mogu efikasno rešiti i Manacher-ovim algoritmom. Međutim, on nije dovoljan za naredne.

Broj različitih palindroma

Po definiciji palindromskog stabla, broj različitih palindroma pozitivne dužine je broj čvorova stabla umanjeno za 2, jer ne računamo čvorove ϵ i η .

Ispitivanje pojavljivanja palindroma

Ukoliko ispitujemo da li se string parne dužine oblika $\bar{u}u$ javlja u s , možemo samo pratiti izlazne grane sa labelama koje odgovaraju slovima stringa u , počev od čvora ϵ . Za neparne palindrome oblika $\bar{u}xu$, krećemo iz η i krećemo se labelama xu .

Literatura

- [1] Knuth D.E, Morris J.H, Pratt V.R. *Fast Pattern Matching in Strings*. SIAM Journal on Computing, 1977, Vol. 6, No. 2 : pp. 323-350
- [2] Gusfield D. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, 1997
- [3] Hirschberg D.S. *A linear space algorithm for computing maximal common subsequences*. Comm. A.C.M. 18(6) p341-343, 1975
- [4] Jurafsky D, Martin J.H. *Speech and Language Processing*. Pearson Education International, 2000
- [5] Aho A.V, Corasick M.J. *Efficient string matching: an aid to bibliographic search*. Comm. A.C.M. 18(6) p333-340, 1975
- [6] Bender M.A., Farach-Colton M. *The LCA Problem Revisited*. Gonnet G.H., Viola A. (eds) LATIN 2000: Theoretical Informatics. LATIN 2000. Lecture Notes in Computer Science, vol 1776. Springer, Berlin, Heidelberg
- [7] Manber U, Myers G. *Suffix arrays: a new method for on-line string searches*. SIAM J Comput. 1993;22(5):935–48.
- [8] Kasai T, Lee G, Arimura H, Arikawa S, Park K. *Linear-Time Longest-Common-Prefix Computation in Suffix Arrays and Its Applications*. Amir A. (eds) Combinatorial Pattern Matching. CPM 2001. Lecture Notes in Computer Science, vol 2089. Springer, Berlin, Heidelberg
- [9] Ukkonen E. *On-line construction of suffix trees*. Algorithmica (1995) 14: 249.
- [10] Blumer A, Blumer J, Haussler D, Ehrenfeucht A, Chen M.T, Seiferas J. *The smallest automation recognizing the subwords of a text*. Theoretical Computer Science Volume 40, 31-55, Elsevier (1985)
- [11] Allouche J-P., Shallit J. *The Ubiquitous Prouhet-Thue-Morse Sequence*. Ding C., Helleseth T., Niederreiter H. (eds) Sequences and their Applications. Discrete Mathematics and Theoretical Computer Science. Springer, London

- [12] Flajolet P, Grabner P.J, Kirschenhofer P, Prodinger H. *On Ramanujan's Q-function*. JCAM, Volume 58, Issue 1, 103-116 (1995)
- [13] Flaxman A.D, Przydatek B. *Solving Medium-Density Subset Sum Problems in Expected Polynomial Time*. STACS 2005. Lecture Notes in Computer Science, vol 3404. Springer, Berlin, Heidelberg
- [14] Manacher G. *A New Linear-Time On-Line Algorithm for Finding the Smallest Initial Palindrome of a String*. JACM, Volume 22, Issue 3, 346-351 (1975)
- [15] <https://codeforces.com/blog/entry/60442>
- [16] <https://en.cppreference.com/w/cpp/algorithm/sort>
- [17] <http://adilet.org/blog/palindromic-tree/>