

# Wyszukiwanie wzorca w tekście

Dane:

- $T[1 \dots n]$  – tekst długości  $n$
- $P[1 \dots m]$  – wzorzec długości  $m$  ( $m < n$ ).

$P$  występuje w  $T$  z przesunięciem  $s$ , jeśli  $T[s + 1 \dots s + m] = P[1 \dots m]$ . (Takie  $s$  jest poprawnym przesunięciem).

Oznaczenia:

$\Sigma^*$  – zbiór wszystkich słów nad alfabetem  $\Sigma$  (łącznie ze słowem pustym  $\epsilon$ ).

$|x|$  – długość słowa  $x$

$xy$  – *konkatenacja* (sklejenie) słów  $x$  i  $y$

$w \sqsubset x$  –  $w$  jest *prefiksem*  $x$  (tzn.  $\exists_{y \in \Sigma^*} x = wy$ )

$w \sqsupset x$  –  $w$  jest *sufiksem*  $x$  (tzn.  $\exists_{y \in \Sigma^*} x = yw$ )

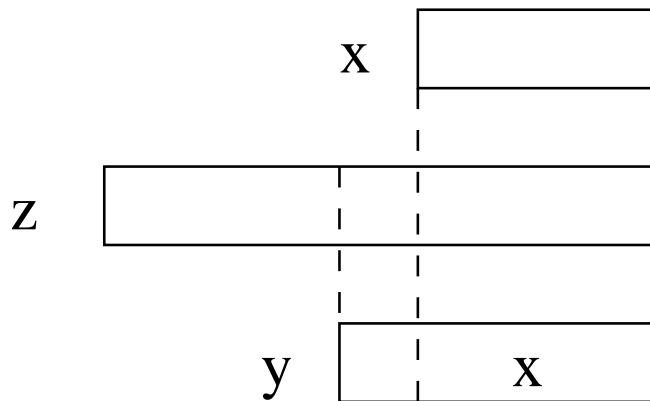
# Lemat o nierozłącznych sufiksach

**Lemat 1. (o nierozłącznych sufiksach)** Niech  $x, y, z \in \Sigma^*$ , takie że  $x \sqsubset z$  i  $y \sqsubset z$ . Wtedy:

- jeśli  $|x| \leq |y|$ , to  $x \sqsubset y$
- jeśli  $|y| \leq |x|$ , to  $y \sqsubset x$
- jeśli  $|y| = |x|$ , to  $y = x$

**D-d.**

Przypadek  $|x| \leq |y|$ :



Pozostałe przypadki – podobnie.  $\square$

# Algorytm naiwny

Koszt sprawdzania warunku:  $P[1 \dots m] = T[s + 1 \dots s + m]$  jest  $\Theta(m)$ .

*Naive-String-Matcher( $T, P$ )*

```
1  $n \leftarrow \text{length}[T]$ 
2  $m \leftarrow \text{length}[P]$ 
3 for  $s \leftarrow 0$  to  $n - m$  do
4   if  $P[1 \dots m] = T[s + 1 \dots s + m]$  then
5     drukuj "występuje z przesunięciem  $s$ "
```

Czas:  $\Theta((n - m + 1) \cdot m)$

( $n - m + 1$  iteracji for,  $\Theta(m)$  – koszt sprawdzenia w if)

Jeśli  $m$  i  $n - m + 1$  są  $\Omega(n)$ , to czas:  $\Theta(n^2)$ .

# Algorytm Karpa – Rabina

Niech  $d = |\Sigma|$ .  $\Sigma$  – cyfry w systemie o podstawie  $d$ . Wzorzec i badany fragment tekstu –  $m$ -cyfrowe liczby w systemie o podstawie  $d$ .

*Przykład.*  $\Sigma = \{0, 1, \dots, 9\}$ . Liczbę  $p$  odpowiadającą wzorcowi  $P[1 \dots m]$  można wyznaczyć w czasie  $O(m)$  wg reguły Hornera:

$$p = P[m] + 10(P[m-1] + 10(P[m-2] + \dots + 10(P[2] + 10P[1]) \dots))$$

Podobnie wyznaczamy  $t_0$  – liczbę reprezentowaną przez  $T[1 \dots m]$ . Dla  $s \geq 0$  wyznaczamy  $t_{s+1}$  reprezentowaną przez  $T[s+2 \dots s+m+1]$  wyznaczamy z  $t_s$ :

$$t_{s+1} = 10(t_s - 10^{m-1}T[s+1]) + T[s+m+1]$$

(Usuujemy cyfrę  $T[s+1]$ , przesuwamy pozostałe w lewo i dodajemy  $T[s+m+1]$ .)

# Algorytm Karpa – Rabina

*Kłopot:*  $p$  i  $t_s$  – DUŻE liczby

*Rozwiązanie:* Obliczenia – modulo  $q$ , gdzie  $q$  liczba (zwykle pierwsza), taka że  $d \cdot q$  mieści się w słowie komputera.

Niech  $h = d^{m-1} \bmod q$ .

$h$  odpowiada liczbie:  $(100 \dots 0)_d$  ( $m - 1$  zer) modulo  $q$ .

Wtedy:

$$t_{s+1} = (d(t_s - T[s+1] \cdot h) + T[s+m+1]) \bmod q$$

Jeśli  $t_s \not\equiv p \pmod{q}$ , to  $T[s+1 \dots s+m] \neq P[1 \dots m]$  (nie musimy sprawdzać!).

Rabin-Karp-Matcher( $T, P, d, q$ )

1  $n \leftarrow \text{length}[T]$

2  $m \leftarrow \text{length}[P]$

...

# Rabin-Karp-Matcher

```
3   $h \leftarrow d^{m-1} \bmod q$  ▷ czas:  $O(m)$ 
4   $p \leftarrow 0$ 
5   $t_0 \leftarrow 0$ 
6  for  $i \leftarrow 1$  to  $m$  do
7       $p \leftarrow (d \cdot p + P[i]) \bmod q$ 
8       $t_0 \leftarrow (d \cdot t_0 + T[i]) \bmod q$ 
9  for  $s \leftarrow 0$  to  $n - m$  do
10     if  $p = t_s$  then
11         if  $P[1 \dots m] = T[s + 1 \dots s + m]$  then
12             drukuj "występuje z przesunięciem  $s$ "
13     if  $s < n - m$  then
14          $t_{s+1} \leftarrow (d(t_s - T[s + 1] \cdot h) + T[s + m + 1]) \bmod q$ 
```

**Czas:**  $O(m) + (\text{liczba błędnych strzałów}) \cdot O(m) +$   
 $(\text{liczba wystąpień } P \text{ w } T) \cdot O(m) + O(n - m + 1)$

# Algorytm Karpa – Rabina

Ponieważ co  $q$ -ta liczba jest  $\equiv p \pmod{q}$ , więc można oczekiwać (przy losowym wyborze  $q$ ), że liczba błędnych strzałów jest  $O(n/q)$ .

Niech  $v$  – liczba wystąpień  $P$  w  $T$ .

Oczekiwany czas:  $O(n) + O(m(v + n/q))$ .

Jeśli  $v$  jest  $O(1)$  a  $q > m$ , to oczekiwany czas :  $O(n + m)$ .

# Automaty skończone

*Automat skończony*  $M = (Q, q_0, A, \Sigma, \delta)$

- $Q$  – skończony zbiór *stanów*
- $q_0 \in Q$  – stan *początkowy*
- $A \subseteq Q$  – zbiór stanów *akceptujących*
- $\delta : Q \times \Sigma \rightarrow Q$  – *funkcja przejść*

*Rozszerzona funkcja przejść*:  $\phi : \Sigma^* \rightarrow Q$  taka, że

$$\phi(\epsilon) = q_0$$

$$\phi(wa) = \delta(\phi(w), a), \quad \text{dla } w \in \Sigma^* \text{ i } a \in \Sigma$$

$(\phi(w)$  – *stan automatu (startującego z  $q_0$ ) po wczytaniu słowa  $w$ .)*

$M$  *akceptuje  $w$ , jeśli  $\phi(w) \in A$ .*



# Automat wyszukiwania wzorca $P$

Niech  $P_k = P[1 \dots k]$  (prefiks  $P$  długości  $k$ ),  $T_i = T[1 \dots i]$  (prefiks  $T$  długości  $i$ ).

funkcja *sufiksowa* dla  $P$  – funkcja  $\sigma : \Sigma^* \rightarrow \{0, \dots, m\}$  taka że:

$$\sigma(x) = \max\{k : P_k \sqsupseteq x\}$$

(Długość najdłuższego prefiksu  $P$ , który jest sufiksem  $x$ )

Automat dla  $P$ :

•  $Q = \{0, \dots, m\}$ ,  $q_0 = 0$ ,  $A = \{m\}$

•  $\delta(q, a) = \sigma(P_q a)$

# Finite-Automaton-Matcher

Finite-Automaton-Matcher( $T, \delta, m$ )

```
1   $n \leftarrow \text{length}[T]$ 
2   $q \leftarrow 0$ 
3  for  $i \leftarrow 1$  to  $n$  do
4       $q \leftarrow \delta(q, T[i])$ 
5      if  $q = m$  then
6           $s \leftarrow i - m$ 
7      drukuj  "występuje z przesunięciem  $s$ "
```

Czas:  $O(n)$ .

# analiza poprawności

**Lemat 2.(Nierówność dla funkcji sufiksowej)** Dla  $x \in \Sigma^*$  i  $a \in \Sigma$  zachodzi  $\sigma(xa) \leq \sigma(x) + 1$ .

**D-d.** Niech  $r = \sigma(xa)$ .

Jeśli  $r = 0$ , to  $r \leq \sigma(x) + 1$ , bo  $\sigma(x) \geq 0$ .

Niech  $r > 0$ . Wtedy  $P_r \sqsubset xa$ . Stąd  $P_{r-1} \sqsubset x$ . Zatem  $r - 1 \leq \sigma(x)$  (z definicji  $\sigma$ ).  $\square$

# analiza poprawności

**Lemat 2.(Nierówność dla funkcji sufiksowej)** Dla  $x \in \Sigma^*$  i  $a \in \Sigma$  zachodzi  $\sigma(xa) \leq \sigma(x) + 1$ .

**D-d.** Niech  $r = \sigma(xa)$ .

Jeśli  $r = 0$ , to  $r \leq \sigma(x) + 1$ , bo  $\sigma(x) \geq 0$ .

Niech  $r > 0$ . Wtedy  $P_r \sqsubset xa$ . Stąd  $P_{r-1} \sqsubset x$ . Zatem  $r - 1 \leq \sigma(x)$  (z definicji  $\sigma$ ).  $\square$

**Lemat 3.(Rekurencja dla funkcji sufiksowej)** Dla  $x \in \Sigma^*$  i  $a \in \Sigma$ , jeśli  $q = \sigma(x)$ , to  $\sigma(xa) = \sigma(P_q a)$ .

**D-d.** Z definicji  $\sigma$ :  $P_q \sqsubset x$ . Stąd  $P_q a \sqsubset xa$ . Niech  $r = \sigma(xa)$ . Z Lematu 2:  $r \leq q + 1$  (czyli:  $|P_r| \leq |P_q a|$ ). Z Lematu 1:  $P_r \sqsubset P_q a$ , ( bo  $P_q a \sqsubset xa$ ,  $P_r \sqsubset xa$  i  $|P_r| \leq |P_q a|$ ). Stąd:  $r \leq \sigma(P_q a)$  (czyli:  $\sigma(xa) \leq \sigma(P_q a)$ ).

Ponieważ  $P_q a \sqsubset xa$ , zachodzi również  $\sigma(P_q a) \leq \sigma(xa)$ .

Zatem  $\sigma(P_q a) = \sigma(xa)$   $\square$

# analiza poprawności

**Tw 4.** Jeśli  $\phi$  jest rozszerzoną funkcją przejść automatu wyszukiwania wzorca  $P$ , a  $T[1 \dots n]$  – tekstem wejściowym automatu, to  $\phi(T_i) = \sigma(T_i)$  dla  $i = 0, \dots, n$ .

**D-d.** Indukcja po  $i$ . Dla  $i = 0$ :  $T_0 = \epsilon$  i  $\phi(\epsilon) = \sigma(\epsilon) = 0$ .

Założmy, że  $\phi(T_i) = \sigma(T_i)$  (założenie indukcyjne). Niech  $q = \phi(T_i)$  oraz  $a = T[i + 1]$ .

$$\begin{aligned}\phi(T_{i+1}) &= \phi(T_i a) && \text{(z def. } T_{i+1} \text{ i } a) \\ &= \delta(\phi(T_i), a) && \text{(z def. } \phi) \\ &= \delta(q, a) && \text{(z def. } q) \\ &= \sigma(P_q a) && \text{(z def. } \delta) \\ &= \sigma(T_i a) && \text{(z Lematu 3 i zał. ind.)} \\ &= \sigma(T_{i+1}) && \text{(z def. } T_{i+1})\end{aligned}$$



# wyznaczenie $\delta$

Compute-Transition-Function( $P, \Sigma$ )

```
1  $m \leftarrow \text{length}[P]$ 
2 for  $q \leftarrow 0$  to  $m$  do
3   for  $a \in \Sigma$  do
4      $k \leftarrow \min\{m + 1, q + 2\}$ 
5     repeat  $k \leftarrow k - 1$ 
6       until  $P_k \sqsupseteq P_q a$ 
7      $\delta(q, a) \leftarrow k$ 
8 return  $\delta$ 
```

Poprawność: bezpośrednio z definicji  $\sigma$  i  $\delta$ .

for z wiersza 2 –  $O(m)$  razy

for z wiersza 3 –  $O(|\Sigma|)$  razy

repeat –  $O(m)$  razy

Sprawdzenie warunku w `until` – czas  $O(m)$ .

Łączny czas:  $O(m^3|\Sigma|)$ .

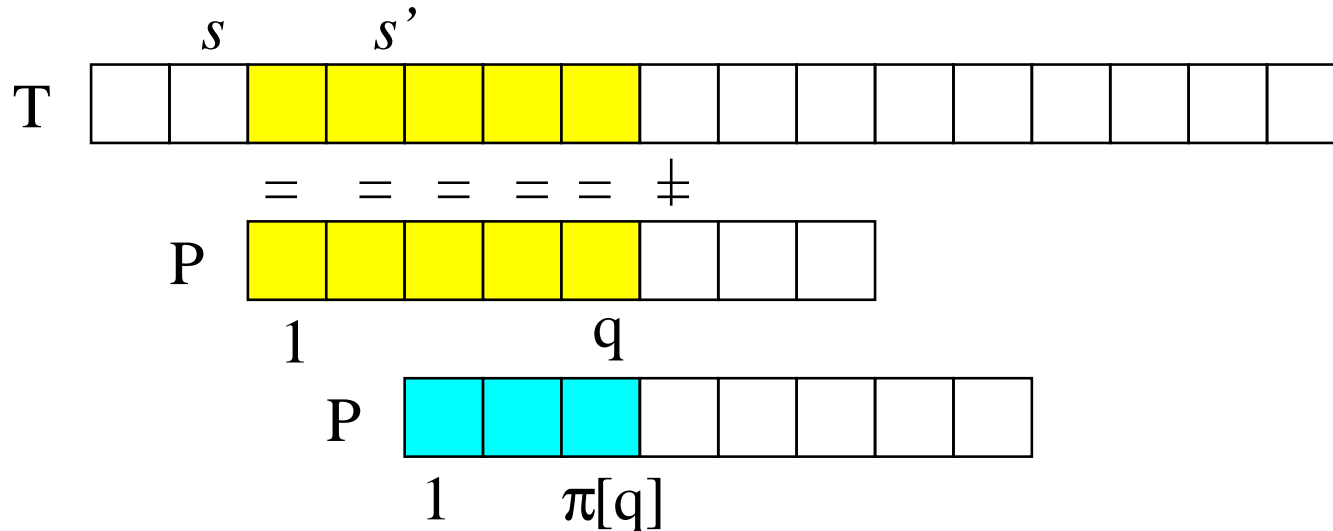
# Algorytm Knutha-Morrisa-Pratta

Funkcja prefiksowa dla  $P$  – funkcja  $\pi : \{1 \dots m\} \rightarrow \{0 \dots m - 1\}$

$$\pi[q] = \max\{k : k < q \wedge P_k \sqsupseteq P_q\}$$

Jeśli  $T[s + 1 \dots s + q] = P[1 \dots q]$  i  $T[s + q + 1] \neq P[q + 1]$ , to następne przesunięcie, które warto sprawdzać:

$$s' = s + q - \pi[q].$$



# KMP-Matcher

KMP-Matcher( $T, P$ )

```
1   $n \leftarrow \text{length}[T]$ 
2   $m \leftarrow \text{length}[P]$ 
3   $\pi \leftarrow \text{Compute-Prefix-Function}(P)$ 
4   $q \leftarrow 0$ 
5  for  $i \leftarrow 1$  to  $n$  do
6      while  $q > 0$  and  $P[q + 1] \neq T[i]$  do
7           $q \leftarrow \pi[q] \triangleright P[1 \dots q] = T[i - q \dots i - 1]$ 
8      if  $P[q + 1] = T[i]$  then
9           $q \leftarrow q + 1$ 
10     if  $q = m$  then
11         drukuj  "występuje z przesunięciem  $i - m$ "
12          $q \leftarrow \pi[q]$ 
```



# Compute-Prefix-Function

Compute-Prefix-Function( $P$ )

```
1   $m \leftarrow \text{length}[P]$ 
2   $\pi[1] \leftarrow 0$ 
3   $k \leftarrow 0$ 
4  for  $q \leftarrow 2$  to  $m$  do
5      while  $k > 0$  and  $P[k + 1] \neq P[q]$  do
6           $k \leftarrow \pi[k] \triangleright P[1 \dots k] = P[q - k \dots q - 1]$ 
7      if  $P[k + 1] = P[q]$  then
8           $k \leftarrow k + 1$ 
9       $\pi[q] \leftarrow k$ 
10 return  $\pi$ 
```

# analiza czasu działania

Czas `Compute-Prefix-Function` – Potencjał:  $k$ .

Początkowo  $k = 0$  (po wierszu 3) oraz zawsze (po każdym wierszu 6 i 8)  $k \geq 0$  bo  $\pi[k] \geq 0$  dla każdego  $k$  (to wynika z poprawności `Compute-Prefix-Function` co pokażemy później).

Każde wykonanie wiersza 6 – zmniejszenie  $k$ , bo  $\pi[k] < k$  (to też wynika z poprawności).

Koszt każdego wykonania pętli `while` – zamortyzowany zmniejszeniem potencjału w wierszu 6.

W wierszu 8 – zwiększamy potencjał o  $\leq 1$ .

Pętla `for` –  $m - 1$  iteracji.

Stąd koszt faktyczny  $\leq$  koszt zamortyzowany  $= O(m)$ .

Koszt `KMP-Matcher`:  $O(m + n)$

(`Compute-Prefix-Function` + podobna analiza z potencjałem:  $q$ ).

# Lemat 5

Niech  $\pi^*[q] = \{\pi[q], \pi^{(2)}[q] \dots, \pi^{(t)}[q]\}$ , gdzie  $\pi^{(0)}[q] = q$ ,  
 $\pi^{(i+1)}[q] = \pi[\pi^{(i)}[q]]$ , oraz  $t$  – najmniejsze, takie że  $\pi^{(t)}[q] = 0$ .

**Lemat 5 (o iteracji funkcji prefiksowej)** Dla  $q = 1, \dots, m$  mamy  
 $\pi^*[q] = \{k : k < q \wedge P_k \sqsupset P_q\}$ .

**D-d. Fakt:**  $i \in \pi^*[q]$  implikuje  $i < q \wedge P_i \sqsupset P_q$ .

*Uzasadnienie Faktu:* Jeśli  $i \in \pi^*[q]$ , to  $\exists_{0 < u \leq t} i = \pi^{(u)}[q]$ .

Indukcja po  $u$ . Dla  $u = 1$ :  $i = \pi[q]$ . Stąd wynika teza, bo  
 $\pi[q] < q \wedge P_{\pi[q]} \sqsupset P_q$ .

Krok indukcyjny: Korzystając z  $P_{\pi[j]} \sqsupset P_j$  i przechodności  $<$   
oraz  $\sqsupset$  otrzymujemy tezę dla wszystkich  $u$ .

Z Faktu wynika  $\pi^*[q] \subseteq \{k : k < q \wedge P_k \sqsupset P_q\}$  (\*).

...

## c.d. dowodu Lematu 5

...

Udowodnimy nie wprost, że  $\{k : k < q \wedge P_k \sqsupset P_q\} \subseteq \pi^*[q]$ .

Założmy, że  $A = \{k : k < q \wedge P_k \sqsupset P_q\} \setminus \pi^*[q] \neq \emptyset$ . Niech

$j = \max(A)$ . Ponieważ

$\pi[q] = \max\{k : k < q \wedge P_k \sqsupset P_q\} \wedge \pi[q] \in \pi^*[q]$ , mamy  $j < \pi[q]$ .

Niech  $j' = \min\{l \in \pi^*[q] : l > j\}$ . Mamy  $P_j \sqsupset P_q$ , bo

$j \in \{k : P_k \sqsupset P_q\}$ , oraz  $P_{j'} \sqsupset P_q$ , bo  $j' \in \pi^*[q]$  i (\*).

Zatem z Lematu 1  $P_j \sqsupset P_{j'}$ . Co więcej  $j = \max\{k : k < j' \wedge P_k \sqsupset P_{j'}\}$ ,

czyli  $j = \pi[j']$ , a ponieważ  $j' \in \pi^*[q]$ , więc  $j = \pi[j'] \in \pi^*[q]$ .

Sprzeczność.  $\square$

# Definicja $E_{q-1}$

Dla  $q = 2, \dots, m$ , zdefiniujmy  $E_{q-1} \subseteq \pi^*[q-1]$  jako

$$E_{q-1} = \{k \geq 0 : k \in \pi^*[q-1] \wedge P[k+1] = P[q]\}$$

**z Lematu 5:**

$$= \{k \geq 0 : k < q-1 \wedge P_k \sqsubset P_{q-1} \wedge P[k+1] = P[q]\}$$

**z definicji  $\sqsubset$ :**

$$= \{k \geq 0 : k < q-1 \wedge P_{k+1} \sqsubset P_q\}$$

# Wniosek 7

**Wn. 7.** Niech  $\pi$  – f-cja prefiksowa dla  $P$ . Wtedy, dla  $q = 2, \dots, m$  zachodzi:

$$\pi[q] = \begin{cases} 0, & \text{jeśli } E_{q-1} = \emptyset \\ 1 + \max(E_{q-1}), & \text{jeśli } E_{q-1} \neq \emptyset \end{cases}$$

**D-d.** Jeśli  $E_{q-1} = \emptyset$  to nie ma  $0 \leq k < q - 1$ , dla którego  $P_{k+1} \sqsupset P_q$ . Stąd  $\pi[q] = 0$  (tylko  $\epsilon = P_0 \sqsupset P_q$ ).

Niech  $E_{q-1} \neq \emptyset$ . Z definicji  $\pi[q]$

$$\begin{aligned} \pi[q] &= \max\{k + 1 \geq 0 : k + 1 < q \wedge P_{k+1} \sqsupset P_q\} \\ &= \max\{k + 1 \geq 0 : k < q - 1 \wedge P_{k+1} \sqsupset P_q\} \\ &= 1 + \max(E_{q-1}) \text{ (bo } E_{q-1} \neq \emptyset). \quad \square \end{aligned}$$

# poprawność Compute-Prefix-Function

Przypisanie w wierszu 2 – poprawne, bo  $\pi[1] = 0$ .

Pokażemy, że wartości  $\pi[q]$  dla  $q > 1$  też są poprawnie wyznaczone.

Na początku każdej iteracji `for` zachodzi:  $k = \pi[q - 1]$  (wymuszone w wierszach 2, 3 oraz w wierszu 9) i **wyznaczone są poprawne wartości  $\pi[j]$ , dla  $j = 1, \dots, q - 1$**  (założenie indukcyjne).

W pętli `while` przeglądamy wartości  $k \in \pi^*[q - 1]$  do chwili znalezienia takiej, że  $P[k + 1] = P[q]$  (czyli  $k \in E_{q-1}$ ). Jeśli takie  $k$  istnieje, to w wierszu 7  $k$  jest największą taką wartością (startujemy od  $k = \pi[q - 1] = \max(\pi^*[q - 1])$  i każde iterowanie  $k \leftarrow \pi[k]$  zmniejsza  $k$ ) i w wierszu 8  $k$  zwiększa się o 1.

Jeśli takie  $k$  nie istnieje to w wierszu 7 i 9  $k = 0 = \pi[q]$ .

Z Wniosku 7: w obu przypadkach  $\pi[q]$  – poprawne.

# poprawność KMP-Matcher

Wynika z poprawności `Finite-Automaton-Matcher`: Pokażemy, że wiersze 6–9 `KMP-Matcher` symulują wiersz 4

`Finite-Automaton-Matcher` t.j.:  $q \leftarrow \delta(q, T[i])$ .

*Fakt:*  $\delta(q, T[i]) - 1 \in \pi^*[q] \cup \{q\}$  albo  $\delta(q, T[i]) = 0$ .

*D-d Faktu:* Niech  $k = \delta(q, T[i])$ . Wtedy  $P_k \sqsupset P_q T[i]$  (z def.  $\delta$  i  $\sigma$ ). Zatem, albo  $k = 0$ , albo  $k \geq 1$  i  $P_{k-1} \sqsupset P_q$  (przez usunięcie ostatniego znaku), czyli  $k - 1 \in \{k : P_k \sqsupset P_q\} = \pi^*[q] \cup \{q\}$ .  $\square$

Niech  $q'$  – wartość  $q$  na początku `while`. Pętla `while` w wierszach 6 i 7 przegląda  $\pi^*[q'] \cup \{q'\}$  w porządku malejącym aż  $P[q + 1] = T[i]$  lub  $q = 0$ . W pierwszym przypadku  $q$  zwiększone o 1 a w przeciwnym pozostaje  $q = 0$ . Z *Faktu*: nic nie pominęliśmy, więc po wierszu 9  $q = \delta(q', T[i])$ .

*Uwaga:* Wiersz 12 – aby uniknąć odwołań do  $P[m + 1]$ . Wiersz 12 nie szkodzi, bo  $\delta(m, a) = \delta(\pi[m], a)$  (równoważnie  $\sigma(Pa) = \sigma(P_{\pi[m]}a)$ ) – ćw.



# Algorytm Boyera-Moore'a

Boyer-Moore-Matcher( $T, P, \Sigma$ )

```
1   $n \leftarrow \text{length}[T]$ 
2   $m \leftarrow \text{length}[P]$ 
3   $\lambda \leftarrow \text{Last-Occurrence-Function}(P, m, \Sigma)$ 
4   $\gamma \leftarrow \text{Good-Suffix-Function}(P, m)$ 
5   $s \leftarrow 0$ 
6  while  $s \leq n - m$  do
7       $j \leftarrow m$ 
8      while  $j > 0$  and  $P[j] = T[s + j]$  do
9           $j \leftarrow j - 1$ 
10     if  $j = 0$  then
11         drukuj  "występuje z przesunięciem  $s$ "
12          $s \leftarrow s + \gamma[0]$ 
13     else  $s \leftarrow s + \max(\gamma[j], j - \lambda[T[s + j]])$ 
```

# Heurystyka niezgodności

Założmy, że  $P[j] \neq T[s + j]$ . Niech  
 $k = \max(\{0\} \cup \{i : 1 \leq i \leq m \wedge T[s + j] = P[i]\})$  Wtedy  
możemy przesunąć wzorzec o  $j - k$ :

- $k = 0$ :  $T[s + j]$  nie występuje w  $P$  i można przesunąć  $P$  o  $j$  nie pomijając żadnego wystąpienia  $P$  w  $T$ .
- $k < j$ : Ostatnie wystąpienie  $T[s + j]$  w  $P$  jest  $j - k$  pozycji na lewo od  $T[s + j]$ . Aby dopasować  $P$  do  $T$  trzeba go przesunąć o  $\geq j - k$  pozycji.
- $k > j$ : Proponowane przesunięcie  $j - k$  jest  $< 0$  i zostanie zignorowane.

# Last-Occurrence-Function

Last-Occurrence-Function( $P, m, \Sigma$ )

```
1 for each  $a \in \Sigma$  do
2    $\lambda[a] = 0$ 
3 for  $j \leftarrow 1$  to  $m$  do
4    $\lambda[P[j]] \leftarrow j$ 
5 return  $\lambda$ 
```

Czas:  $O(|\Sigma| + m)$ .

# Heurystyka dobrego sufiksu

Dla tekstów  $Q$  i  $R$  definiujemy relację

$$Q \sim R \equiv Q \sqsupset R \vee R \sqsupset Q.$$

Zachodzą własności:  $Q \sim R \equiv R \sim Q$  (symetria  $\sim$ ) oraz  $Q \sqsupset R$  i  $S \sqsupset R$  implikuje  $Q \sim S$  (Lemat 1).

Heurystyka: Jeśli  $P[j + 1 \dots m] = T[s + j + 1 \dots s + m]$  i  $P[j] \neq T[s + j]$ , to można przesunąć  $P$  o:

$$\gamma[j] = m - \max\{k : 0 \leq k < m \wedge P[j + 1 \dots m] \sim P_k\}$$

(najmniejsze dodatnie przesunięcie  $(m - k)$  gwarantujące, że  $T[s + j + 1 \dots s + m] \sim P_k$ .)

$\gamma$  (funkcja dobrego sufiksu) – dobrze zdefiniowana, gdyż  $P[j + 1 \dots m] \sim P_0$ .

Ponadto: Zawsze  $\gamma[j] \geq 1$  (zawsze dodatnie przesunięcie w algorytmie).

# funkcja dobrego sufiksu

**Fakt:**  $\gamma[j] \leq m - \pi[m]$  (bo:  $P_{\pi[m]} \sqsupset P$  oraz  $P[j + 1 \dots m] \sqsupset P$ , co implikuje  $P_{\pi[m]} \sim P[j + 1 \dots m]$ , czyli  $\max\{k : 0 \leq k < m \wedge P[j + 1 \dots m] \sim P_k\} \geq \pi[m]$ ).

**Wniosek:**

$\gamma[j] = m - \max\{k : \pi[m] \leq k < m \wedge P[j + 1 \dots m] \sim P_k\}$ .

Ponadto:  $P_k \sqsupset P[j + 1 \dots m]$  i  $k < m$  implikuje  $P_k \sqsupset P_m$  i (z definicji  $\pi$ )  $k \leq \pi[m]$ .

Mamy zatem:

$$\gamma[j] = m - \max(\{\pi[m]\} \cup \{k : 0 \leq k < m \wedge P[j + 1 \dots m] \sqsupset P_k\})(**)$$

(t.j. nie warto sprawdzać  $\{k : 0 \leq k < m \wedge P_k \sqsupset P[j + 1 \dots m]\}$ )

# funkcja dobrego sufiksu

Niech  $P'$  – odwrócony  $P$  (t.j.  $\forall_{1 \leq i \leq m} P'[i] = P[m - i + 1]$ ).

Niech  $\pi'$  – f-cja prefiksowa dla  $P'$ .

**Fakt:** Jeśli  $k = \max\{k : 0 \leq k < m \wedge P[j + 1 \dots m] \sqsupseteq P_k\}$ , to

$\pi'[l] = m - j$ , gdzie  $l = (m - k) + (m - j)$

(równoważnie:  $j = m - \pi'[l]$  i  $k = m - l + \pi'[l]$  (!)).

**Uwaga:**  $l \leq m$  (bo  $m - j \leq k$ , bo  $P[j + 1 \dots m] \sqsupseteq P_k$ ) oraz  $l \geq 1$  (bo  $k < m$  i  $j \leq m$ ). Zatem  $\pi'[l]$  – dobrze zdefiniowane.

**D-d Faktu:**  $P[j + 1 \dots m]$  w  $P$  odpowiada  $P'_{m-j}$  w  $P'$ .

$P[j + 1 \dots m] \sqsupseteq P_k$  odpowiada

$P'_{m-j} = P'[m - k + 1 \dots m - k + (m - j)]$ , czyli

$P'_{m-j} \sqsupseteq P'_{(m-k)+(m-j)} = P'_l$ . Stąd  $\pi'[l] \geq m - j$ , bo  $m - j < l$ .

...

# funkcja dobrego sufiksu

Założmy, że  $p = \pi'[l] > m - j$ . Wtedy z def.  $\pi'$ :  $P'_p \sqsupset P'_l$ , czyli

$P'[1 \dots p] = P'[l - p + 1 \dots l]$  co odpowiada:

$$P[m - p + 1 \dots m] = P[m - l + 1 \dots m - l + p].$$

Podstawiając  $l = 2m - k - j$  dostajemy:

$$P[m - p + 1 \dots m] = P[k - m + j + 1 \dots k - m + j + p].$$

Stąd  $P[m - p + 1 \dots m] \sqsupset P_{k-m+j+p}$ .

Stąd  $P[j + 1 \dots m] \sqsupset P_{k-m+j+p}$ , bo

$$m - p + 1 < m - (m - j) + 1 = j + 1.$$

Mamy  $k - m + j + p > k$ , bo  $p > m - j$ , oraz

$$k - m + j + p < k - m + j + l = k - m + j + (2m - k - j) = m,$$

bo  $p = \pi'[l] < l$ . Niech  $k' = k - m + j + p$ . Zatem  $k < k' < m$  i

$P[j + 1 \dots m] \sqsupset P_{k'}$  – sprzeczność z wyborem  $k$ . Stąd

$$\pi'[l] \leq m - j. \quad \square$$

# funkcja dobrego sufiksu

Korzystając z (!) oraz implikacji (ćw.):

$$(j = m - \pi'[l]) \Rightarrow (P[j + 1 \dots m] \sqsupset P_{m-l+\pi'[l]}) (!!)$$

mamy:

$$\begin{aligned}\gamma[j] &= m - \max(\{\pi[m]\} \cup \\ &\quad \{k : 0 \leq k < m \wedge P[j + 1 \dots m] \sqsupset P_k\}) (**) \\ &= m - \max(\{\pi[m]\} \cup \\ &\quad \{m - l + \pi'[l] : 1 \leq l \leq m \wedge j = m - \pi'[l]\}) \text{ z (!) i (!!)} \\ &= \min(\{m - \pi[m]\} \cup \\ &\quad \{l - \pi'[l] : 1 \leq l \leq m \wedge j = m - \pi'[l]\})\end{aligned}$$



# Good-Suffix-Function

Good-Suffix-Function( $P, m$ )

```
1  $\pi \leftarrow \text{Compute-Prefix-Function}(P)$ 
2  $P' \leftarrow \text{odwrócone } P$ 
3  $\pi' \leftarrow \text{Compute-Prefix-Function}(P')$ 
4 for  $j \leftarrow 0$  to  $m$  do
5    $\gamma[j] \leftarrow m - \pi[m]$ 
6 for  $l \leftarrow 1$  to  $m$  do
7    $j \leftarrow m - \pi'[l]$ 
8   if  $\gamma[j] > l - \pi'[l]$  then
9      $\gamma[j] \leftarrow l - \pi'[l]$ 
10 return  $\gamma$ 
```

Czas:  $O(m)$ .

Czas algorytmu Boyera-Moore'a:  $O((n - m + 1)m + |\Sigma|)$ ,  
jednak w praktyce dużo przesunięć może być pominiętych.