

Lecture #4

collection of texts

• Basic problem: search P as a substring of T

• Old applic = Editors / Spellcheck

Now = { Anti-virus $\Rightarrow \Sigma = \text{ASCII}$ (words)
 Search engines (words)

FINGERPRINT

Useful tool in many contexts

$$h(s) = \sum_{i=1}^{|s|} \frac{1}{2^i} s[i]$$

just as a binary representation

of course $s' = s'' \iff h(s') = h(s'')$

• Where is the problem? We cannot manage pointers larger than 4/8 bytes.

\rightarrow We need to square $h(s')$ in a square way; without introducing many suffixes so that we can "count" them.

idea resort modular arithmetic: $h_p(s) = h(s) \bmod p$

① p is chosen as a prime many properties

② bit of Horner rule to manage space intermediary results

③ incremental computation is possible



n -long substring

$$h_p(T_{r+1, r+u}) = 2^u \times h_p(T_{r+1, r}) - 2^u T[r] + T[r+u]$$

$O(1)$ time

FACT

Computing the $N-u+1$ fingerprints of n -long substrings takes $O(N)$ time, which is independent of n .

Prop 1

#primes smaller than n are $\left[\frac{n}{\ln n}, \frac{1.26 n}{\ln n} \right]$

done

Prop 2. If $n \geq 29$, the product of all primes $\leq n$ is larger than 2^n

Cor 3. If $n \geq 29$ and $x \leq 2^n$ then x has fewer than $\pi(n)$ prime divisors (distinct).

TEO. $|P| = m, |T| = n$ with $n \cdot m \geq 29$

p prime $\leq I$, where I will be at most

$$\mathbb{P}(\text{false match}) \leq \frac{\pi(n \cdot m)}{\pi(I)}$$

τ false match $\Leftrightarrow P \neq T_r$ but $h_p(P) = h_p(T_r) \Rightarrow p$ divides $h(P) - h(T_r)$

R all ~~false match positions~~ $\Leftrightarrow \forall p$ divides $\prod_{r \in R} [h(P) - h(T_r)]$

$\Rightarrow p$ is one of the prime divisors of that number

\Rightarrow How large is it? $|h(P) - h(T_r)| \leq 2^n$

$$\text{hence } |\prod \dots| \leq 2^{n \cdot m}$$

$\Rightarrow \#$ prime divisors (distinct) $\leq \pi(n \cdot m)$

$$\mathbb{P}(\text{false match}) = \mathbb{P}(p \text{ is one of those divisors}) \leq \frac{\pi(n \cdot m)}{\pi(I)} \quad \leftarrow p \text{ selected} \leq I$$

what about k primes? All k primes will give a false match $\Rightarrow \left[\frac{\pi(n \cdot m)}{\pi(I)} \right]^k$

a more refined argument shows that if τ is a false match then

the k primes divide $|h(P) - h(T_r)| \leq 2^n$, and they are ones in $\pi(n)$

$$\mathbb{P}(\) \leq m \cdot \left[\frac{\pi(n)}{\pi(I)} \right]^k$$

$$\Rightarrow \frac{n \cdot m}{I \cdot m} / \frac{n^k \cdot \log n}{2^k n + I \cdot m + I \cdot k \cdot n} \leq \frac{2 \log n}{n \cdot I \cdot m}$$

what about I ?

We must fit every arithmetic calculation in our RAM word

$$\text{Typically } I = m \cdot n^2 \Rightarrow \frac{\pi(n \cdot m)}{\pi(n^2 m)} \leq \frac{2.53}{n} = O\left(\frac{1}{n}\right)$$

$$\text{Sec Ops} \leq 4(\log n + \log m)$$

$\Rightarrow \#$ sig false matches is $O(1) \Rightarrow$ checking them takes $O(m)$

so we get a Las Vegas algorithm

$$\text{what about checking } p \text{ at each false match? } \mathbb{P}(t \text{ errors}) \leq \left(\frac{2.53}{n}\right)^t$$

$next(i) =$ longest prefix of $P[1, i-1]$ which is a suffix of that string.

HOW TO COMPUTE IT



verify longer if $P[j] \neq P[i-1]$

we try $next(i) = j-1$ known by induction

Not necessarily OK

base $\begin{cases} next(1) = -1 \\ next(2) = 0 \end{cases}$

$next(i) = \begin{cases} j & \text{if } P[i] = P[j] \text{ \& } j-1 = next(i-1) \\ \text{ITERATE ON } j & \text{else} \end{cases}$

Compute-next (P)

```

next(1) = -1; next(2) = 0;
for (i = 3; i ≤ n; i++)
    j = next(i-1) + 1;
    while ((T[i] ≠ T[j]) & j > 0)
        j = next(j) + 1;
    next(i) = j;

```

We must try all shifts:

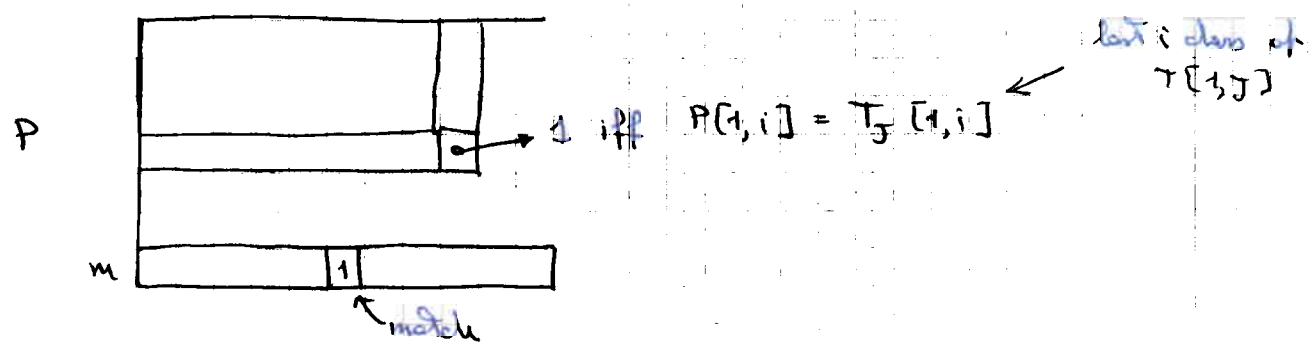
	18	19	20	21	22	
P =	a	b	a	b	a	b
	a	b	a	b	a	b
	a	b	a	b	a	b
	a	b	a	b	a	b

$next(22) = ?$
 $next(21) = 9$ $j=10$
 $next(10) = 4$ $c \neq a$
 $next(5) = 2$ $d \neq a$
 $a = a$ $j=3$

$\Rightarrow next(22) = 3$

AGREP

Turn everything into shift + AND ops.



IDEA ① Compute M column-by-column rightward.

② $U[i][\sigma] = 1$ iff $P[i] = \sigma$

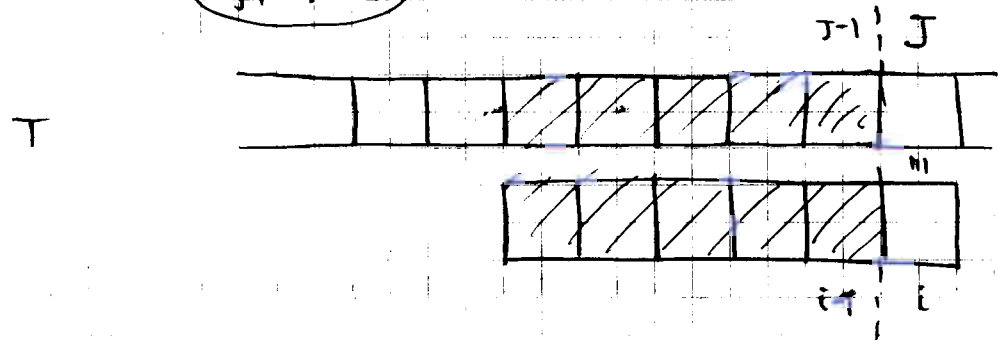
	A	B	C	D
A	1	0	0	0
B	0	1	0	0
P	A	1	0	0
A	1	0	0	0
D	0	0	0	1

characteristic matrix for P

③ Bit shift $(v) = [1|v_{m-1}]$ (drop left bit)

$M[i][j] = 1$ iff $M[i-1, j-1] = 1$ && $P[i] = T[j]$

$\begin{matrix} P[i] \\ \vdots \\ T_j[1, i] \end{matrix}$
 $\begin{matrix} P[1, i-1] \\ \vdots \\ T_{j-1}[1, i-1] \end{matrix}$
 $U[i][T[j]] = 1$



This holds for all $\ominus \oplus$ I do not want to merge these

Bit shift $(M[i][j-1])$ & $U[i][T[j]]$

First bit = 1 because this way we correctly manage the first bit of a column, which is 1 iff $U[1][T[j]]$ (first char of P matches $T[j]$)

bc while we have no carry from above

	b	x	a	b	a	a	c
a					1	1	0
b					0	0	0
a					1	0	0
a					0	1	0
d					0	0	0

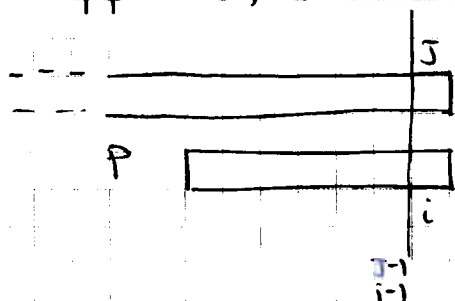
\uparrow $j-1$ \uparrow j

$$M[j][s] = \text{BitShift}(M[j][s]) \& U[j][a]$$

$$\text{BitShift} \left(\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right) \& \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \& \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

How to extend to k ~~mismatches~~ k mismatches

$M^h[i][j] = 1$ iff $P[1, i]$ matches $T_j[1, i]$ with at most h mismatches



RECURSIVE?

Two cases:

① $P[i] = T[j]$ ② $\leq h$ ~~mismatches~~ ^{mismatches} are in $P[1, i-1]$, $T_{j+1}[1, i-1]$

② $P[i] \neq T[j]$ ③ $\leq h-1$ mismatches are in \uparrow
no case

$$M^h[j][s] = (\text{BitShift}(M^h[j][s-1]) \& U[j][T[j]])$$

$$\parallel (\text{BitShift}(M^{h-1}[j][s-1]))$$

no case about $P[i] \neq T[j]$

	B	X	A	B	A	A	C
A					1	1	0
B					0	0	0
A					1	0	0
A					0	1	0
D					0	0	0

M^0 J

	B	X	A	B	A	A	C
A					1	1	1
B					0	1	1
A					1	0	0
A					0	1	0
D					0	0	1

M^1 J

$$M^1[i][j] = \left(\text{BitShift} \left(\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right) \& \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \right) \parallel \text{BitShift} \left(\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right)$$

$M^1[i][j-1]$ $U[i]$ $M^0[i][j-1]$

Complexity

$$n * k * \frac{m}{w} \stackrel{\text{Typically}}{=} O(u * k)$$

$$M^1[i][j+1] = \left(\text{BitShift} \left(\begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \right) \& \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \right) \parallel \text{BitShift} \left(\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \right)$$

case of match @ j
and k errors before is not possible

$\begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

How do we manage sets of chars?

→ Union of $U[i]$, V in the set.

