

Signature File Partitioning (cont.)

- Fixed prefix

Use the initial part as a key

We want to consider both sequential and parallel processing

Partition activation ratio (PAR), signature activation ratio (SAR) (covered last time)

S1: 0111 1000 $k=2$
S2: 1000 1011 (Use 2-prefix)
S3: 0011 1100
S4: 1100 0011
S5: 0110 1100
S6: 1001 0011
S7: 0000 1111



Q1: 1110 0001
Q2: 0000 1111
Q3: 0110 0011

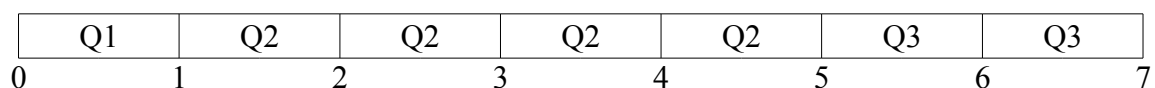
We select the pages of which the prefixes P_i satisfy the rule $P_i \& Q_i = Q_i$.

$k = 2$

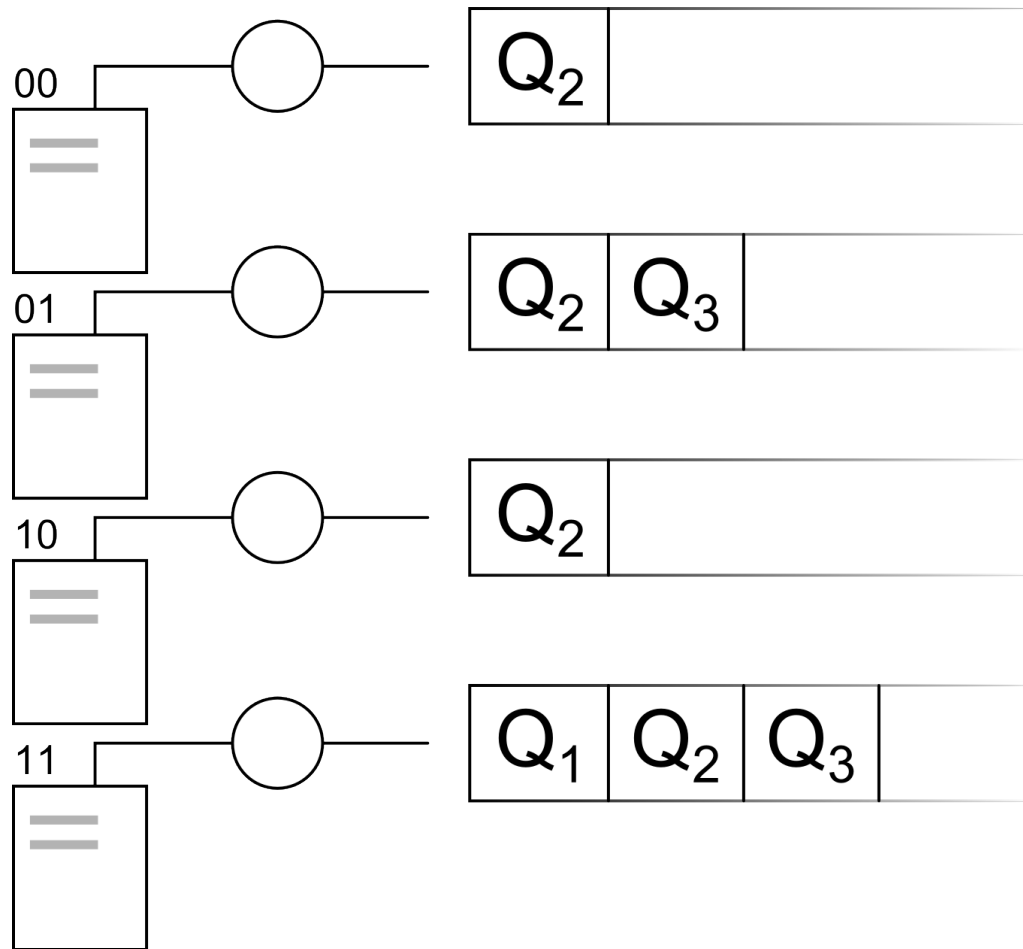
Q1: 111... 11 (**1 page selected**)
 Q2: 000... 00, 01, 10, 11 (**4**)
 Q3: 011... 01, 11 (**2**)

Sequential processing

All queries arrive at $t=0$



Parallel processing



Turnaround time = Time of completion – Time of arrival

	Time of arrival	
	Seq.	Par.
Q1	1	1
Q2	5	2
Q3	7	3

Serial average turnaround time = $(1+5+7) = 3 \text{ tu}$ (tu: Time units)

Parallel average turnaround time = $(1+2+3) = 2 \text{ tu}$

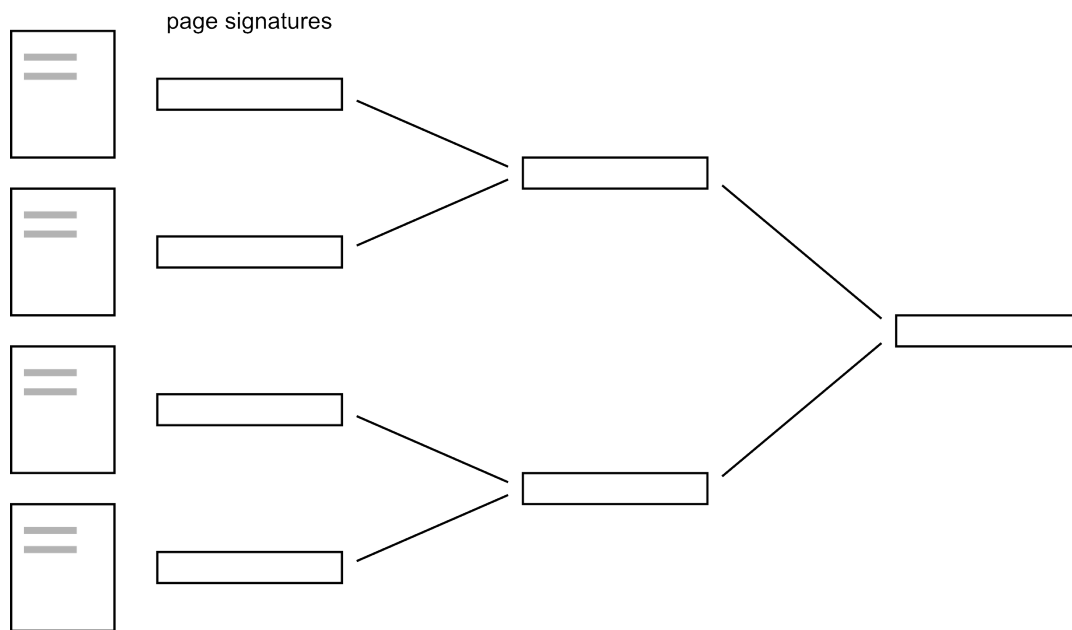
Throughput (T): No. of jobs completed per unit time

$$T_s = 3/7$$

$$T_p = 3/3$$

$$T_p > T_s$$

Signature Tree Structure



Page signatures are ANDed (“superposed”) to obtain superpage signatures. Any branch that do not satisfy the query condition is simply pruned.

Information Filtering

- Also known as “selective dissemination of information”
- It is the mirror image of information retrieval
- Change the roles of queries and documents

User profiles \equiv Queries

We receive **documents**, and we match these **documents** with **user profiles**. That is, we send the incoming document to the matching user profiles. Therefore the process is reversed.

We can convert the documents to vectors, and calculate similarity

Or, if an incoming document contains all user profile terms, send it to the owner of the profile (we use the AND operand).

(see: Yan, Tak W., and Héctor García-Molina. "Index structures for selective dissemination of information under the boolean model." *ACM Transactions on Database Systems (TODS)* 19.2 (1994): 332-364.)

There are several approaches described in the paper:

Brute force approach (Sequential comparison):

We compare the incoming document with each user profile one by one

Counting method:

Sample profiles:

p_1 : (a, b)

p_2 : (a, d)

p_3 : (a, d, e)

p_4 : (b, f)

p_5 : (c, d, e, f)

Sample document: a, c, a, f, b, c

Unique terms: {a, b, c, f}

Construct an inverted index:

Directory (in memory)		Inverted list (on disk)		
a	→	p ₁	p ₂	p ₃
b	→	p ₁	p ₄	
c	→	p ₅		
d	→	p ₂	p ₃	p ₅
e	→	p ₃	p ₅	
f	→	p ₄	p ₅	

Before processing a document set all count entries equal to 0.

Take the unique terms of the incoming document, visit the posting list and increment the counts for matching profiles.

	Total		a	b	c	f
p ₁	2	p ₁	0	1	2	
p ₂	2	p ₂	0	1		
p ₃	3	p ₃	0	1		
p ₄	2	p ₄	0		1	2
p ₅	4	p ₅	0		1	2

In this example, we obtained the required counts for p₁ and p₄; so they will receive the document.

Assume that we know the occurrence frequency of the terms in documents:

A profile appears in one of the inverted lists only in its term that appears most infrequently in documents:

The expectation: The term that appears less frequently in documents will appear more frequently in user profiles.

1. Distribute the signatures below into pages with (a) k=2 and (b) k=3 prefix.

a. For $k=2$, pages have these prefixes: $P_0=00$, $P_1=01$, $P_2=10$, $P_3=11$.

Then $P_0 = \{S_1, S_5\}$; $P_1 = \{S_3, S_8\}$, $P_2 = \{S_2, S_7\}$, $P_3 = \{S_4, S_6\}$

b. For $k=2$, pages are $P_0=000$, $P_1=001$, $P_2=010$, $P_3=011$, $P_4=100$, $P_5=101$, $P_6=110$, $P_7=111$.

Then $P_0 = \{S_5\}$, $P_1 = \{S_1\}$, $P_2 = \{S_3\}$, $P_3 = \{S_8\}$, $P_4 = \{S_7\}$, $P_5 = \{S_2\}$, $P_6 = \{S_6\}$, $P_7 = \{S_4\}$.

2. For both cases, which page gets retrieved the most? Why?

For $k=2$, P3 gets retrieved every time. For $k=3$, P7 gets retrieved every time. The reason is that these prefixes are all 1s, and $P_i \& Q_i = Q_i$ is satisfied for any query.