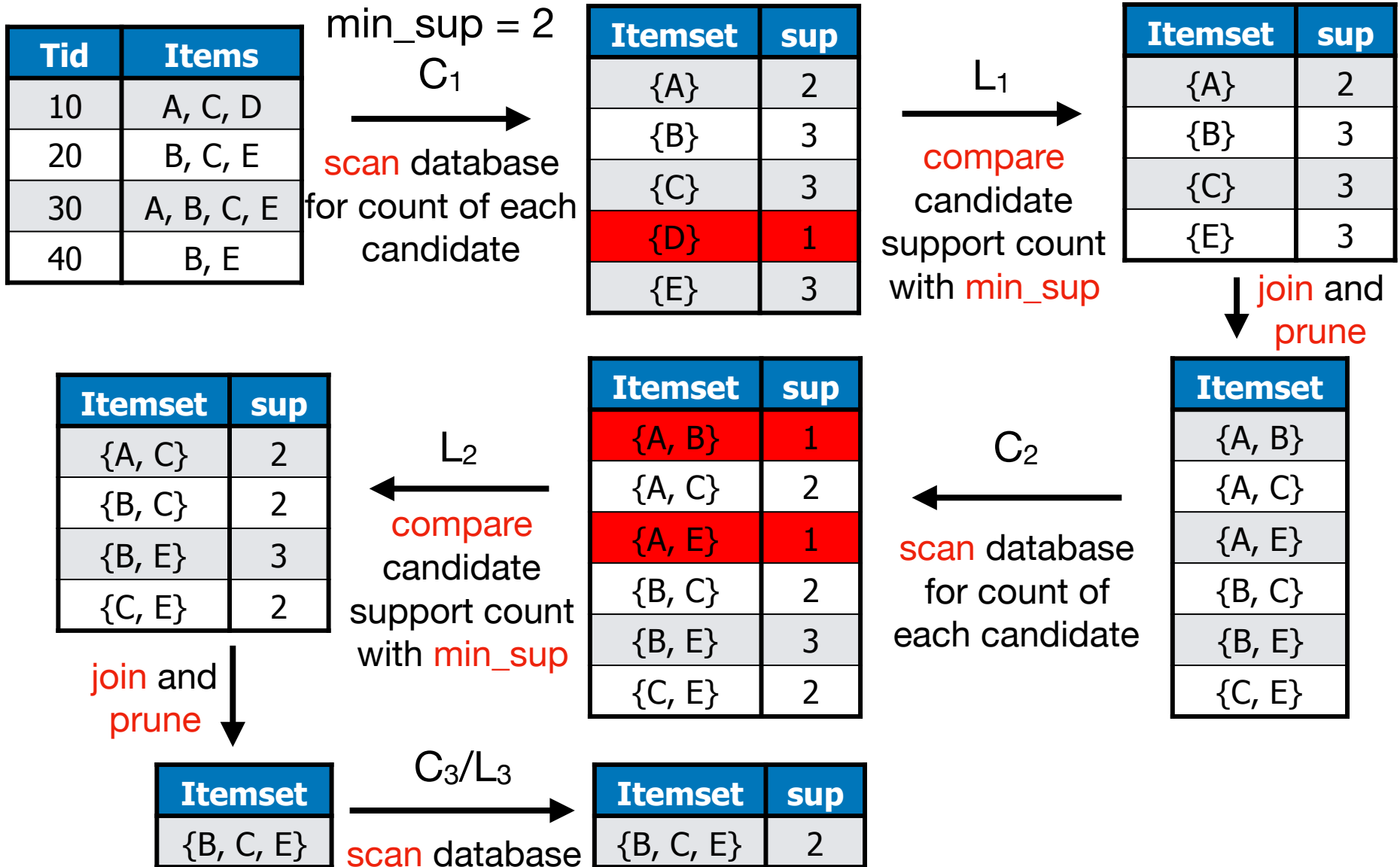


# Apriori

- How to generate candidates?
  - Step 1: self-joining  $L_k$
  - Step 2: pruning
- Example of Candidate-generation
  1.  $L_3 = \{abc, abd, acd, ace, bcd\}$
  2. Self-joining  $L_3 \otimes L_3$ : abcd from abc and abd; acde from acd and ace
  3. Pruning: acde is removed because ade is not in  $L_3$
  4.  $C_4 = \{abcd\}$

# Apriori



# Apriori

$C_k$ : Candidate itemset of size  $k$

$L_k$ : Frequent itemset of size  $k$

$L_1 = \{1\text{-frequent items}\};$

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **do begin**

$C_{k+1}$  = candidates generated from  $L_k$ ;

**for each** transaction  $t$  in database **do**

        increment the count of all candidates in  $C_{k+1}$  that are  
        contained in  $t$

**end**

$L_{k+1}$  = candidates in  $C_{k+1}$  with  $\text{min\_sup}$

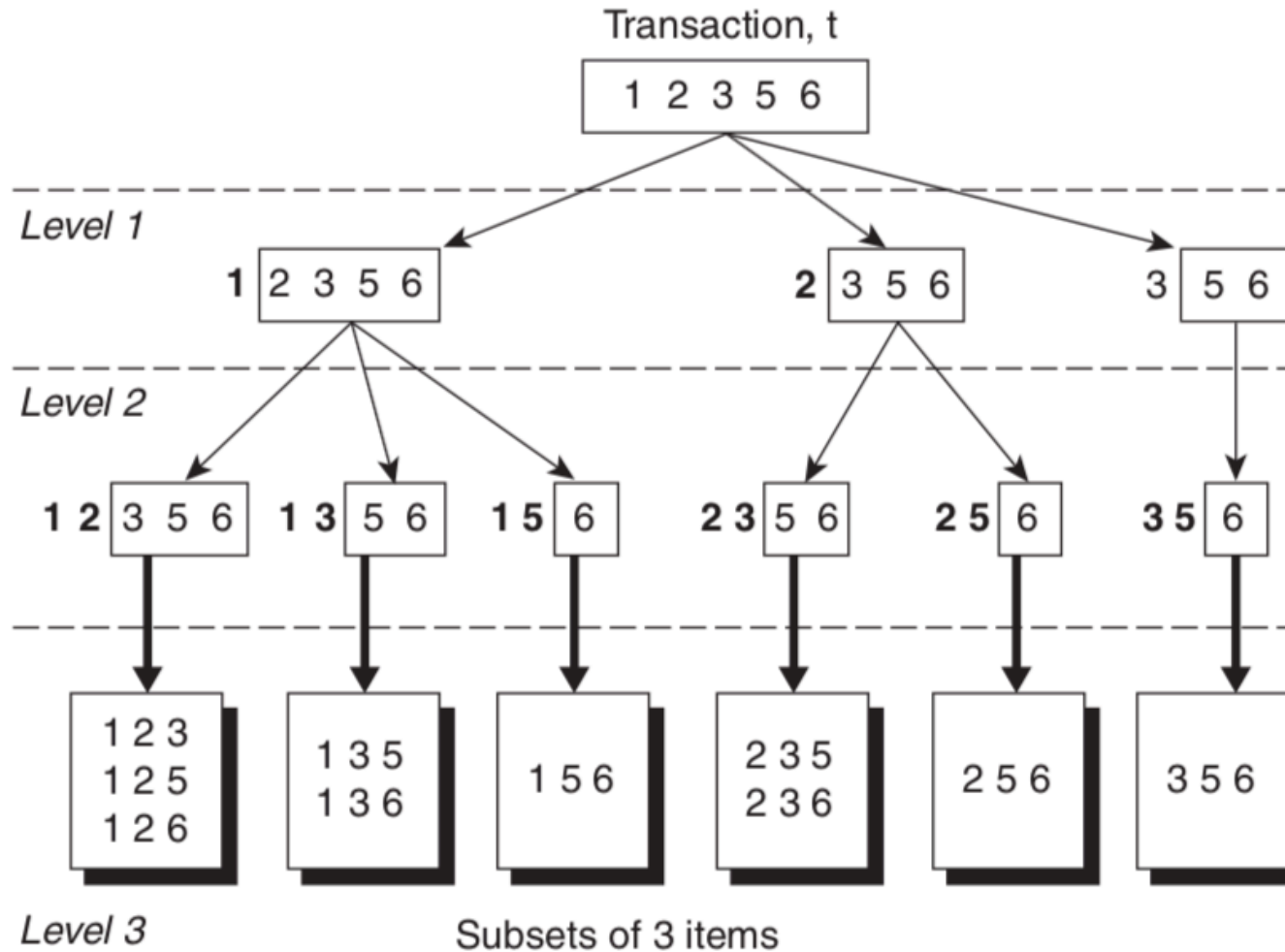
**end**

**return**  $\cup_k L_k$ ;

# Apriori

- How to count supports of each candidate?
  - The total number of candidates can be huge
  - One transaction may contain many candidates
  - **Support Counting Method:**
    - store candidate itemsets in a **hash-tree**
    - **leaf node** of hash-tree contains a list of itemsets and counts
    - **interior node** contains a hash table

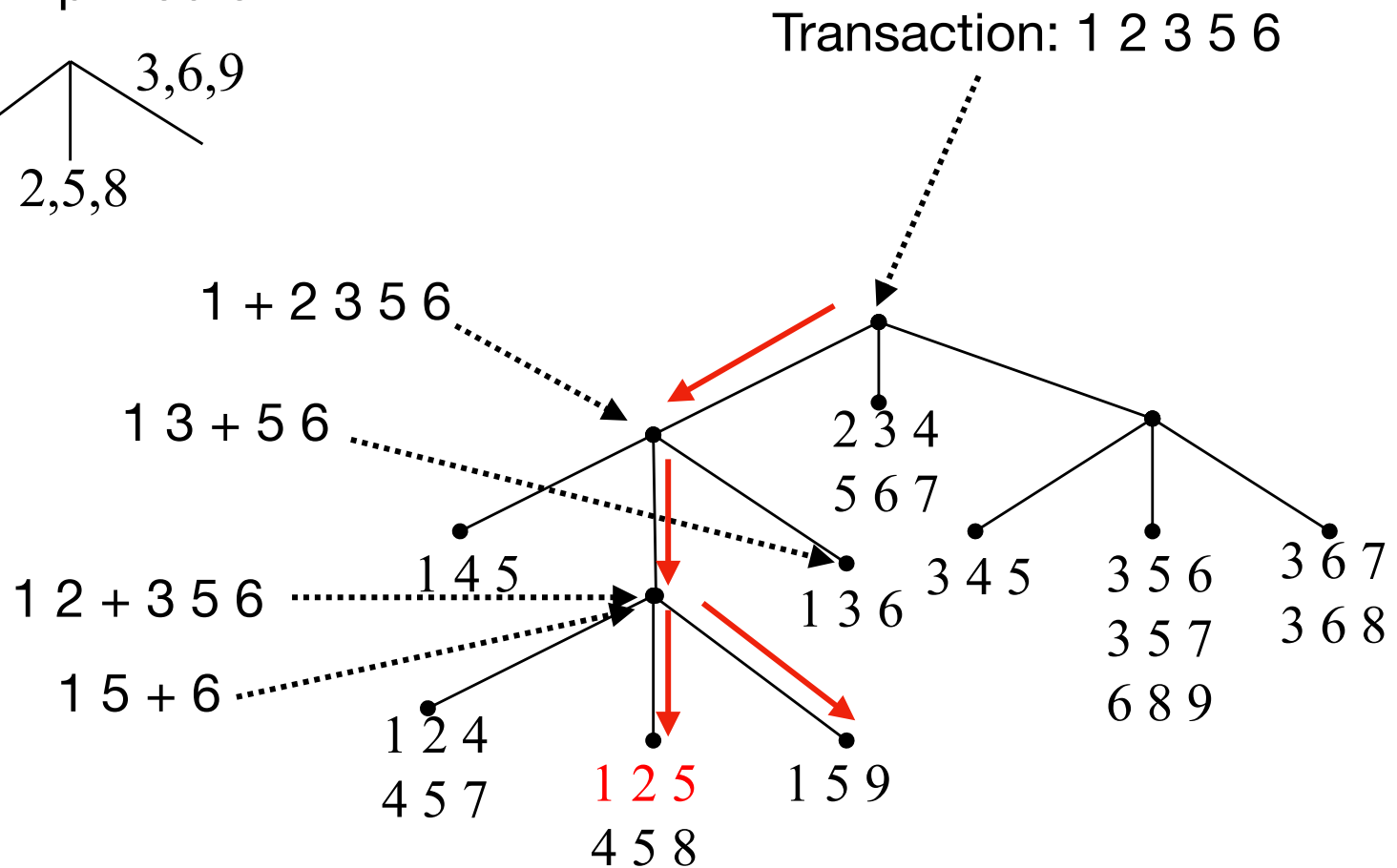
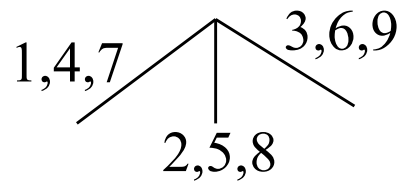
# Apriori



Prefix structure enumerating 3-itemset in Transaction t

# Apriori

Hash function  
 $h(p) = p \bmod 3$

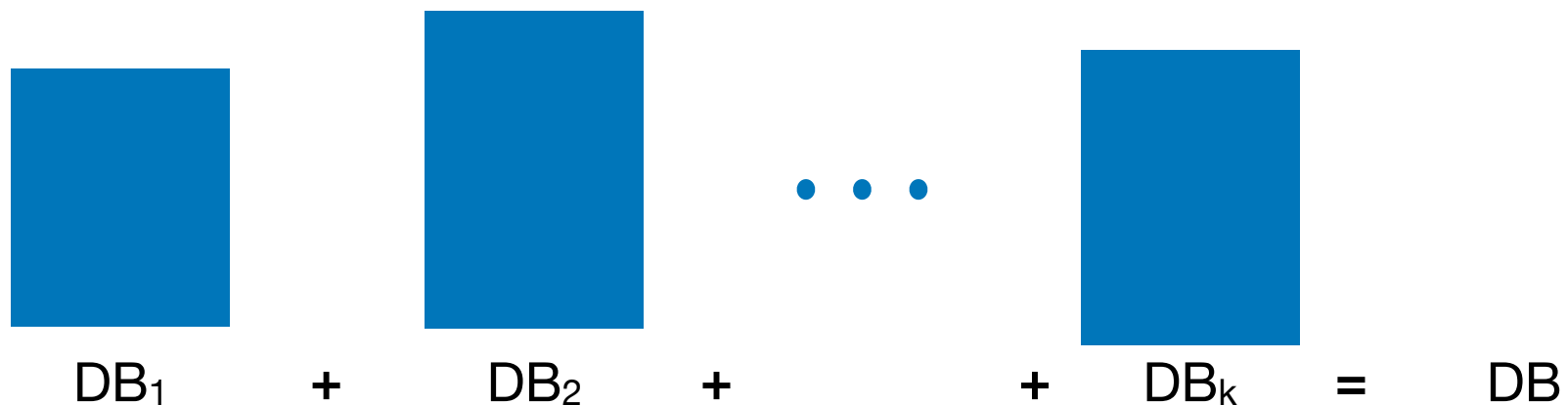


# Improving the Efficiency of Apriori

- Challenges:
  - Multiple scans of transaction database
  - Huge number of candidates
  - Support counting for candidates
- Improving the Efficiency of Apriori
  - Reduce passes of transaction database scans
  - Shrink number of candidates
  - Facilitate support counting of candidates

# Improving the Efficiency of Apriori

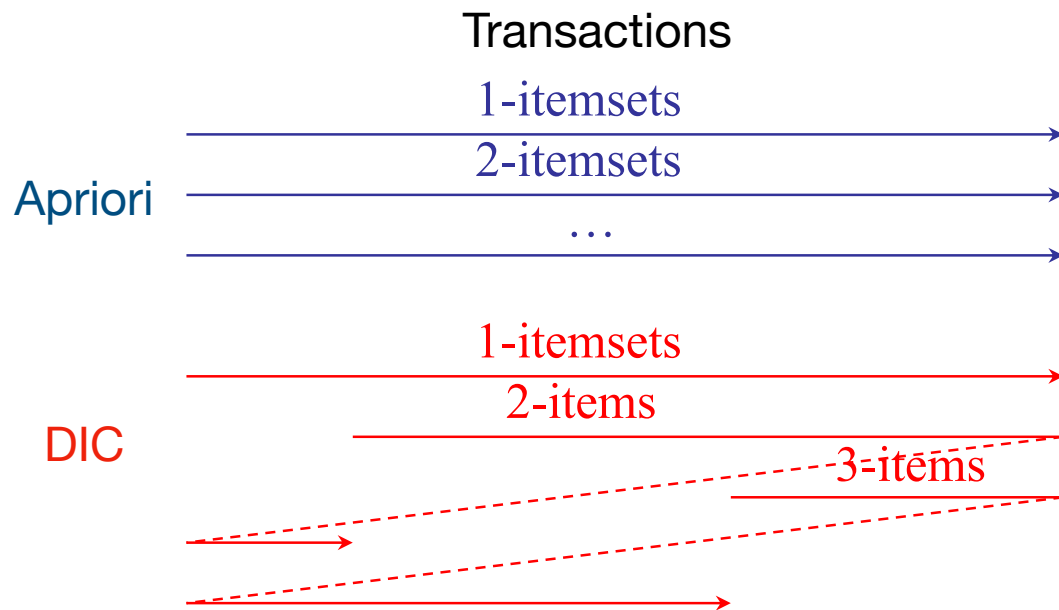
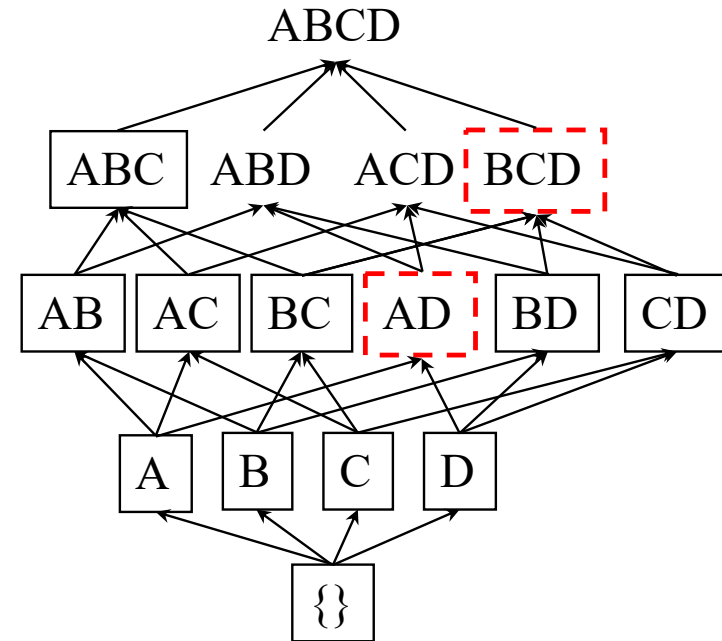
- **Partition** (reduce scans): partition data to find candidate itemsets
  - Any itemset that is potentially frequent (relative support  $\geq \text{min\_sup}$ ) must be frequent (relative support in the partition  $\geq \text{min\_sup}$ ) in at least one of the partition
  - Scan 1: partition database and find local frequent patterns
  - Scan 2: assess the actual support of each candidate to determine the global frequent itemsets





# Improving the Efficiency of Apriori

- **Dynamic itemset counting** (reduce scans): adding candidate itemsets at different points during a scan
- new candidate itemsets can be added at any start point (rather than determined only before scan)



- once both A and D are determined frequent, the counting of AD begins
- Once all length 2 subsets of BCD are determined frequent, the counting of BCD begins

# Improving the Efficiency of Apriori

- **Hash-based Technique** (shrink number of candidates): hashing itemsets into corresponding buckets
- A k-itemset whose corresponding hashing bucket count is below  $\text{min\_sup}$  cannot be frequent

$\text{min\_sup} = 3$

$H_2$

bucket address	0	1	2	3	4	5	6
bucket count	2	2	4	2	2	4	4
bucket contents	{I1, I4} {I3, I5}	{I1, I5} {I1, I5}	{I2, I3} {I2, I3} {I2, I3}	{I2, I4} {I2, I4}	{I2, I5} {I2, I5}	{I1, I2} {I1, I2}	{I1, I3} {I1, I3}

Create hash table  $H_2$   
 using hash function  
 $h(x, y) = ((\text{order of } x) \times 10 + (\text{order of } y)) \bmod 7$

$$h(1, 4) = 1 * 10 + 4 = 0 \bmod 7$$

$$h(3, 5) = 3 * 10 + 5 = 0 \bmod 7$$

# Improving the Efficiency of Apriori

- **Sampling**: mining on a subset of the given data
  - Trade off some degree of accuracy against efficiency
  - Select sample  $S$  of original database, mine frequent patterns within  $S$  (a lower support threshold) instead of the entire database  $\rightarrow$  the set of frequent itemsets local to  $S = L_S$
  - Scan the rest of database once to compute the actual frequencies of each itemset in  $L_S$
  - If  $L_S$  actually contains all the frequent itemsets, stop; otherwise
  - Scan database again for possible missing frequent itemsets

# A Frequent-Pattern Growth Approach

- Bottlenecks of Apriori
  - Breadth-first (i.e., level-wise) search
  - Candidate generation and test, often generates a huge number of candidates
- FP-Growth
  - Depth-first search
  - Avoid explicit candidate generation
  - Grow long patterns from short ones using local frequent items
    - “abc” is a frequent pattern
    - Get all transactions having “abc,” i.e., project database  $D$  on abc:  $D \mid abc$
    - “d” is a local frequent item in  $D \mid abc \rightarrow abcd$  is a frequent pattern

# A Frequent-Pattern Growth Approach

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>	
100	{ <i>f, a, c, d, g, i, m, p</i> }	{ <i>f, c, a, m, p</i> }	min_sup = 3
200	{ <i>a, b, c, f, l, m, o</i> }	{ <i>f, c, a, b, m</i> }	
300	{ <i>b, f, h, j, o, w</i> }	{ <i>f, b</i> }	F-list = f-c-a-b-m-p
400	{ <i>b, c, k, s, p</i> }	{ <i>c, b, p</i> }	
500	{ <i>a, f, c, e, l, p, m, n</i> }	{ <i>f, c, a, m, p</i> }	

1. Scan database once, find frequent 1-itemset
2. Sort frequent items in frequency **descending** order  
—> F-list

<b>Header Table</b>	
<u><i>Item</i></u>	<u><i>frequency head</i></u>
<i>f</i>	4
<i>c</i>	4
<i>a</i>	3
<i>b</i>	3
<i>m</i>	3
<i>p</i>	3

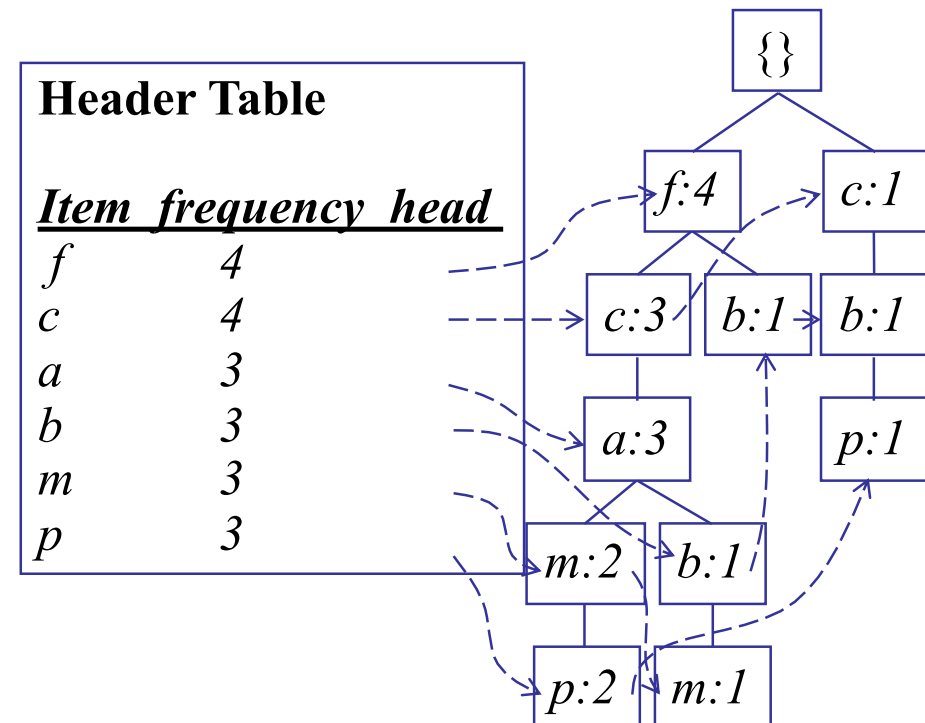
# A Frequent-Pattern Growth Approach

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{ <i>f, a, c, d, g, i, m, p</i> }	{ <i>f, c, a, m, p</i> }
200	{ <i>a, b, c, f, l, m, o</i> }	{ <i>f, c, a, b, m</i> }
300	{ <i>b, f, h, j, o, w</i> }	{ <i>f, b</i> }
400	{ <i>b, c, k, s, p</i> }	{ <i>c, b, p</i> }
500	{ <i>a, f, c, e, l, p, m, n</i> }	{ <i>f, c, a, m, p</i> }

min\_sup = 3

F-list = f-c-a-b-m-p

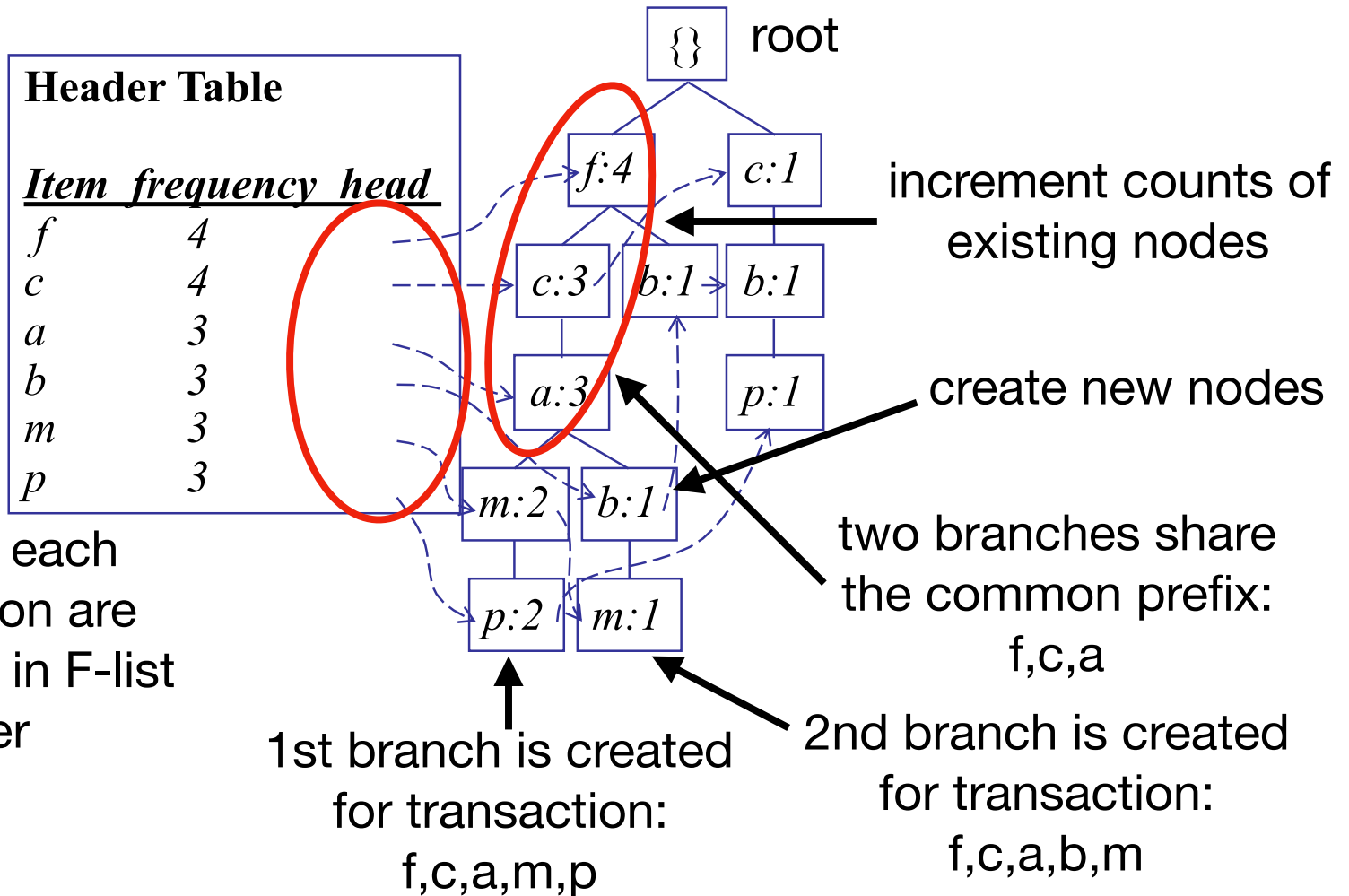
1. Scan database once, find frequent 1-itemset
2. Sort frequent items in frequency **descending** order → F-list
3. Scan database again, construct FP-tree
4. Mine FP-tree



# How to Construct FP-tree?

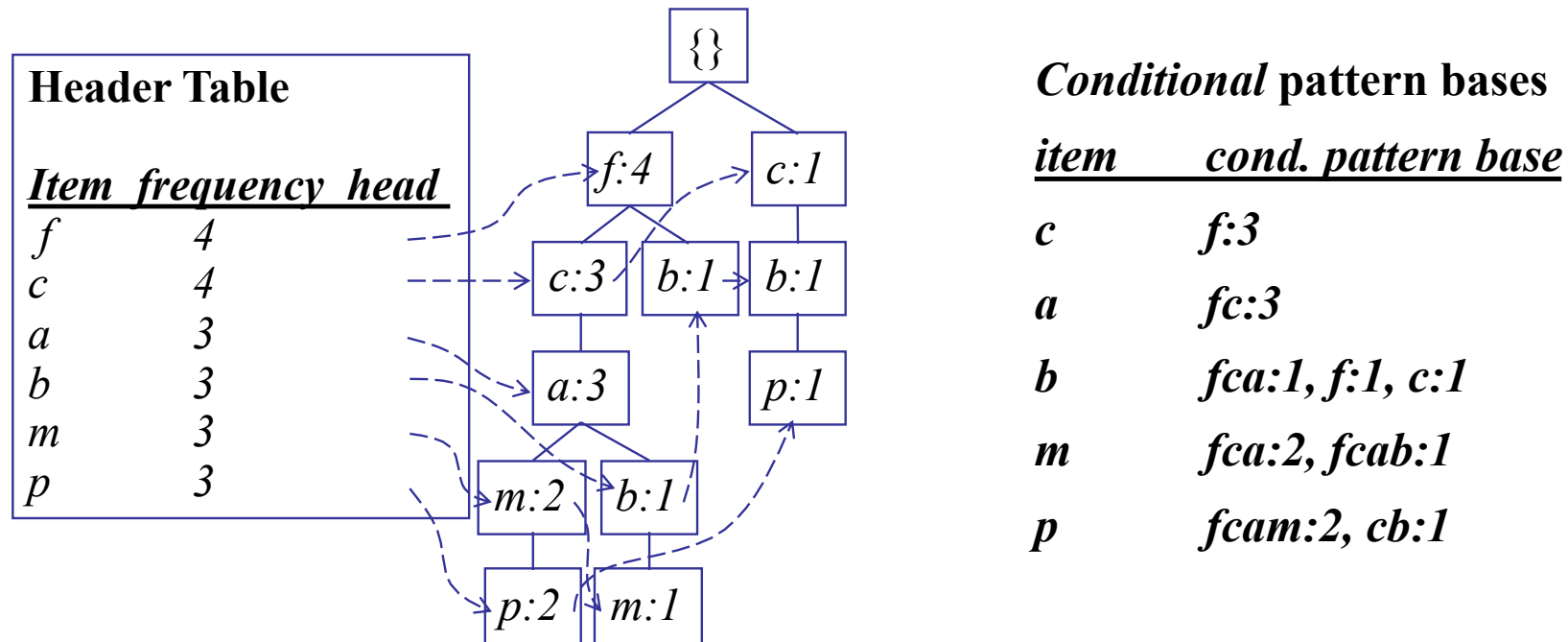
**FP-tree**: a **compressed** representation of database.  
It retains the itemset **association** information.

To facilitate tree traversal, each item points to its occurrence in the tree via **node-link**



# How to Mine FP-tree?

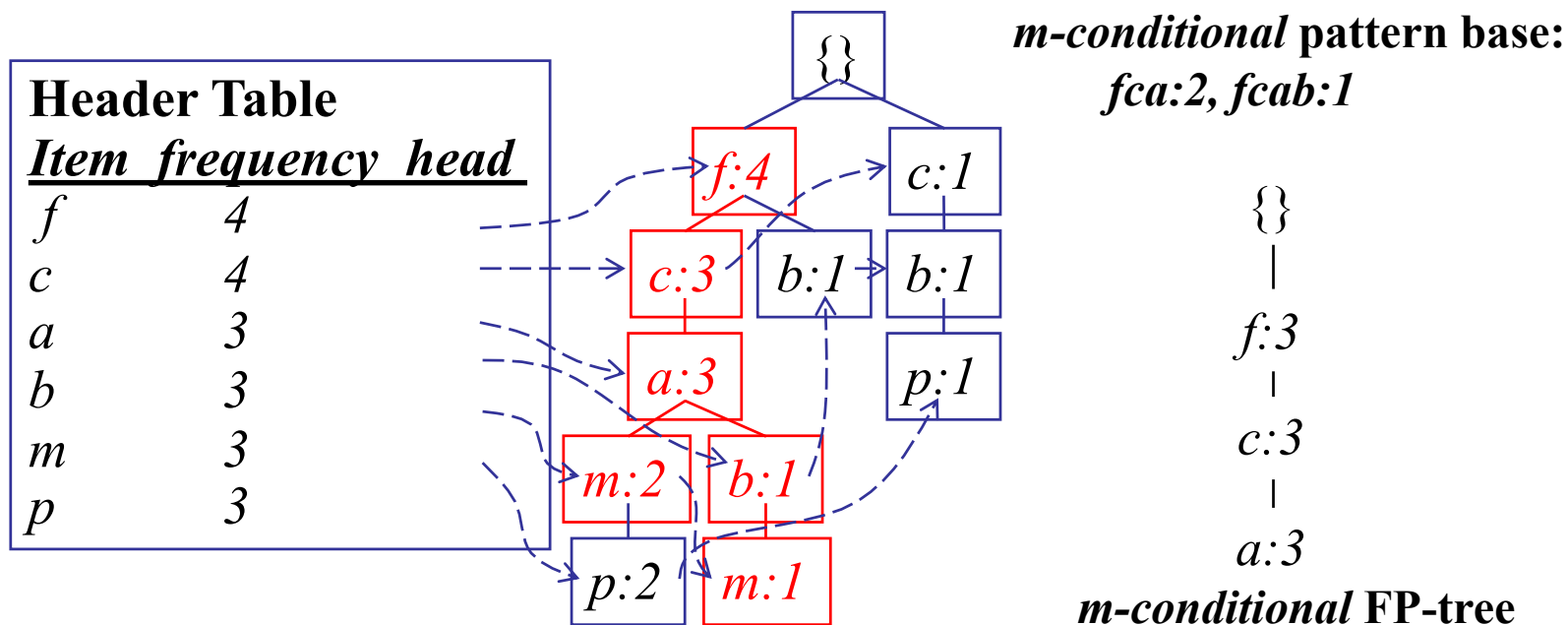
1. Start from each frequent length-1 pattern (**suffix pattern, usually the last item in F-list**) to construct its **conditional pattern base** (**prefix paths** co-occurring with the suffix)





# How to Mine FP-tree?

1. Start from each frequent length-1 pattern (**suffix pattern, usually the last item in F-list**) to construct its **conditional pattern base**
2. Construct the **conditional FP-tree** based on the conditional pattern base

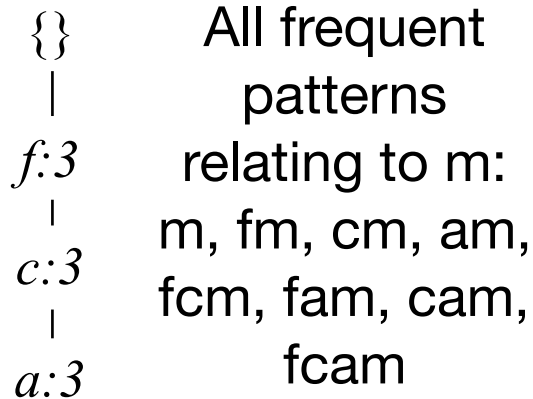


# How to Mine FP-tree?

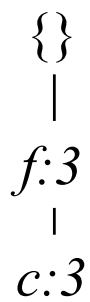
1. Start from each frequent length-1 pattern (**suffix pattern, usually the last item in F-list**) to construct its **conditional pattern base**
2. Construct the **conditional FP-tree** based on the conditional pattern base
3. Mining **recursively** on each conditional FP-tree until the resulting FP-tree is **empty**, or it contains only **a single path** — which will generate frequent patterns out of all combinations of its sub-paths

*m*-conditional pattern base:

*fca:2, fcab:1*

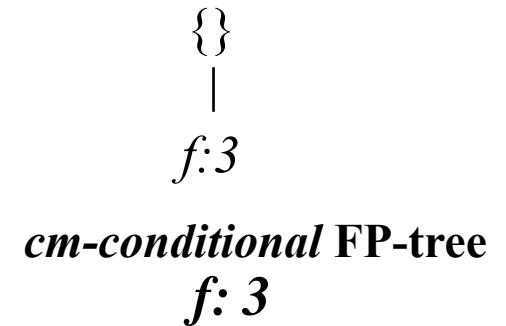


*m*-conditional FP-tree



*am*-conditional FP-tree

*fc: 3*



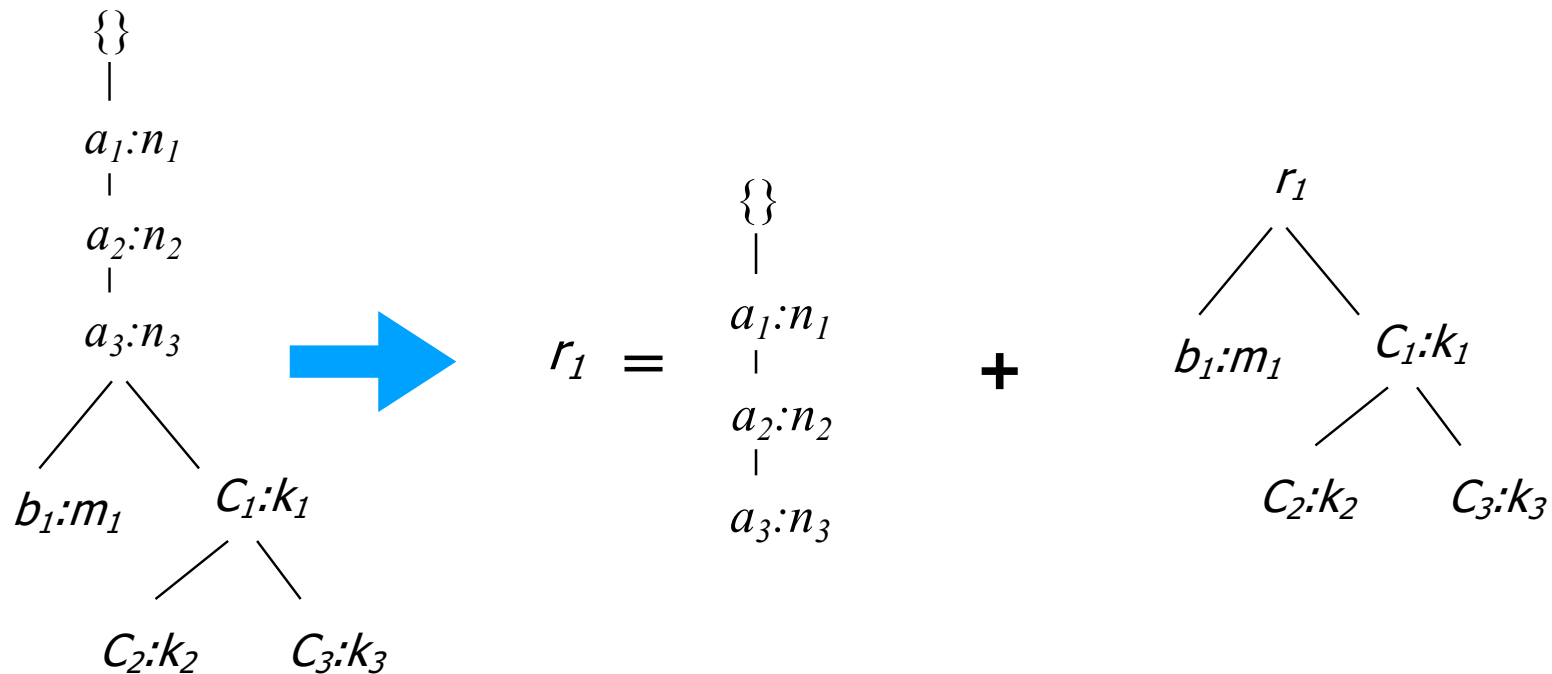
*cm*-conditional FP-tree

*cam*-conditional FP-tree

*f: 3*

# Single Prefix Path in FP-tree

- Suppose a (conditional) FP-tree has a shared **single prefix-path**
- Mining can be decomposed into two parts
  - Reduction of the single prefix path into one node
  - Concatenation of the mining results of the two parts



# Scaling FP-Growth

- What if FP-tree cannot fit into memory?
  - **Database projection**: partition a database into a set of projected databases, then construct and mine FP-tree for each projected database
  - **Parallel projection**:
    - project the database in parallel for each frequent item
    - all partitions are processed in parallel
    - space costly
  - **Partition projection**:
    - project a transaction to a frequent item  $x$  if there is no any other item after  $x$  in the list of frequent items appearing in the transaction
    - a transaction is projected to only one projected database

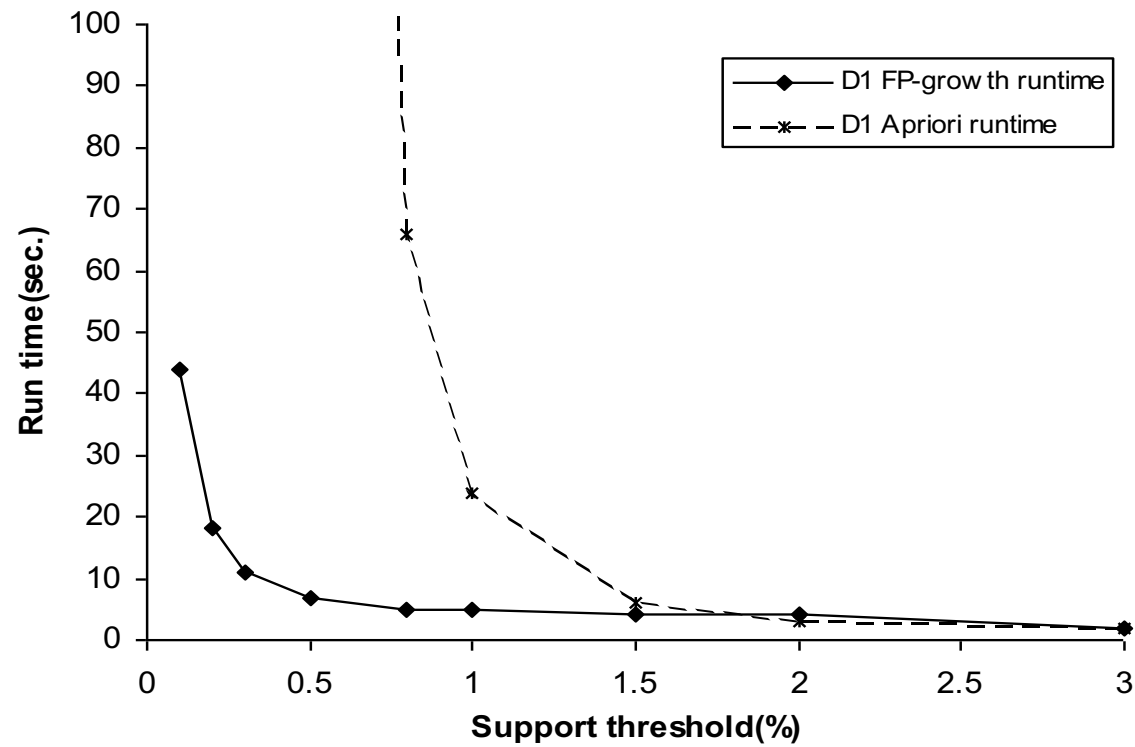
# Benefits of FP-tree

- Completeness
  - Preserve complete information for frequent pattern mining
  - Never break a long pattern of any transaction
- Compactness
  - Reduce irrelevant info — infrequent items are gone
  - Items in frequency descending order: occurs more frequently, the more likely to be shared
  - Never be larger than the original database (not including node-links and the count fields)

# Benefits of FP-Growth

- Divide-and-conquer:
  - Decompose both the mining task and database according to the frequent patterns obtained so far
  - Lead to focused search of smaller databases
- Other factors:
  - No candidate generation, no candidate test
  - Compressed database: FP-tree
  - No repeated scan of the entire database
  - Basic operations: count local frequent items and build sub FP-tree, no pattern search and matching

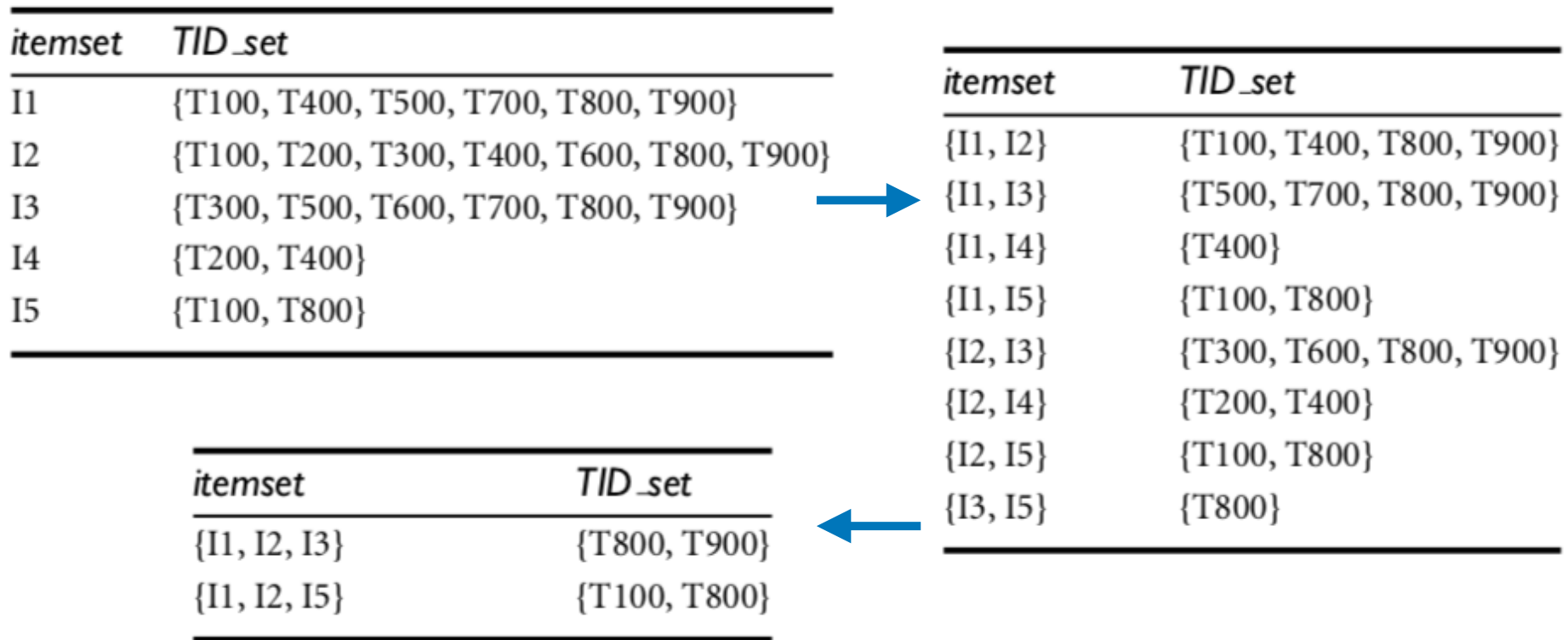
# Performance of FP-Growth in Large Datasets



FP-Growth vs. Apriori

# ECLAT: Frequent Pattern Mining with Vertical Data Format

- Vertical data format: itemset — transID\_set
  - transID\_set: a set of transaction IDs containing the itemset
- Derive frequent patterns based on the **intersections** of transID\_set





# ECLAT: Frequent Pattern Mining with Vertical Data Format

- Vertical data format: itemset — transID\_set
  - transID\_set: a set of transaction IDs containing the itemset
- Derive frequent patterns based on the **intersections** of transID\_set
- Use **diffset** to reduce the cost of storing long transID\_set
  - $\{I1\} = \{T100, T400, T500, T700, T800, T900\}$
  - $\{I1, I2\} = \{T100, T400, T800, T900\}$
  - $\text{diffset}(\{I1\}, \{I1, I2\}) = \{T500, T700\}$

# Summary

- Frequent itemset mining methods:
  - Apriori: candidate generation-and-test
  - Improving efficiency of Apriori: partition, dynamic item counting, hash-based technique, sampling
  - FP-Growth: depth-first search
  - Scaling of FP-Growth: database projection
  - Frequent pattern mining with vertical data format

# Outline

- Basic Concepts in Frequent Pattern Mining
- Frequent Itemset Mining Methods
- Pattern Evaluation Methods

# Pattern Evaluation Methods: Correlations

- play basketball  $\Rightarrow$  eat cereal [40%, 66.7%] is misleading
  - the overall % of students eating cereal is 75% > 66.7%
- play basketball  $\Rightarrow$  not eat cereal [20%, 33.3%] is more accurate
- **Lift**: a measure of dependent/correlated event

$$\text{lift} = \frac{P(A \cup B)}{P(A)P(B)} = \frac{P(B|A)}{P(B)}$$

	Basketball	Not	Sum
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

$$\text{lift}(\text{Basketball}, \text{Cereal}) = \frac{2000/5000}{(3000/5000) \times (3750/5000)} = 0.89 \quad < 1, \text{ negatively correlated}$$

$$\text{lift}(\text{Basketball}, \text{Notcereal}) = \frac{1000/5000}{(3000/5000) \times (1250/5000)} = 1.33$$

# Other Pattern Evaluation Methods

- $\chi^2$  measure, all\_confidence measure, max\_confidence measure, Kulczynski measure, ...