

Part 6

Association Analysis: Basic Concepts and Algorithms

NGURAH AGUS SANJAYA ER, S.KOM, M.KOM, PH.D
E-mail: agus.sanjaya@cs.unud.ac.id



Association Rule Mining

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Association Rules

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\}$,
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\}$,
 $\{\text{Beer, Bread}\} \rightarrow \{\text{Milk}\}$,

Definition: Frequent Itemset

- **Itemset**
 - A collection of one or more items
 - Example: {Milk, Bread, Diaper}
 - k-itemset
 - An itemset that contains k items
- **Support count (σ)**
 - Frequency of occurrence of an itemset
 - E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- **Support**
 - Fraction of transactions that contain an itemset
 - E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$
- **Frequent Itemset**
 - An itemset whose support is greater than or equal to a $minsup$ threshold

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Definition: Association Rule



- Association Rule
 - An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
 - $X \cap Y = \emptyset$
 - Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$
- Rule Evaluation Metrics
 - Support (s)
 - Fraction of transactions that contain both X and Y
 - Confidence (c)
 - Measures how often items in Y appear in transactions that contain X

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example:

$$\{\text{Milk, Diaper}\} \Rightarrow \text{Beer}$$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

Association Rule Mining Task



- Given a set of transactions T , the goal of association rule mining is to find all rules having
 - support $\geq \text{minsup}$ threshold
 - confidence $\geq \text{minconf}$ threshold
- Brute-force approach:
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Prune rules that fail the minsup and minconf thresholds
⇒ Computationally prohibitive!

Mining Association Rules



TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Rules:

$\{\text{Milk}, \text{Diaper}\} \rightarrow \{\text{Beer}\}$ ($s=0.4, c=0.67$)
 $\{\text{Milk}, \text{Beer}\} \rightarrow \{\text{Diaper}\}$ ($s=0.4, c=1.0$)
 $\{\text{Diaper}, \text{Beer}\} \rightarrow \{\text{Milk}\}$ ($s=0.4, c=0.67$)
 $\{\text{Beer}\} \rightarrow \{\text{Milk}, \text{Diaper}\}$ ($s=0.4, c=0.67$)
 $\{\text{Diaper}\} \rightarrow \{\text{Milk}, \text{Beer}\}$ ($s=0.4, c=0.5$)
 $\{\text{Milk}\} \rightarrow \{\text{Diaper}, \text{Beer}\}$ ($s=0.4, c=0.5$)

Observations:

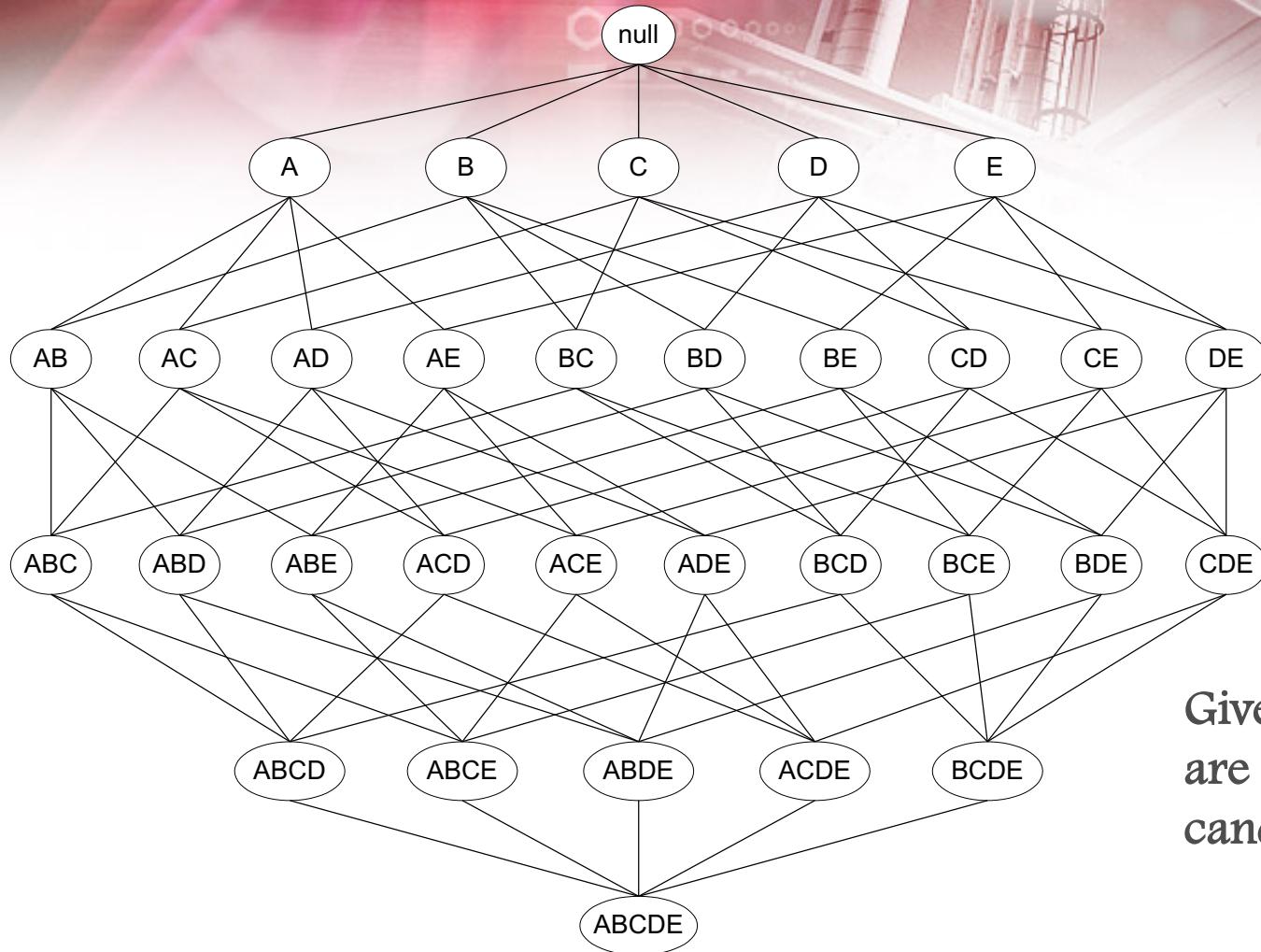
- All the above rules are binary partitions of the same itemset:
 $\{\text{Milk}, \text{Diaper}, \text{Beer}\}$
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

Mining Association Rules



- Two-step approach:
 - **Frequent Itemset Generation**
 - Generate all itemsets whose support $\geq \text{minsup}$
 - **Rule Generation**
 - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally expensive

Frequent Itemset Generation

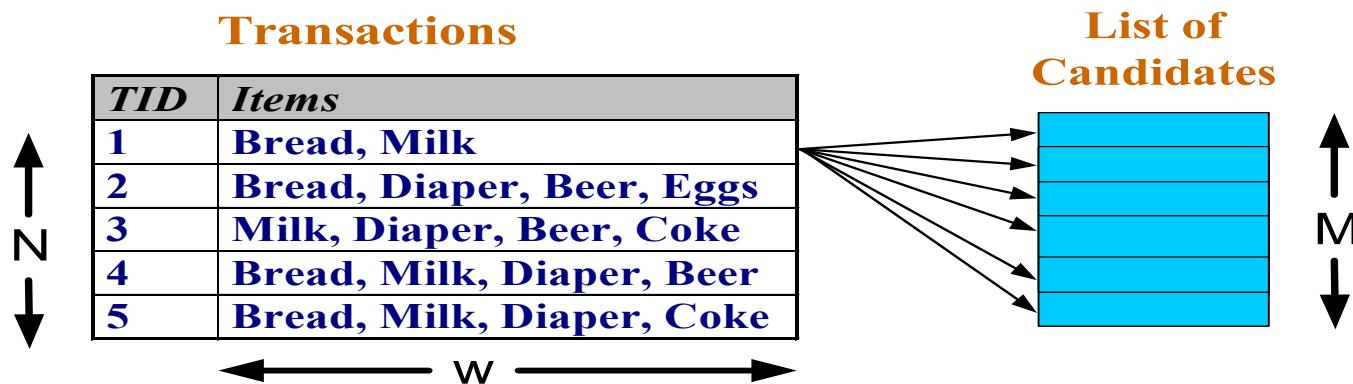


Given d items, there
are 2^d possible
candidate itemsets

Frequent Itemset Generation



- Brute-force approach:
 - Each itemset in the lattice is a candidate frequent itemset
 - Count the support of each candidate by scanning the database

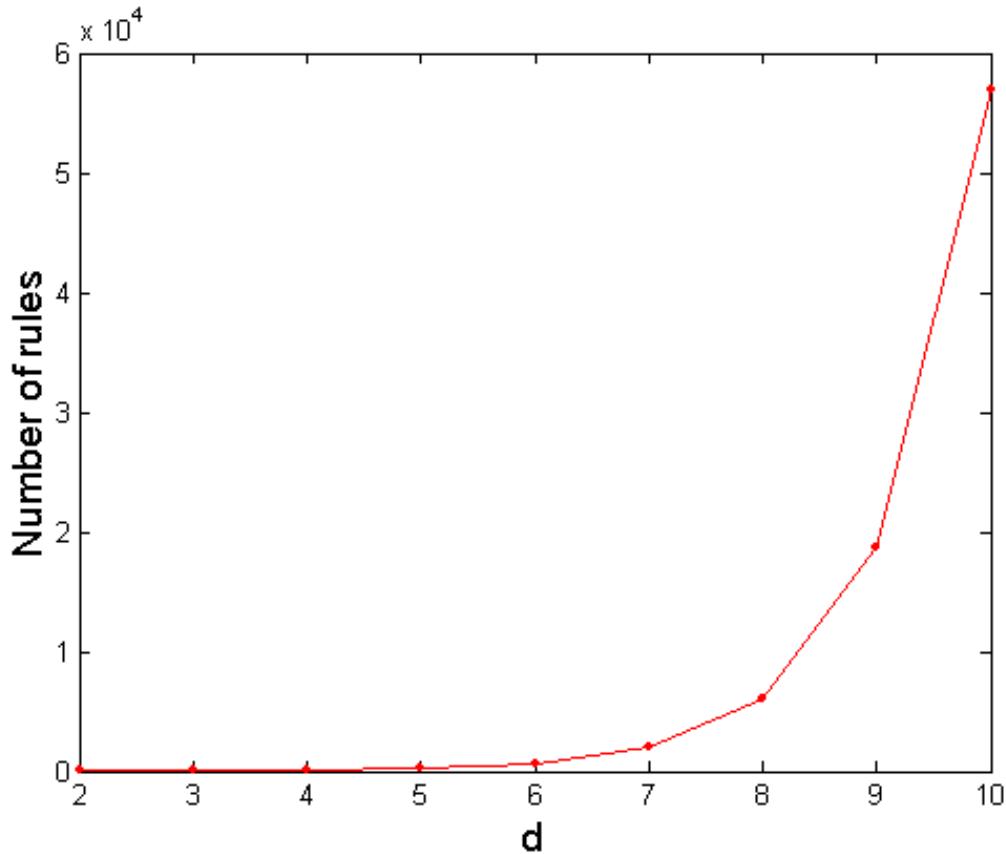


- Match each transaction against every candidate
- Complexity $\sim O(NMw) \Rightarrow$ Expensive since $M = 2^d$!!!

Computational Complexity



- Given d unique items:
 - Total number of itemsets = 2^d
 - Total number of possible association rules:



$$\begin{aligned}R &= \sum_{k=1}^{d-1} \binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \\&= 3^d - 2^{d+1} + 1\end{aligned}$$

If $d=6$, $R = 602$ rules

Frequent Itemset Generation Strategies



- Reduce the number of candidates (M)
 - Complete search: $M = 2^d$
 - Use pruning techniques to reduce M
- Reduce the number of transactions (N)
 - Reduce size of N as the size of itemset increases
 - Used by DHP and vertical-based mining algorithms
- Reduce the number of comparisons (NM)
 - Use efficient data structures to store the candidates or transactions
 - No need to match every candidate against every transaction

Reducing Number of Candidates

- Apriori principle:
 - If an itemset is frequent, then all of its subsets must also be frequent
- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the anti-monotone property of support



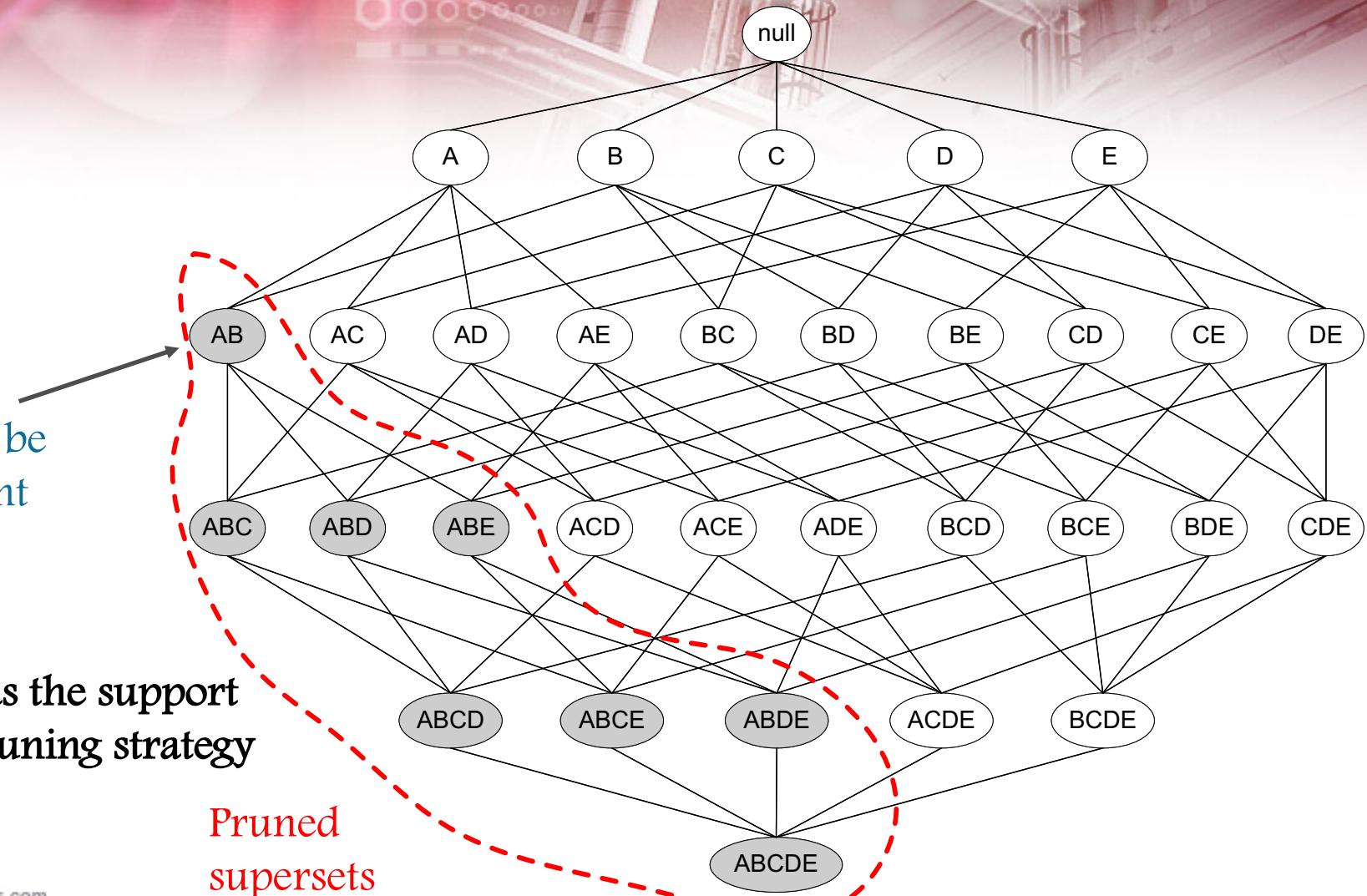
Illustrating Apriori Principle



Found to be
Infrequent

Known as the support
based pruning strategy

Pruned
supersets





Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$$

With support-based pruning,

$$6 + 6 + 1 = 13$$



Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	3



Apriori Algorithm



- Method:
 - Let $k=1$
 - Generate frequent itemsets of length 1
 - Repeat until no new frequent itemsets are identified
 - Generate length $(k+1)$ candidate itemsets from length k frequent itemsets
 - Prune candidate itemsets containing subsets of length k that are infrequent
 - Count the support of each candidate by scanning the DB
 - Eliminate candidates that are infrequent, leaving only those that are frequent

Candidate Generation



- Apriori algorithm generates candidate itemsets by performing:
 1. Candidate generation → generates new candidate k-itemsets based on the frequent (k-1)-itemsets found in the previous iteration
 2. Candidate pruning → eliminates some of the candidate k-itemsets using the support-based strategy pruning
- Consider k-itemset $X = \{i_1, i_2, \dots, i_k\}$. The algorithm must determine whether all of its proper subsets $X - \{i_j\}$ ($\forall j = 1, 2, \dots, k$), are frequent.
- The complexity for this operation is $O(k)$ for each candidate k-itemset.

Candidate Generation (2)



- List of requirement for an effective candidate generation procedure:
 - It should avoid generating too many unnecessary candidates. A candidate itemset is unnecessary if at least one of its subsets is infrequent. For example: $\{a,b\}$ is a candidate 2-itemset if $\{a\}$ is infrequent then $\{a,b\}$ is unnecessary because it is guaranteed to be infrequent by the anti-monotone property of support.
 - It must ensure that the candidate set is complete.
 - It should not generate the same candidate itemset more than once. For example: the candidate itemset $\{a,b,c,d\}$ can be generated in many ways, i.e merging $\{a,b,c\}$ with $\{d\}$, $\{b,d\}$ with $\{a,c\}$, $\{c\}$ with $\{a,b,d\}$

Candidate Generation Procedures



1. Brute-Force Method → considers every k -itemset as a potential candidate then applies candidate pruning step to remove any unnecessary candidate. The number of candidate generated is $\binom{d}{k}$ where d is the total number of items
2. $F_{k-1} \times F_1$ → extend $(k-1)$ -itemset with frequent 1-itemset to form k -itemset. Avoid duplicate candidates by sorting items in each frequent itemset lexicographically. Each frequent $(k-1)$ -itemset X is then extended with frequent items that are lexicographically larger than the items in X . Reduce the number of unnecessary candidate by ensuring that every item in the candidate must be contained in at least $k-1$ of the frequent $(k-1)$ items.
3. $F_{k-1} \times F_{k-1}$ → merges a pair of frequent $(k-1)$ -itemsets only if their first $k-2$ items are identical. Let $A = \{a_1, a_2, \dots, a_{k-1}\}$ and $B = \{b_1, b_2, \dots, b_{k-1}\}$ be a pair of frequent $(k-1)$ -itemsets. A and B are merged if they satisfy the following conditions: $a_i = b_i$ (for $i=1,2,\dots,k-2$) and $a_{k-1} \neq b_{k-1}$

Candidate Generation using $F_{k-1} \times F_1$



Frequent 2-itemset

Itemset
{Beer,Diapers}
{Bread,Diapers}
{Bread, Milk}
{Diapers, Milk}

Frequent 1-itemset

Itemset
Beer
Bread
Diapers
Milk

Candidate generation

Itemset
{Beer, Diapers, Bread}
{Beer, Diapers, Milk}
{Bread, Diapers, Milk}
{Bread, Milk, Beer}

Candidate pruning

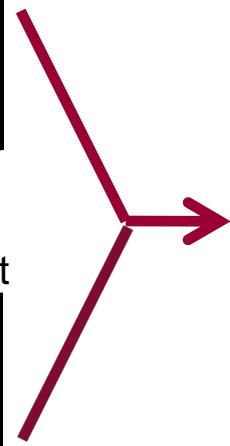
Itemset
{Bread, Diapers, Milk}

Candidate Generation using $F_{k-1} \times F_{k-1}$



Frequent 2-itemset	
	Itemset
	{Beer,Diapers}
	{Bread,Diapers}
	{Bread, Milk}
	{Diapers, Milk}

Frequent 2-itemset	
	Itemset
	{Beer,Diapers}
	{Bread,Diapers}
	{Bread, Milk}
	{Diapers, Milk}



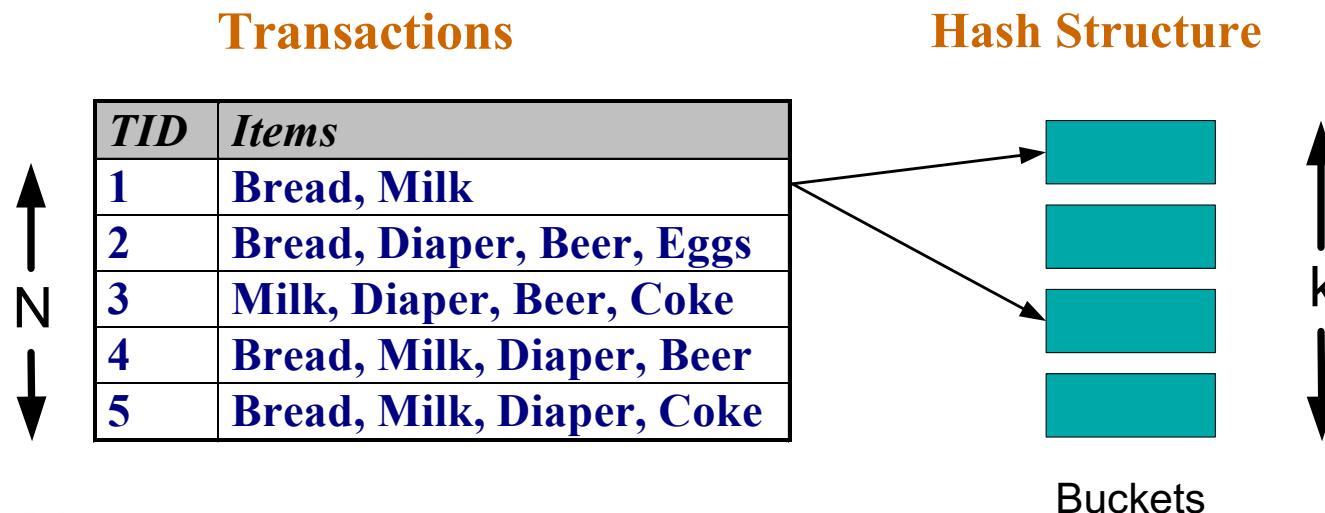
Candidate generation	
	Itemset
	{Bread, Diapers, Milk}

Candidate pruning	
	Itemset
	{Bread, Diapers, Milk}

Reducing Number of Comparisons



- Candidate counting:
 - Scan the database of transactions to determine the support of each candidate itemset
 - To reduce the number of comparisons, store the candidates in a hash structure
 - Instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets



Generate Hash Tree



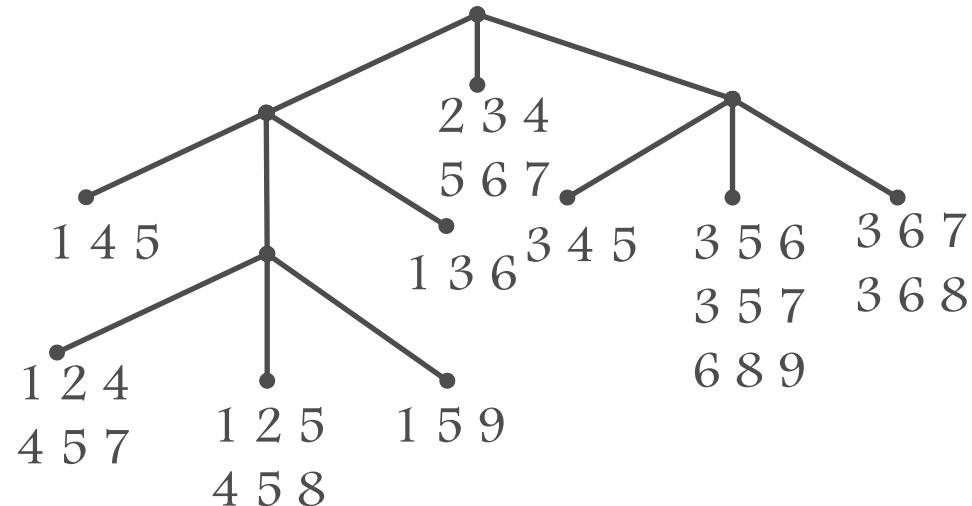
Suppose you have 15 candidate itemsets of length 3:

$\{1\ 4\ 5\}$, $\{1\ 2\ 4\}$, $\{4\ 5\ 7\}$, $\{1\ 2\ 5\}$, $\{4\ 5\ 8\}$, $\{1\ 5\ 9\}$, $\{1\ 3\ 6\}$, $\{2\ 3\ 4\}$, $\{5\ 6\ 7\}$,
 $\{3\ 4\ 5\}$, $\{3\ 5\ 6\}$, $\{3\ 5\ 7\}$, $\{6\ 8\ 9\}$, $\{3\ 6\ 7\}$, $\{3\ 6\ 8\}$

You need:

- Hash function
- **Max leaf size:** max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)

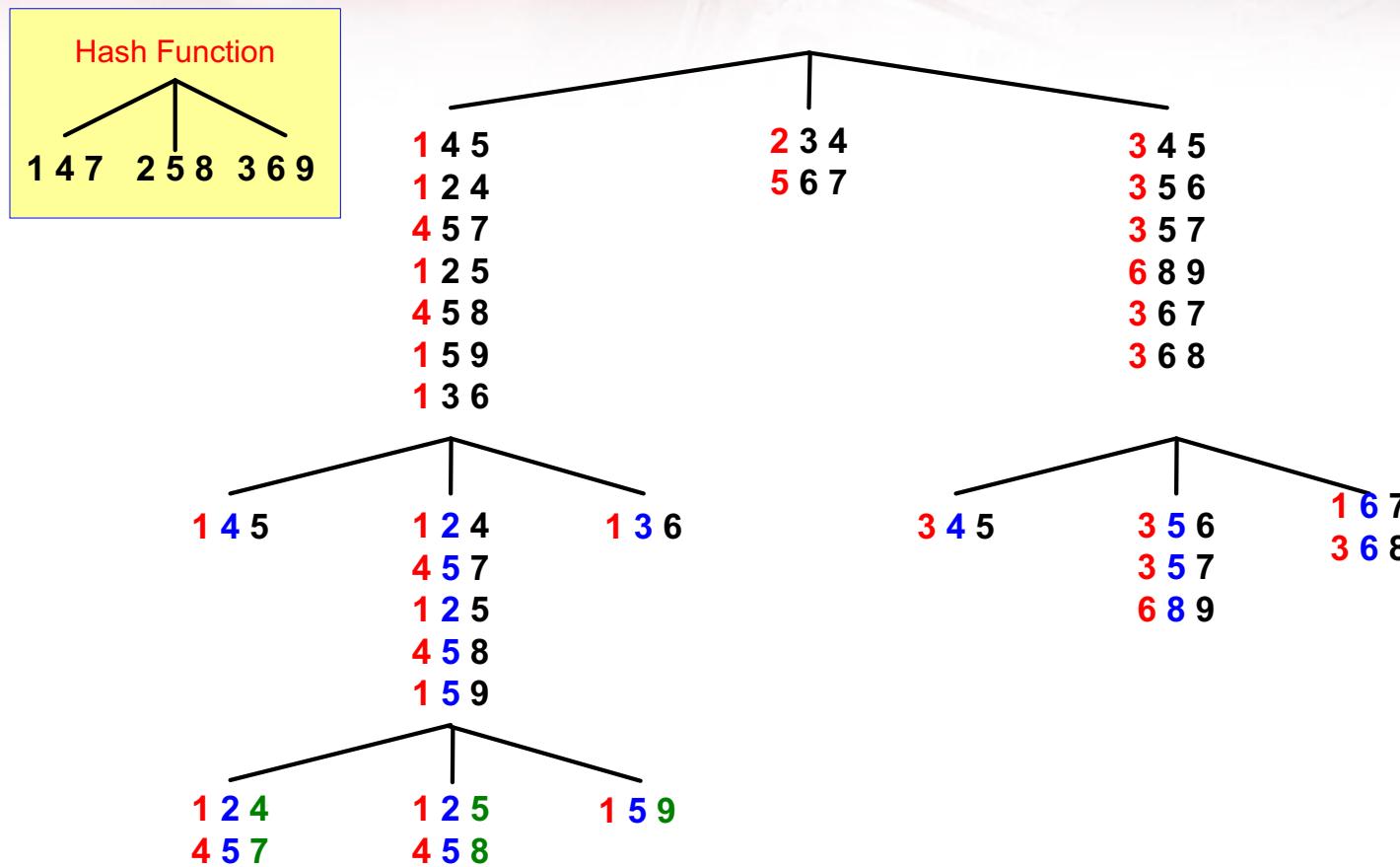
Hash function



Generate Hash Tree



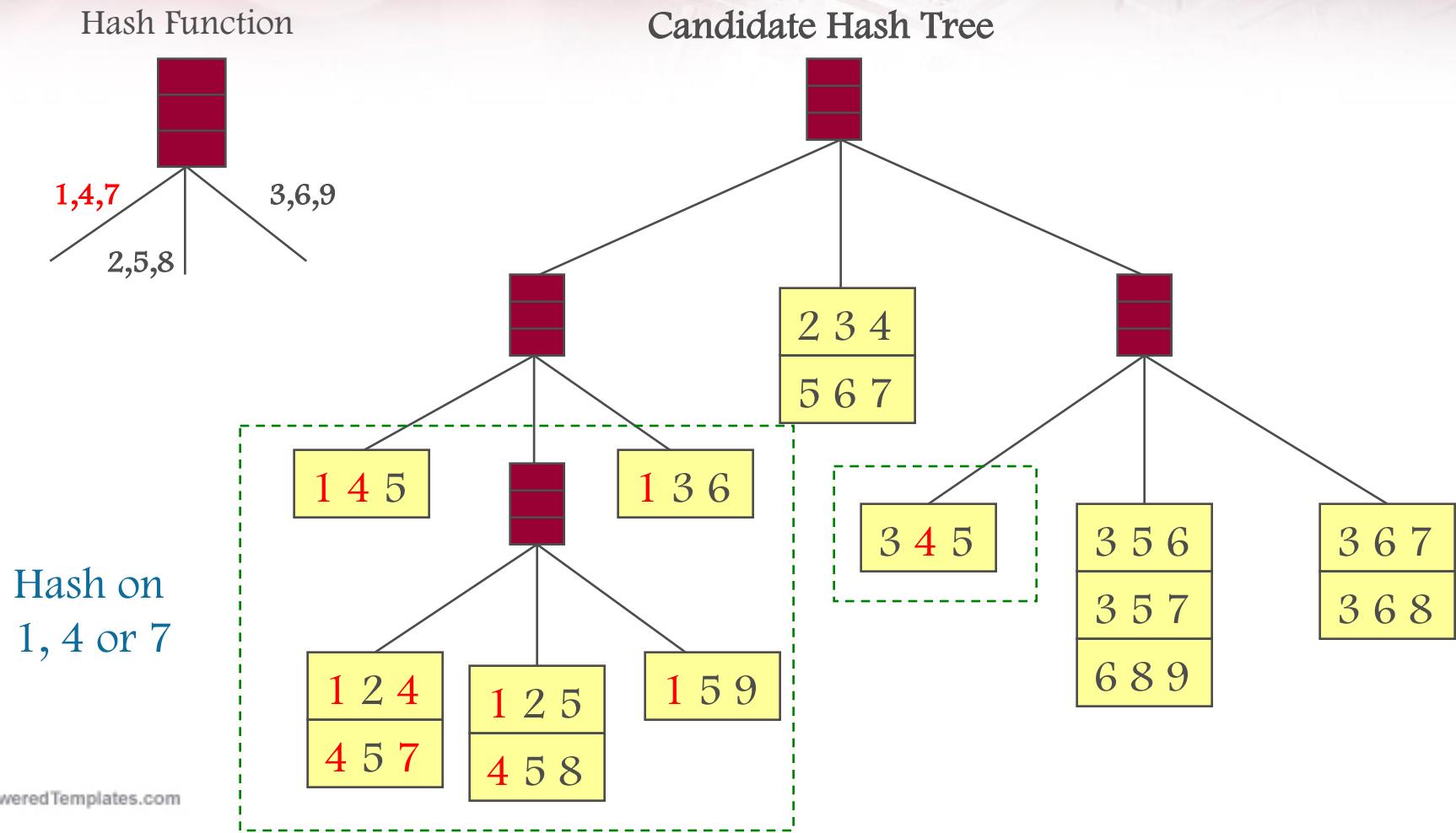
{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4},
{5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}



Association Rule Discovery: Hash tree



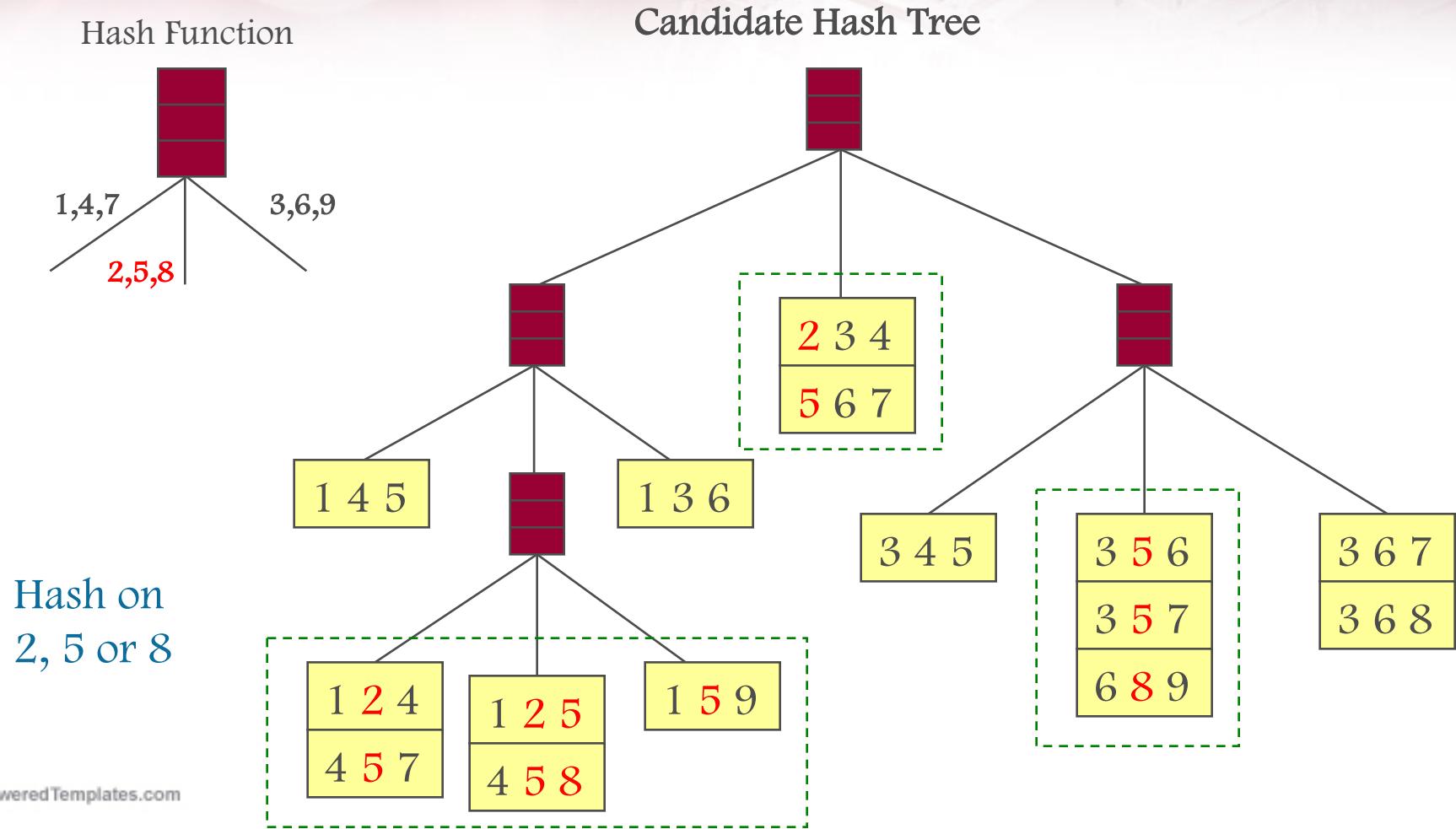
{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4},
{5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}



Association Rule Discovery: Hash tree



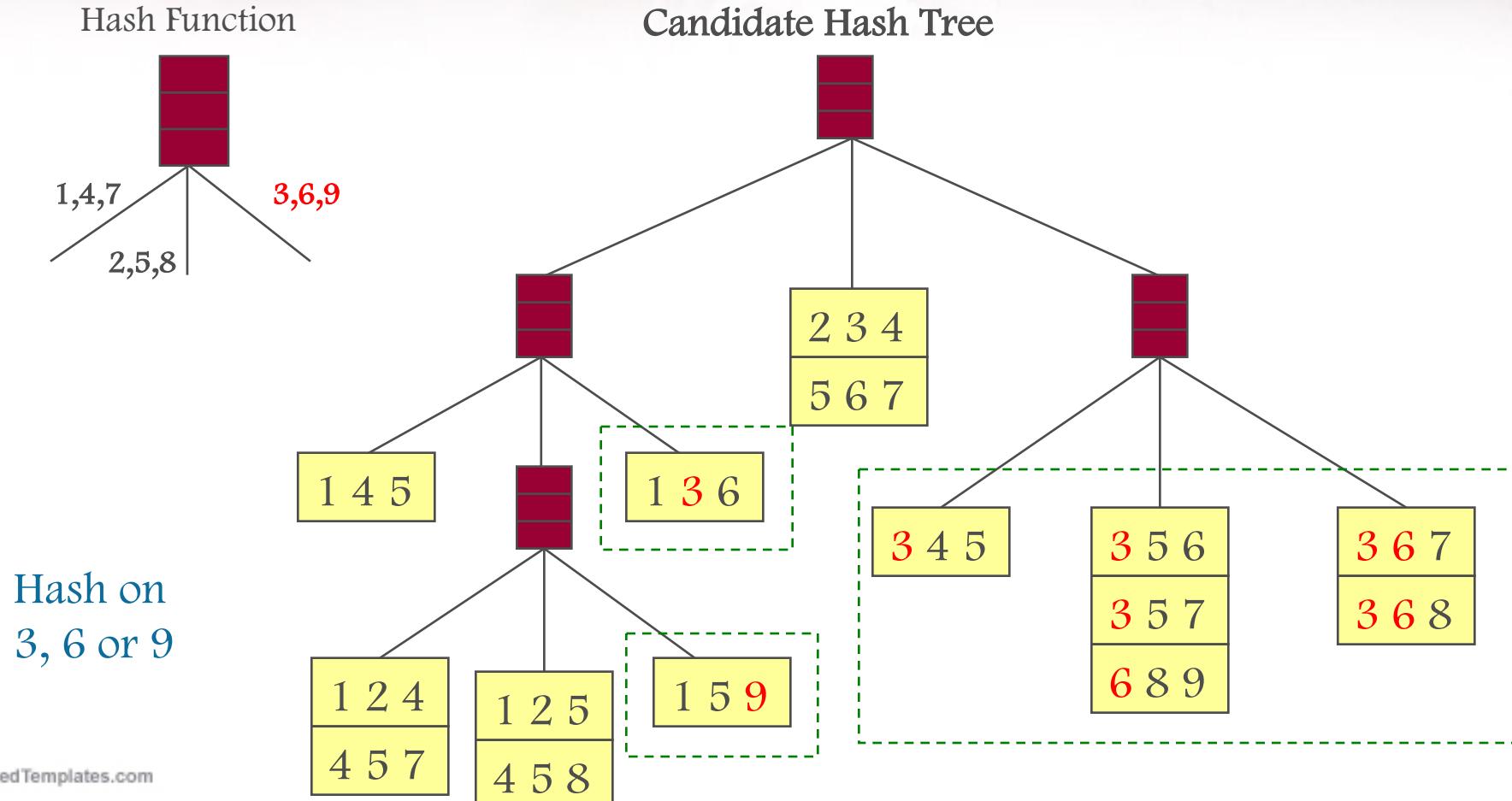
{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4},
{5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}



Association Rule Discovery: Hash tree



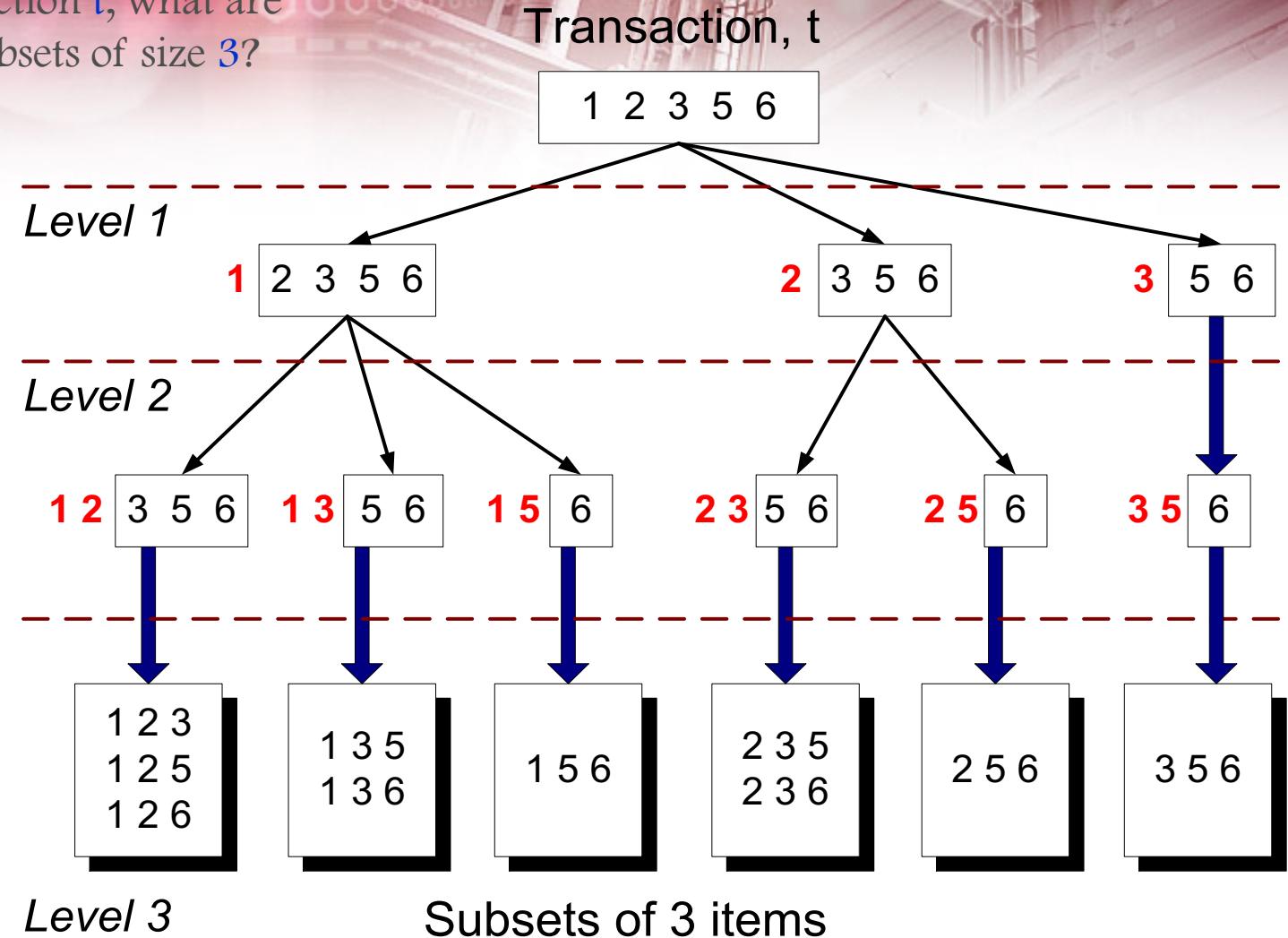
{1 4 5}, {1 2 4}, {4 5 7}, {1 2 5}, {4 5 8}, {1 5 9}, {1 3 6}, {2 3 4},
{5 6 7}, {3 4 5}, {3 5 6}, {3 5 7}, {6 8 9}, {3 6 7}, {3 6 8}



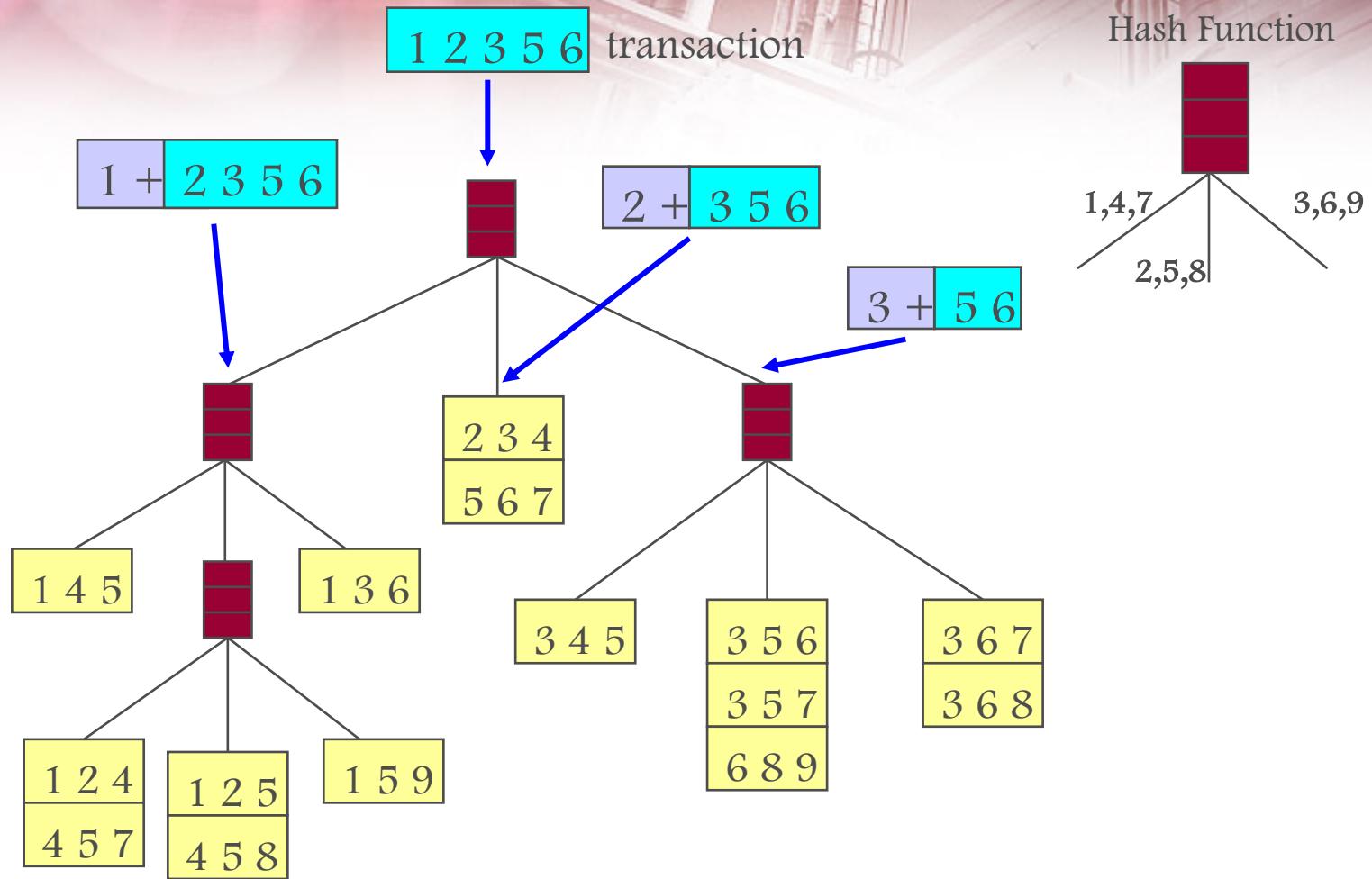
Subset Operation



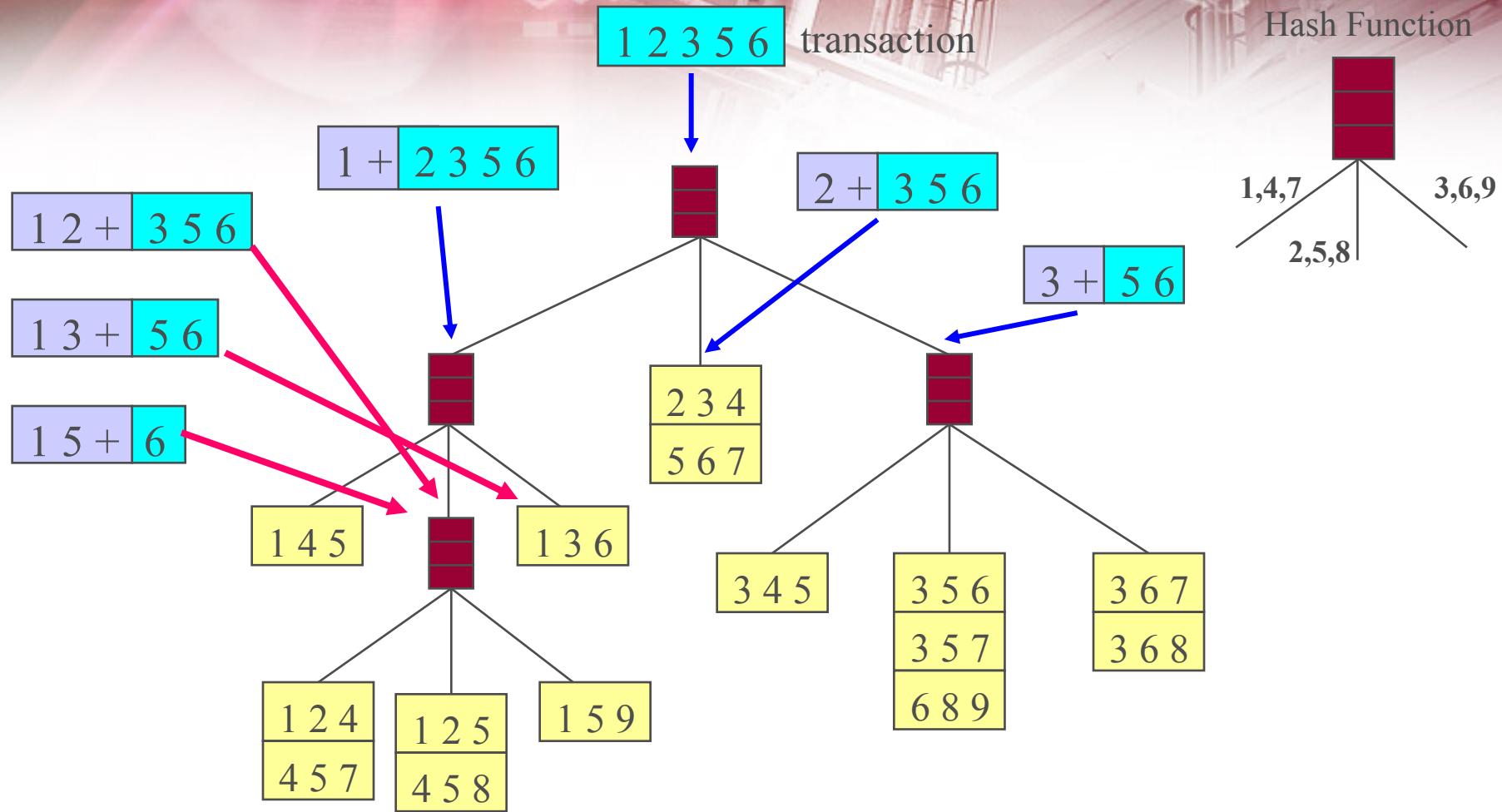
Given a transaction t , what are the possible subsets of size 3?



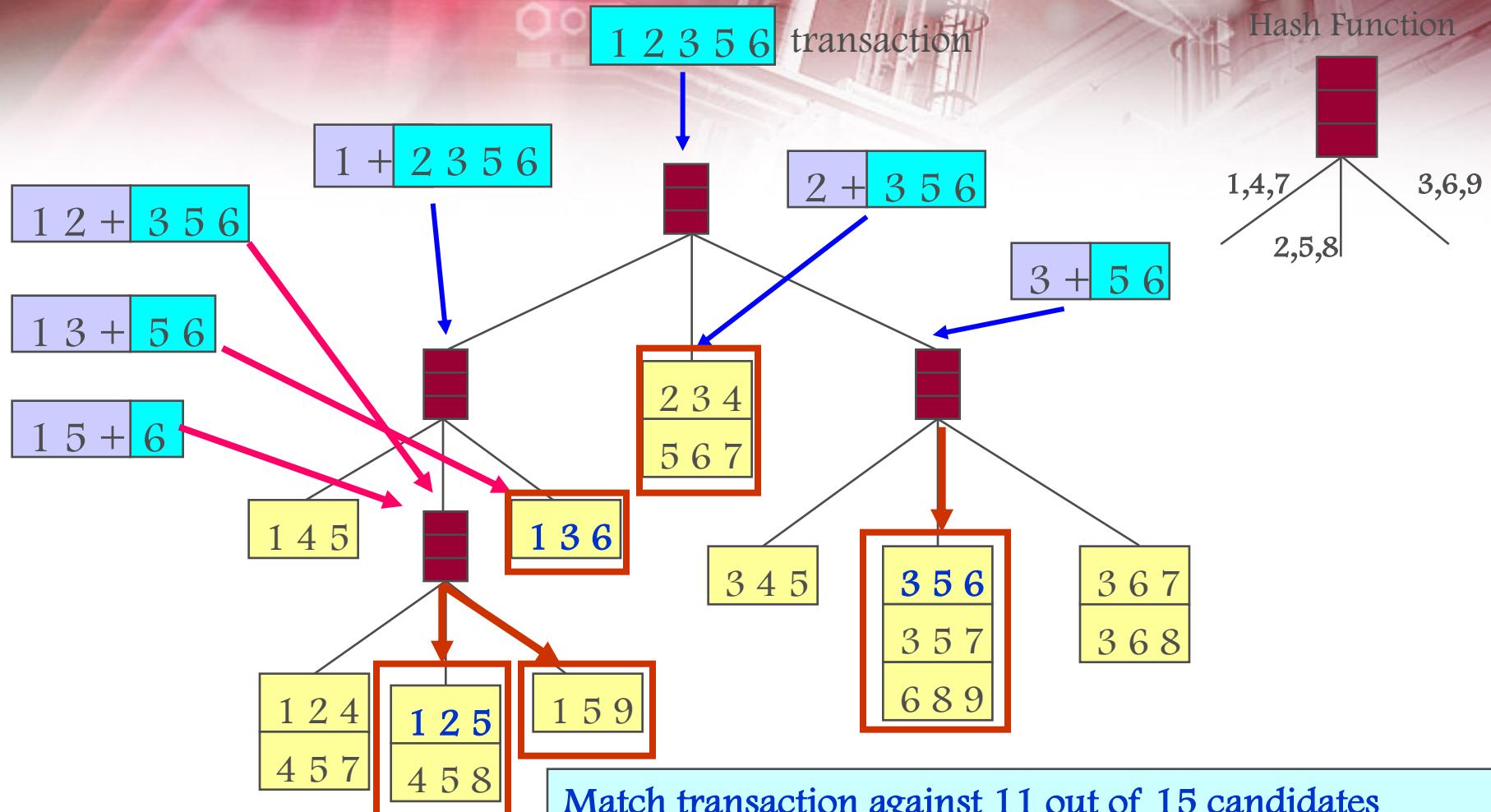
Subset Operation Using Hash Tree



Subset Operation Using Hash Tree



Subset Operation Using Hash Tree



Factors Affecting Complexity



- Choice of minimum support threshold
 - Lowering support threshold results in more frequent itemsets
 - This may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
 - More space is needed to store support count of each item
 - If number of frequent items also increases, both computation and I/O costs may also increase
- Size of database
 - Since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- Average transaction width
 - Transaction width increases with denser data sets
 - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

Compact Representation of Frequent Itemsets



- Some itemsets are redundant because they have identical support as their supersets

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	

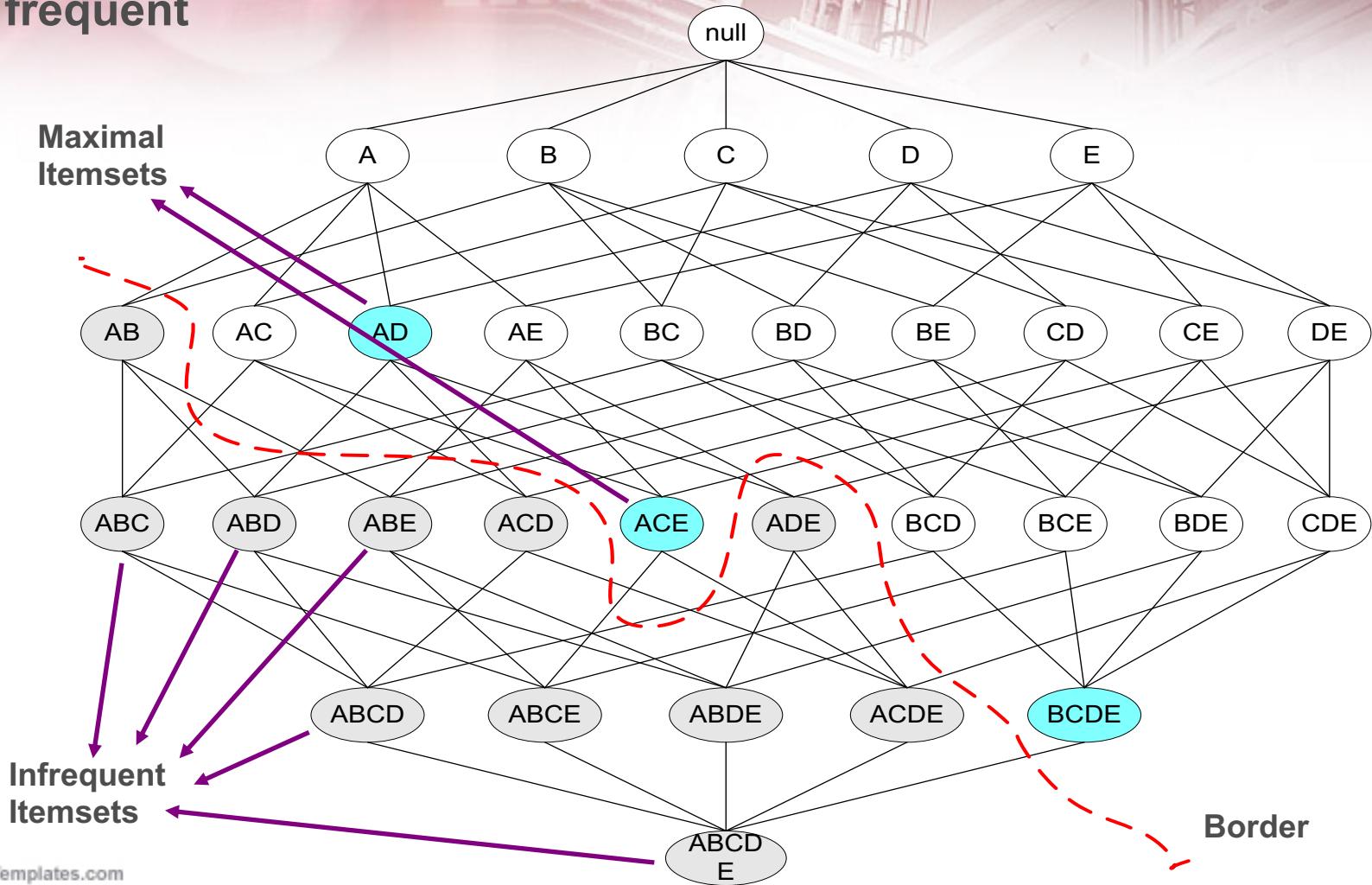
- Number of frequent itemsets
- Need a compact representation

$$= 3 \times \sum_{k=1}^{10} \binom{10}{k}$$

Maximal Frequent Itemset



An itemset is **maximal frequent** if none of its immediate supersets is frequent



Closed Itemset



- An itemset is *closed* if none of its immediate supersets has the same support as the itemset, *closed frequent* if it is closed and its support is greater than or equal to minimum support

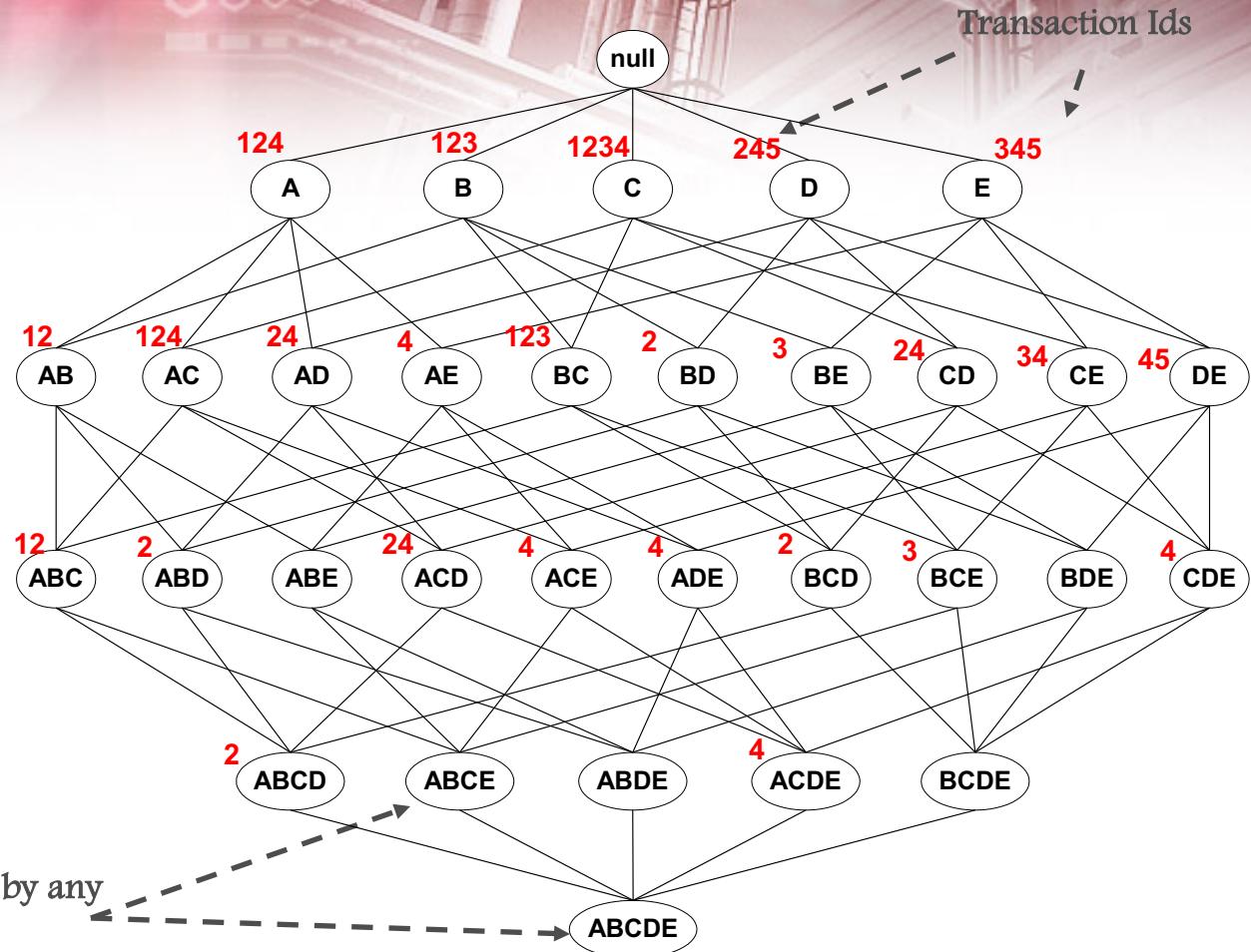
TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

Itemset	Support
{A,B,C}	2
{A,B,D}	3
{A,C,D}	2
{B,C,D}	3
{A,B,C,D}	2

Maximal vs Closed Itemsets

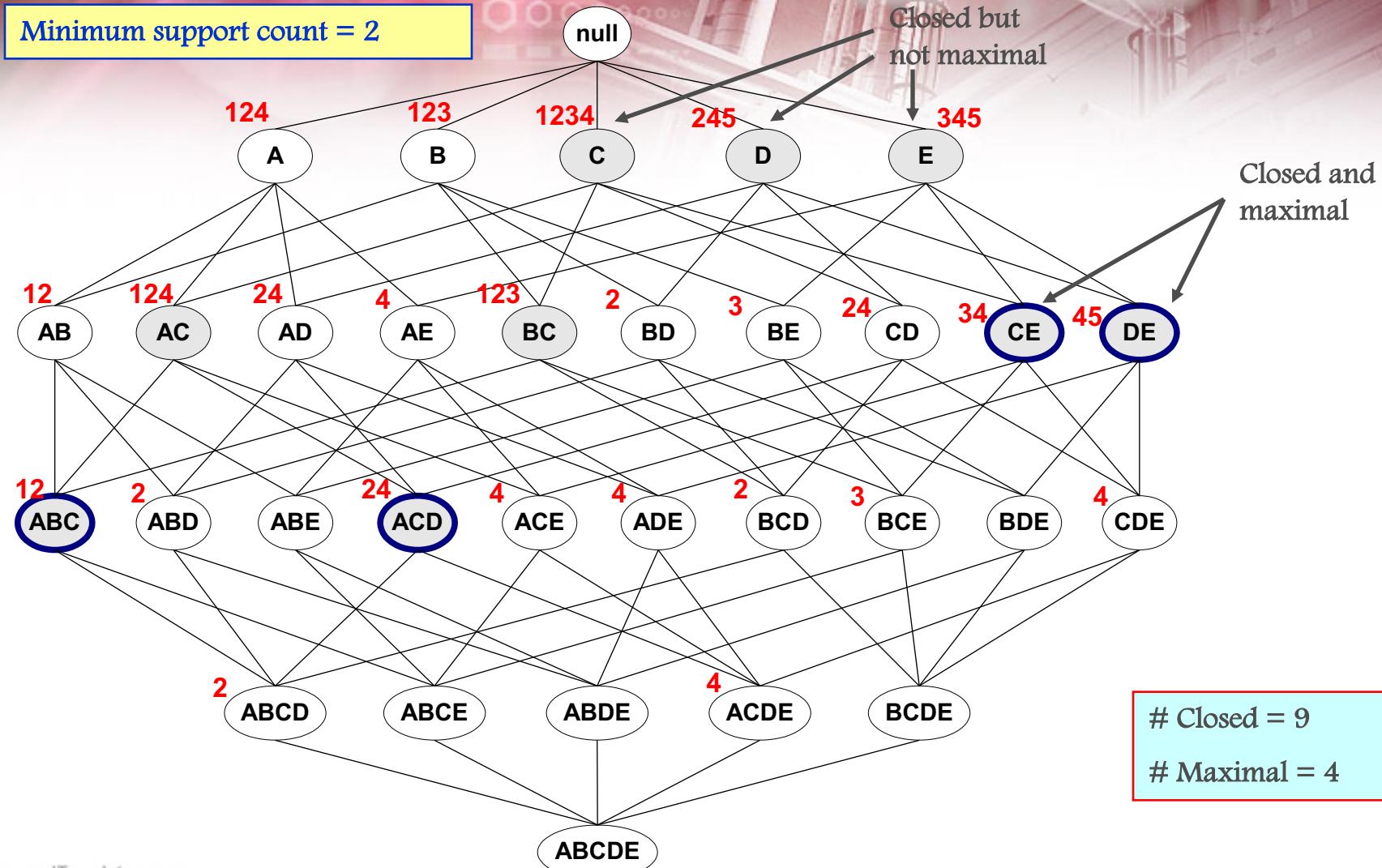
TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE



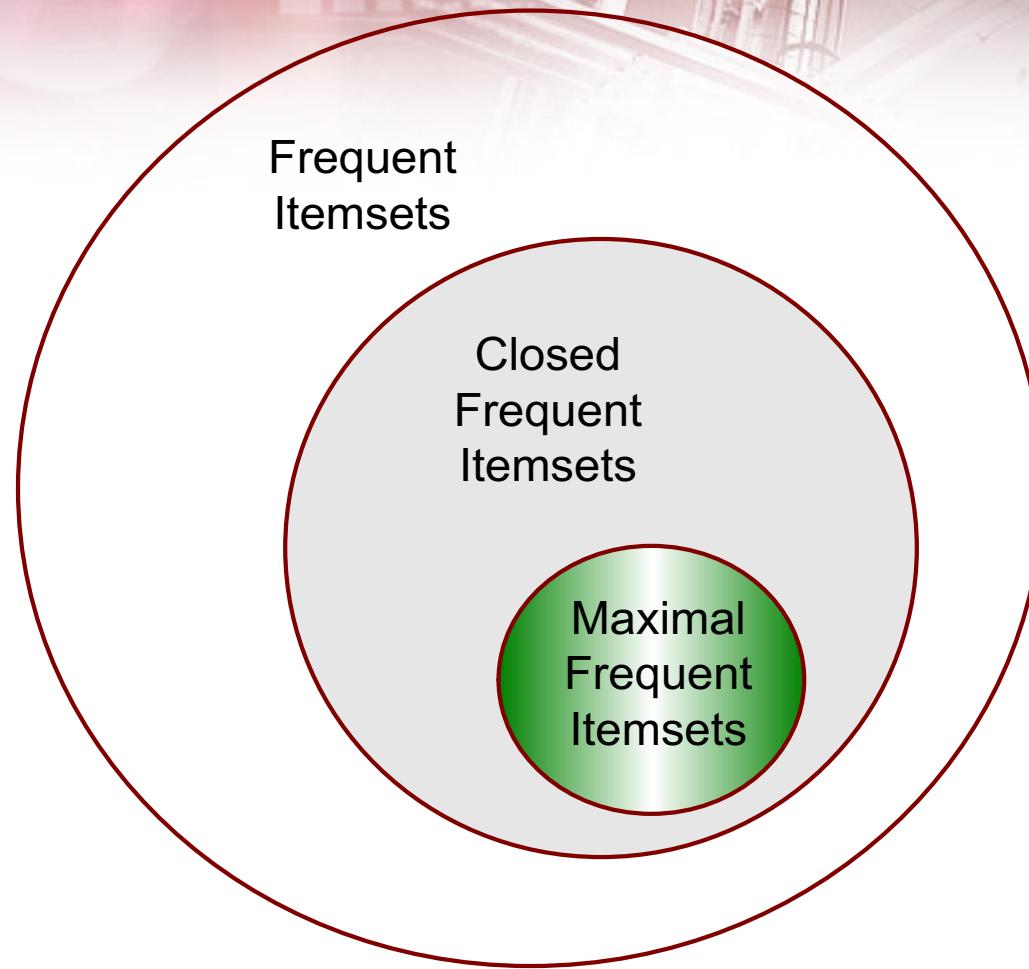
Maximal vs Closed Frequent Itemsets



Minimum support count = 2



Maximal vs Closed Itemsets



Alternative Methods for Frequent Itemset Generation

- Representation of Database
 - horizontal vs vertical data layout

Horizontal
Data Layout

TID	Items
1	A,B,E
2	B,C,D
3	C,E
4	A,C,D
5	A,B,C,D
6	A,E
7	A,B
8	A,B,C
9	A,C,D
10	B

Vertical Data Layout

A	B	C	D	E
1	1	2	2	1
4	2	3	4	3
5	5	4	5	6
6	7	8	9	
7	8	9		
8	10			
9				

FP-growth Algorithm



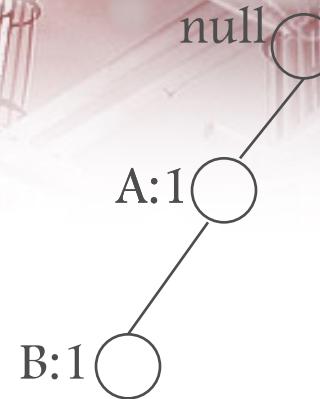
- Use a compressed representation of the database using an FP-tree
- Once an FP-tree has been constructed, it uses a recursive divide-and-conquer approach to mine the frequent itemsets

FP-tree construction

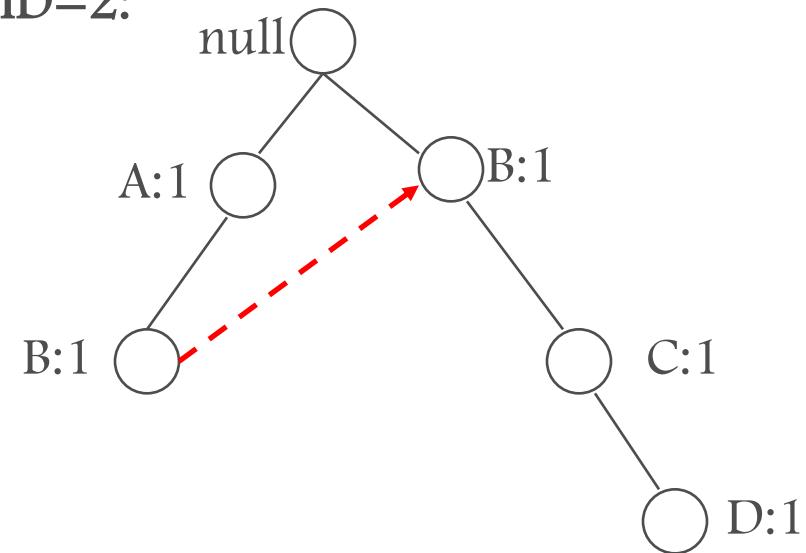


TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

After reading TID=1:



After reading TID=2:



FP-Tree Construction

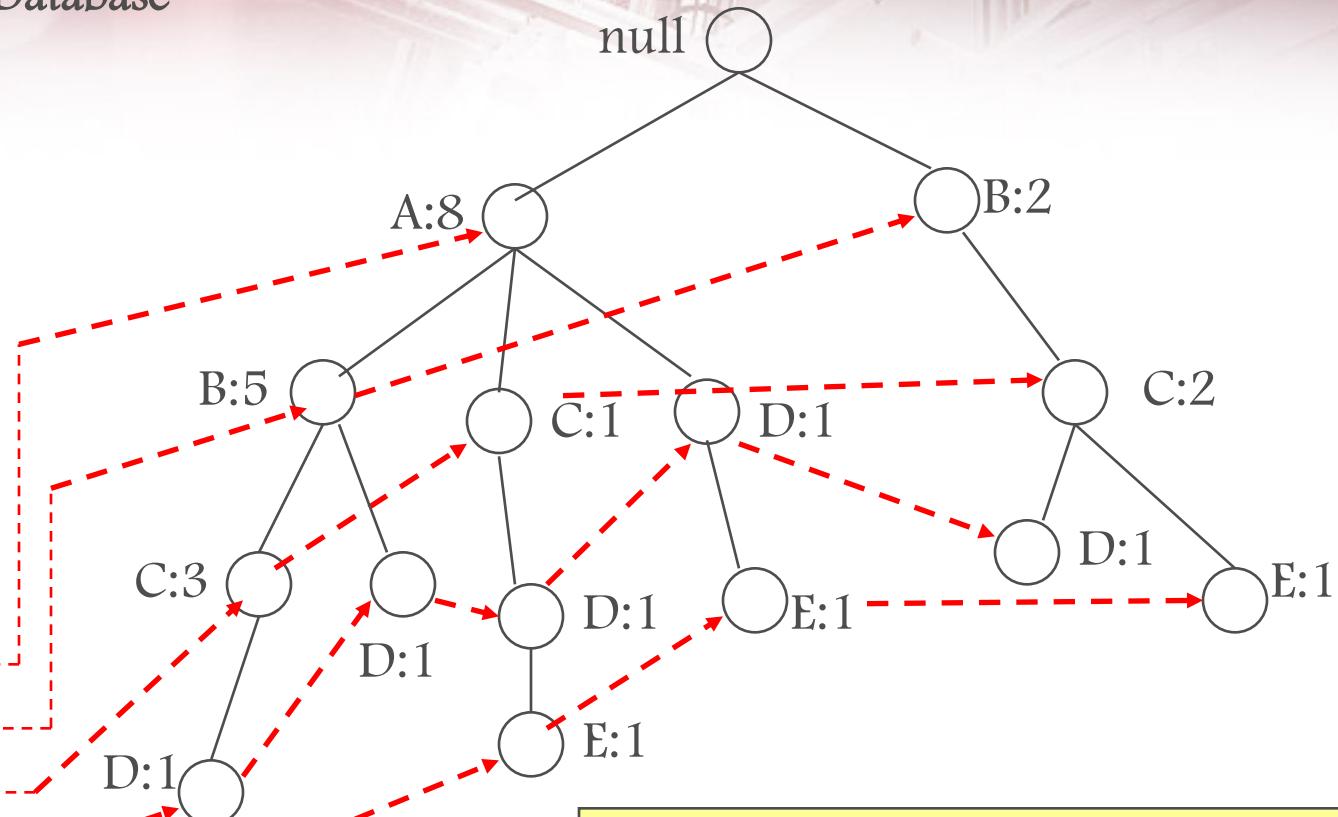
TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{A}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Transaction Database

Pointers are used to assist frequent itemset generation

Header table

Item	Pointer
A	
B	
C	
D	
E	



Items in header table are sorted in decreasing order of the support count of the items

FP-growth



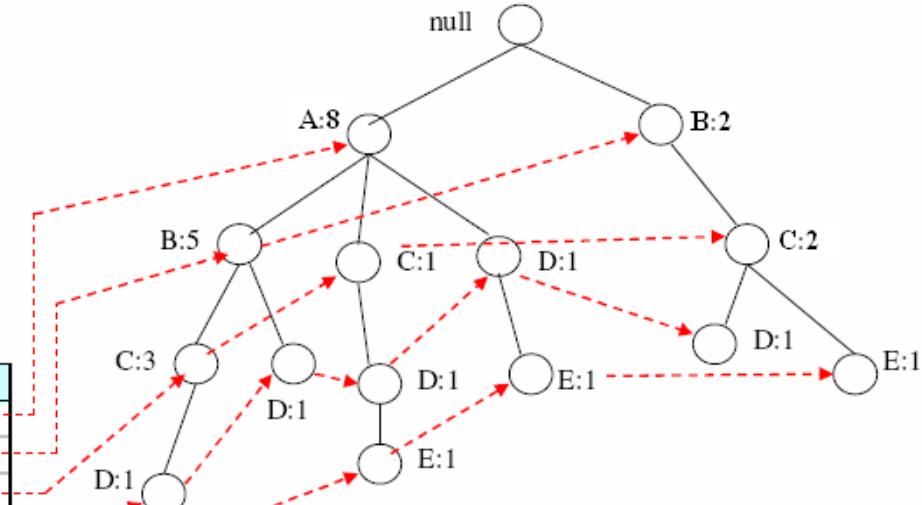
A
B
AB

C
BC, AC
ABC

D
CD, BD, AD
BCD, ACD, ABD
ABCD

E
DE, CE, BE, AE
CDE, BDE, ADE, BCE, ...
BCDE, ACDE, ABCE, ...

Item	Pointer
A	
B	
C	
D	
E	

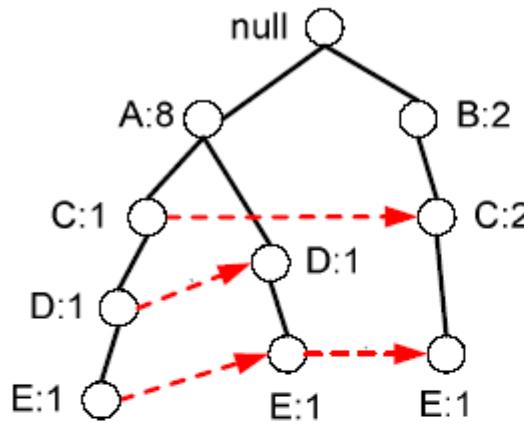


Generate frequent itemsets using divide-and-conquer approach

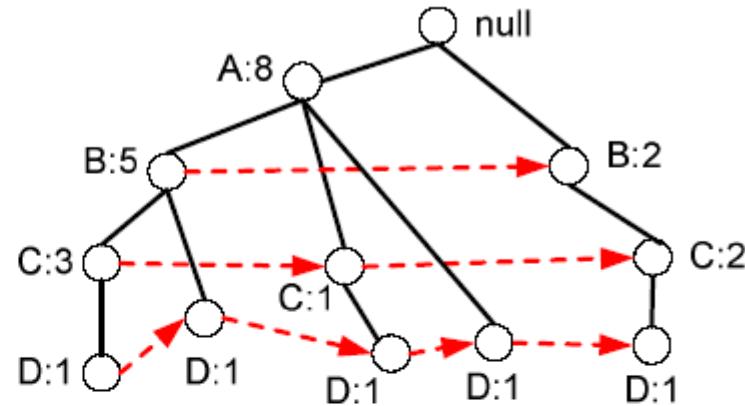
FP-growth



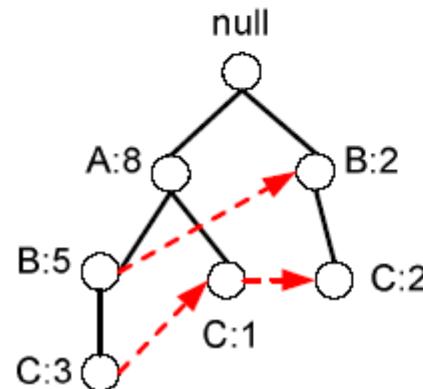
(a) Paths containing node E



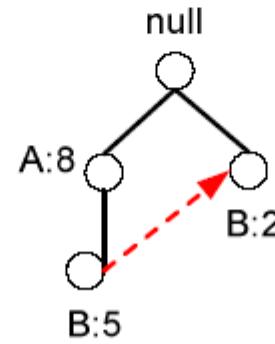
(b) Paths containing node D



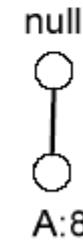
(c) Paths containing node C



(d) Paths containing node B



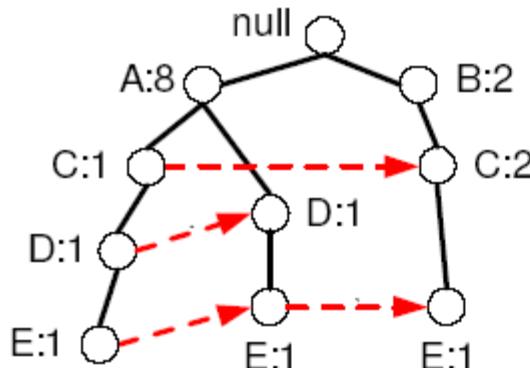
(e) Paths containing node A



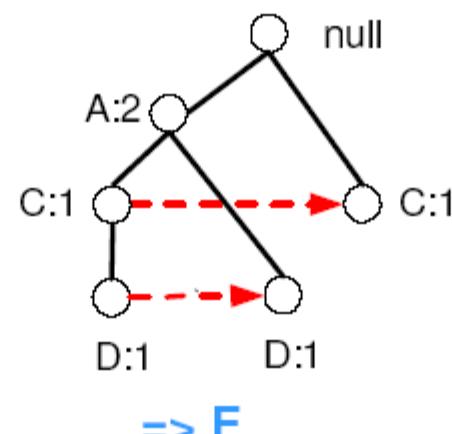
Conditional FP-Tree



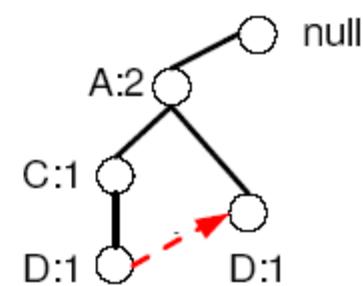
(a) Prefix paths ending in E



(b) Conditional FP-tree for E



(c) Prefix paths ending in DE

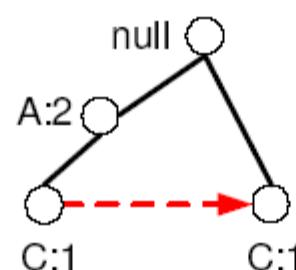


(d) Conditional FP-tree for DE

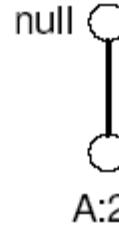


=> DE, ADE

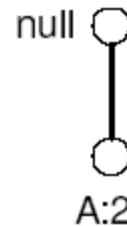
(e) Prefix paths ending in CE



(f) Conditional FP-tree for CE



(g) Prefix paths ending in AE



=> AE