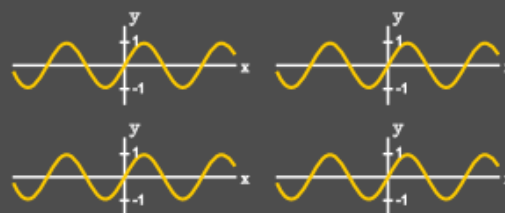


CAN YOU FIGURE OUT THESE MOVIE TITLES?

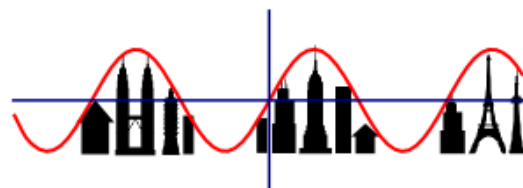
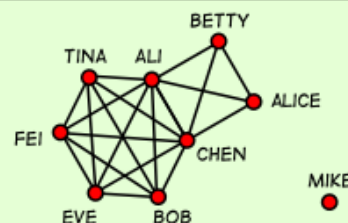
$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$



$$B_{\text{🔥}}(p) = \{x \in M \mid d(x, p) < \text{🔥}\}$$

$$\sup\{\text{🦆}\}$$

Fe X Fe



$$\frac{\partial u}{\partial t} - \alpha \nabla^2 u = 0$$

$$a+bi$$

$$e^{i\pi} + 1 = 0$$

and

$$666$$

CREATED BY SPIKEDMATH.COM

spikedmath.com
© 2011

Data Mining: Association Analysis

Laura Brown

Some slides adapted from G. Piatetsky-Shapiro;
Han, Kamber, & Pei; Tan, Steinbach, & Kumar

Association Rule Mining

- Given a set of transaction, 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

{Diaper} -> {Beer},
{Milk, Bread} -> {Eggs,Coke},
{Beer, Bread} -> {Milk},

Implication means co-occurrence,
not causality!

Applications

- Example 1 – text mining
 - baskets = sentences
 - items = words in those sentences
 - find words that appear together unusually frequently, i.e., linked concepts
- Example 2 – document mining
 - baskets = sentences
 - items = documents containing those sentences
 - items that appear together too often could represent plagiarism

Applications

- Example 3 – healthcare mining
 - baskets = people
 - items = genes or blood-chemistry factors
 - detect combinations of genes that results in a disease
 - requires extension: absence of an item needs to be observed as well as presence

Terminology

Association Analysis: Frequent itemset mining

Terminology

- Itemsets – a set of items (collection of one or more items, $X \subseteq I$)
 - *Items* : $I = \{x_1, x_2, \dots, x_m\}$
 - *k-itemset* – an itemset with k items
 - Ex. $X = \{Milk, Bread, Beer\}$
- Tidsets – a set of *tids*, $T \subseteq \mathcal{T}$
 - *Transaction identifiers* or *tids* : $\mathcal{T} = \{t_1, t_2, \dots, t_n\}$
- Transaction – a transaction is a tuple of the form $\langle t, X \rangle$, where $t \in \mathcal{T}$ is a unique *tid* and X is an itemset

Database

D	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

Binary Database

<i>t</i>	i (<i>t</i>)
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

Transaction Database

t (<i>x</i>)				
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
1	1	2	1	1
3	2	4	3	2
4	3	5	5	3
5	4	6	6	4
	5			5
	6			

Vertical Database

- The database **D** has 5 items, $I = \{A, B, C, D, E\}$ and 6 *tids* $\mathcal{T} = \{1, 2, 3, 4, 5, 6\}$
- When describing a transaction we can drop the set notation: $\langle 1, \{A, B, D, E\} \rangle \rightarrow \langle t, ABDE \rangle$

Support and Frequent Itemsets

- The **support** of an itemset X , $sup(X)$ or $\sigma(X)$, in \mathbf{D} is the number of transactions in \mathbf{D} that contain X
 - Ex. $sup(\{\text{Milk, Bread, Diaper}\}) = 2$

<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

- The **relative support** of an itemset X , $rsup(X)$ or $s(X)$, is the fraction of transactions that contain X
$$rsup(X) = \frac{sup(X)}{|\mathbf{D}|}, \quad \text{Ex. } rsup(\{\text{Milk, Bread, Diaper}\}) = 2 / 5$$
- An itemset X is **frequent** in \mathbf{D} if $sup(X) \geq minsup$, where $minsup$ is a user defined *minimum support threshold*

Frequent Itemsets

minsup = 3

t	$i(t)$
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

Transaction Database

sup	itemsets
6	<i>B</i>
5	<i>E, BE</i>
4	<i>A, C, D, AB, AE, BC, BD, ABE</i>
3	<i>AD, CE, DE, ABD, ADE, BCE, BDE, ABDE</i>

Frequent Itemsets

- The set \mathcal{F} is the set of all frequent itemsets, and $\mathcal{F}^{(k)}$ is the set of frequent k -itemsets
- The 19 frequent itemsets comprise the set \mathcal{F}

$$\mathcal{F}^{(1)} = \{A, B, C, D, E\}$$

$$\mathcal{F}^{(2)} = \{AB, AD, AE, BC, BD, BE, CE, DE\}$$

$$\mathcal{F}^{(3)} = \{ABD, ABE, ADE, BCE, BDE\}$$

$$\mathcal{F}^{(4)} = \{ABDE\}$$

Example: Frequent Itemsets

- Items = { milk, coke, pepsi, beer, juice }
- MinSupport = 3 baskets

$$t_1 = \{m, c, b\}$$

$$t_2 = \{m, p, j\}$$

$$t_3 = \{m, b\}$$

$$t_4 = \{c, j\}$$

$$t_5 = \{m, p, b\}$$

$$t_6 = \{m, c, b, j\}$$

$$t_7 = \{c, b, j\}$$

$$t_8 = \{b, c\}$$

- Frequent itemsets:

Example: Frequent Itemsets

- Items = { milk, coke, pepsi, beer, juice }
- MinSupport = 3 baskets

$$t_1 = \{m, c, b\}$$

$$t_2 = \{m, p, j\}$$

$$t_3 = \{m, b\}$$

$$t_4 = \{c, j\}$$

$$t_5 = \{m, p, b\}$$

$$t_6 = \{m, c, b, j\}$$

$$t_7 = \{c, b, j\}$$

$$t_8 = \{b, c\}$$

- Frequent itemsets:
 - $\{m\}, \{c\}, \{b\}, \{j\}, \{m, b\}, \{b, c\}, \{c, j\}$
 - $\mathcal{F}^{(1)} = \{m, c, b, j\}$
 - $\mathcal{F}^{(2)} = \{mb, bc, cj\}$

Association Rules

- An **association rule** is an expression of the form

$$X \rightarrow Y$$

where X and Y are disjoint itemsets. Denote $X \cup Y$ as XY

Ex. $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

- The **support** of a rule is the number of transactions in which X and Y co-occur

$$s = \text{sup}(X \rightarrow Y) = \text{sup}(XY)$$

- The **relative support** of a rule is the fraction of transactions in which X and Y co-occur

$$\text{rsup}(X \rightarrow Y) = \text{sup}(XY) / |\mathbf{D}| = P(X \wedge Y)$$

- The **confidence** of a rule is the conditional probability that a transaction contains Y given that it contains X

$$c = \text{conf}(X \rightarrow Y) = P(Y \mid X) = P(X \wedge Y) / P(X) = \text{sup}(XY) / \text{sup}(X)$$

Frequent Itemset Mining

Association Analysis

Mining Association Rules

Two-step approach

1. Frequent Itemset Generation

generate all itemsets whose support $\geq \textit{minsup}$

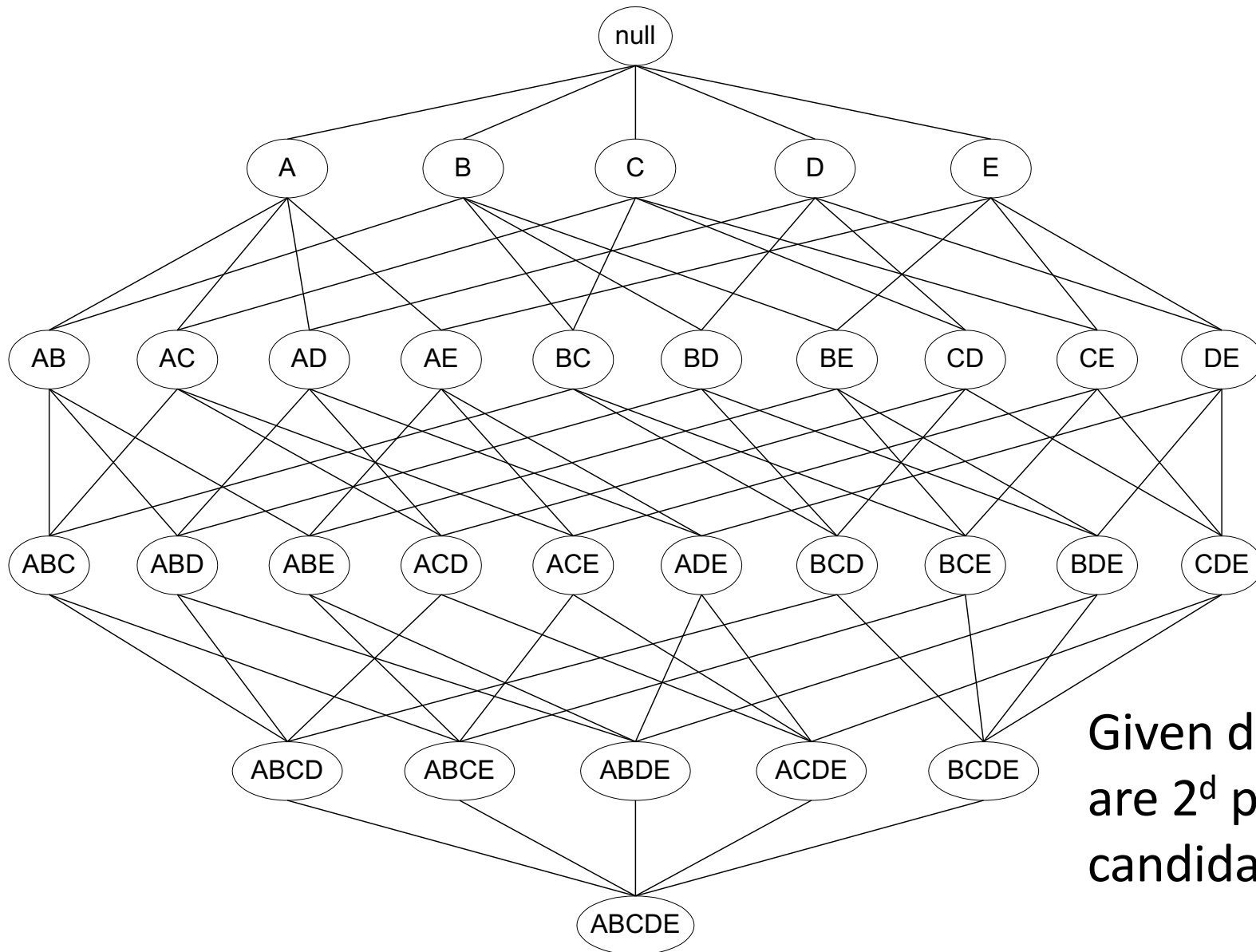
2. Rule Generation

generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

Frequent Itemset Generation

- Frequent itemset generation is computationally expensive
- How many itemsets are potentially to be generated in the worst case?
 - number is sensitive to the *minsup* threshold
 - when *minsup* is low, there exists potentially an exponential number of frequent itemsets

Frequent Itemset Generation



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

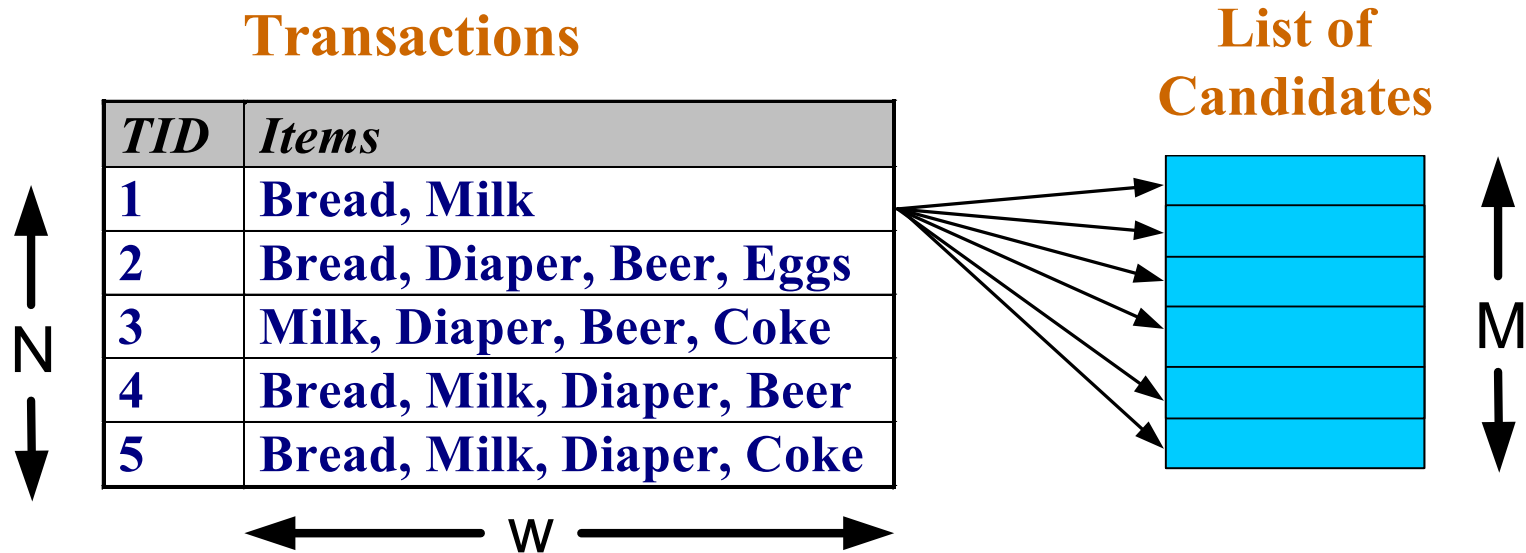
Subset Property

- Every subset of a frequent set is frequent!
 - If {A, B} is frequent. Each occurrence of A, B includes both A and B, then both A and B alone must also be frequent
- A long pattern (itemsets) contains a combinatorial number of sub-patterns (itemsets)
 - A frequent set with 100 items contains

$$\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1$$

Brute Force Algorithm

- 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 $\sim O(NMw)$ -> expensive $M = 2^d$

Frequent Itemset Generation Strategies

- Reduce the **number of candidates** (M)
 - complete search: $M=2^d$
 - use pruning methods 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 candidates or transactions
 - no need to match every candidate against every transaction

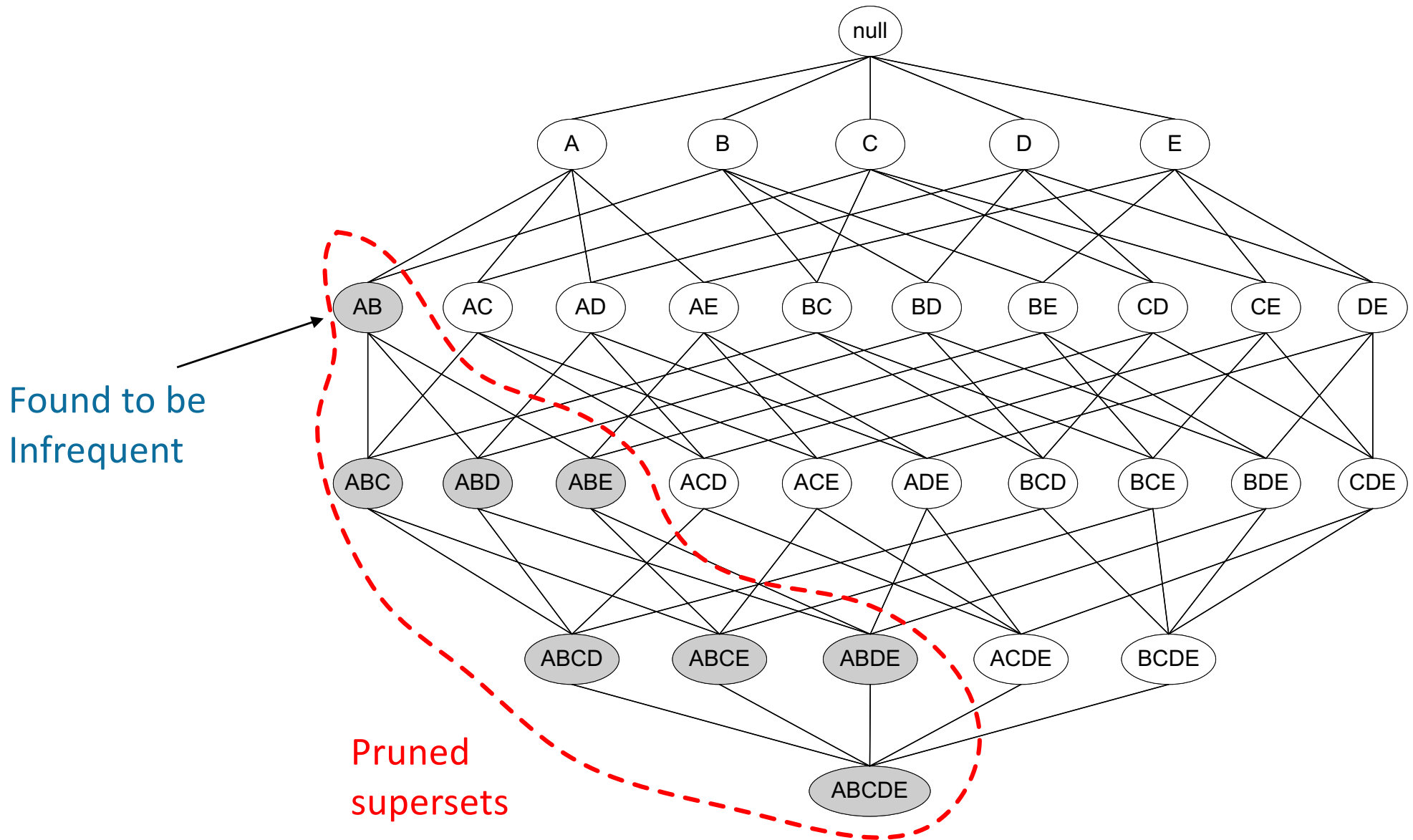
Apriori Frequent Itemset Mining

Association Analysis

Apriori

- Apriori principle:
 1. If an itemset X is frequent, then all of its subsets Y , $Y \subseteq X$ must also be frequent
 2. If an itemset X is not frequent, its supersets Y , $Y \supseteq X$ are also not frequent
- Apriori uses these principles to improve the efficiency of the search for frequent itemsets
- Anti-monotone property of support
 - support of an itemset never exceeds the support of its subsets
$$\forall X, Y : (X \subseteq Y) \Rightarrow \text{sup}(X) \geq \text{sup}(Y)$$

Illustrating Apriori Principle



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:

$$C_{6,1} + C_{6,2} + C_{6,3} = 41$$

With support-based pruning

$$6 + 6 + 1 = 13$$



Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	3

Apriori Pseudocode

C_k : Candidate itemset of size k

L_k : frequent itemset of size k

C_1 = scan database to find support for each item;

L_1 = {frequent items};

for ($k = 1$; $L_k \neq \text{emptyset}$; $k++$) **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

L_{k+1} = candidates in C_{k+1} with min_support

end

return $\bigcup_k L_k$;

Apriori Algorithm - Example

$\text{Sup}_{\min} = 2$

Database TDB

Tid	Items
10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E

1st scan

C_1

Itemset	sup
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

L_1

Itemset	sup
{A}	2
{B}	3
{C}	3
{E}	3

L_2

Itemset	sup
{A, C}	2
{B, C}	2
{B, E}	3
{C, E}	2

C_2

Itemset	sup
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2

2nd scan

C_2

Itemset	sup
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2

C_3

Itemset	sup
{B, C, E}	2

3rd scan

L_3

Itemset	sup
{B, C, E}	2

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)

Improvements to Apriori

- Partition DB, find local frequent patterns, consolidate to global patterns
 - Savasere, Omiecinski, and Navathe, VLDB, 1995
- Reduce number of candidates with DHP
 - Park, Chen, and Yu, SIGMOD, 1995
- Sampling for frequent patterns, verify pattern in db
 - Toivonen, VLDB, 1996.
- Dynamic Itemset counting (DIC)
 - Brin, Motwani, Ullman, Tsur, SIGMOD, 1997

Apriori Example

Association Analysis

Apriori Algorithm Example

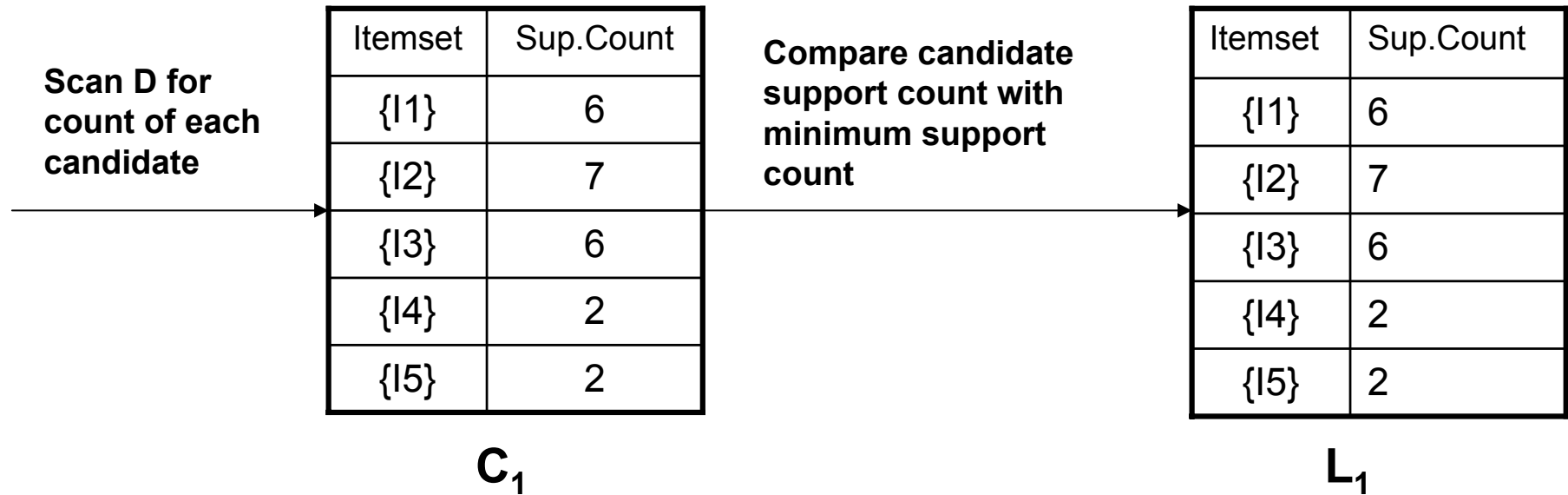
- Consider DB of 9 transactions

- $minsup = 2$

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

Step 1: Generate 1-itemset patterns

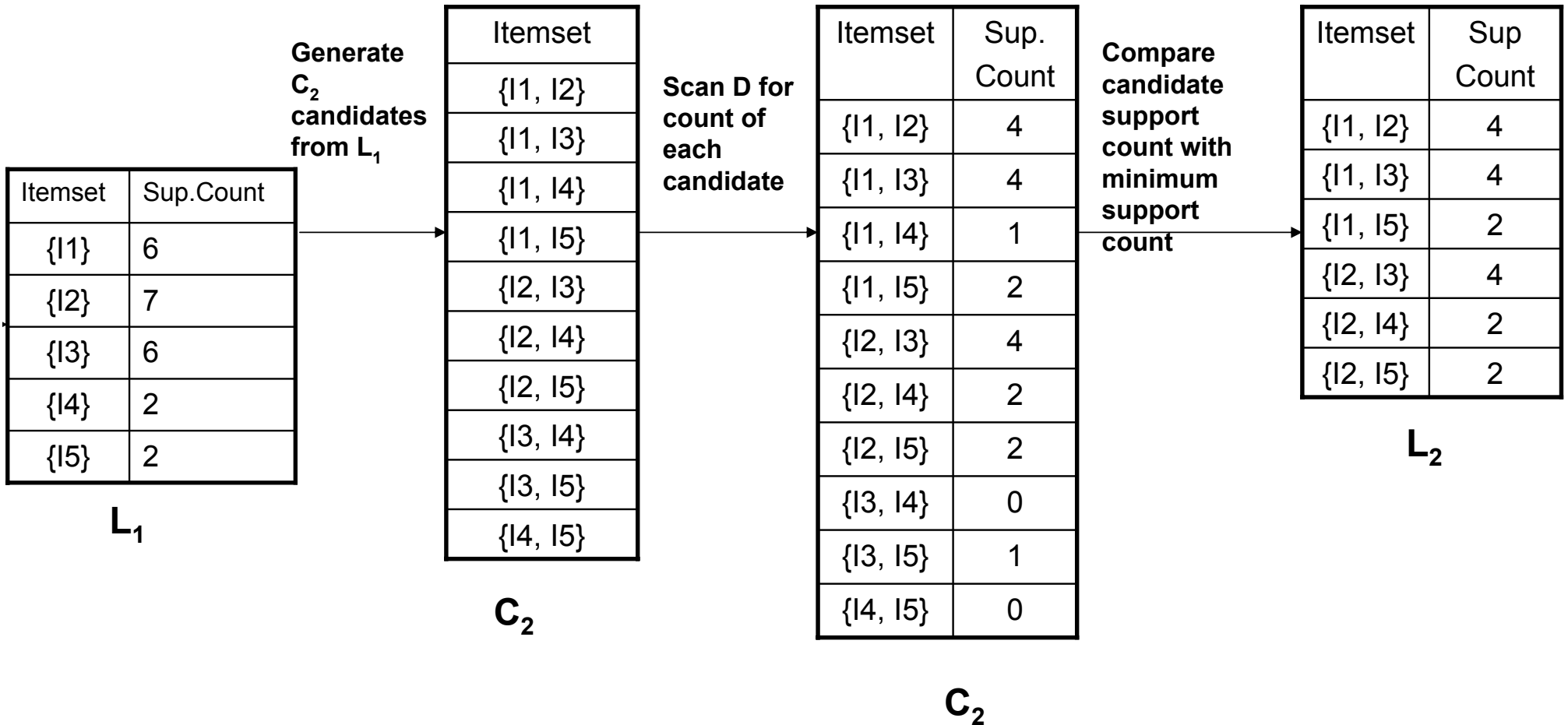
minsup = 2



- In first iteration, each item is a member of the set of candidates
- Set of frequent 1-itemset, L₁, consists of candidate 1-itemsets satisfying minimum support

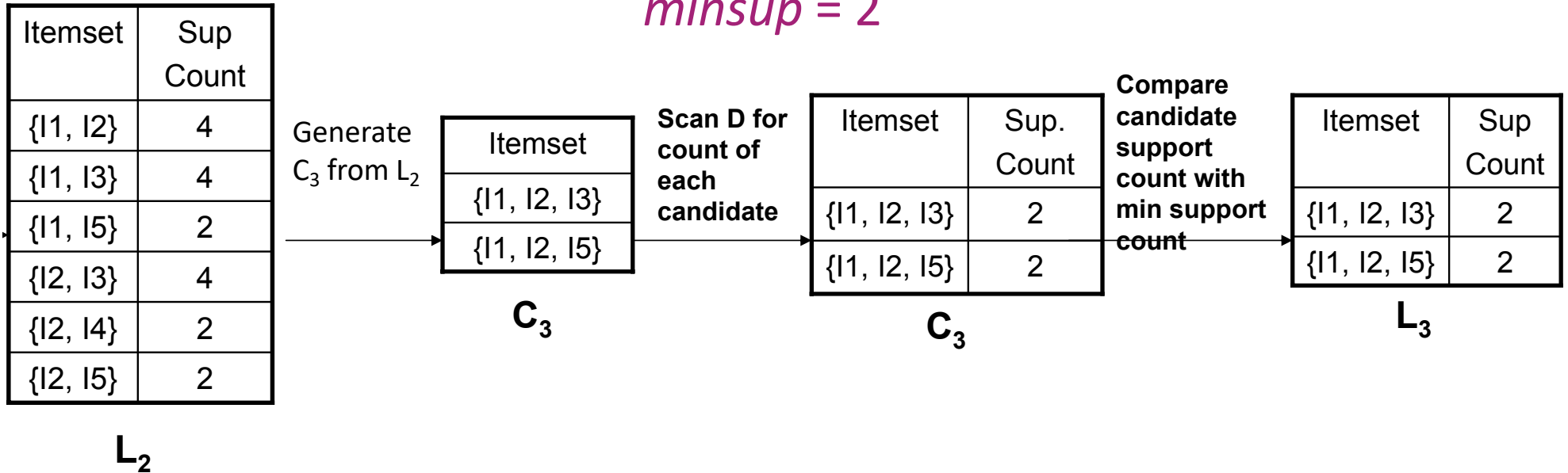
Step 2: Generate 2-itemset patterns

minsup = 2



Step 3: Generate 3-itemset patterns

minsup = 2



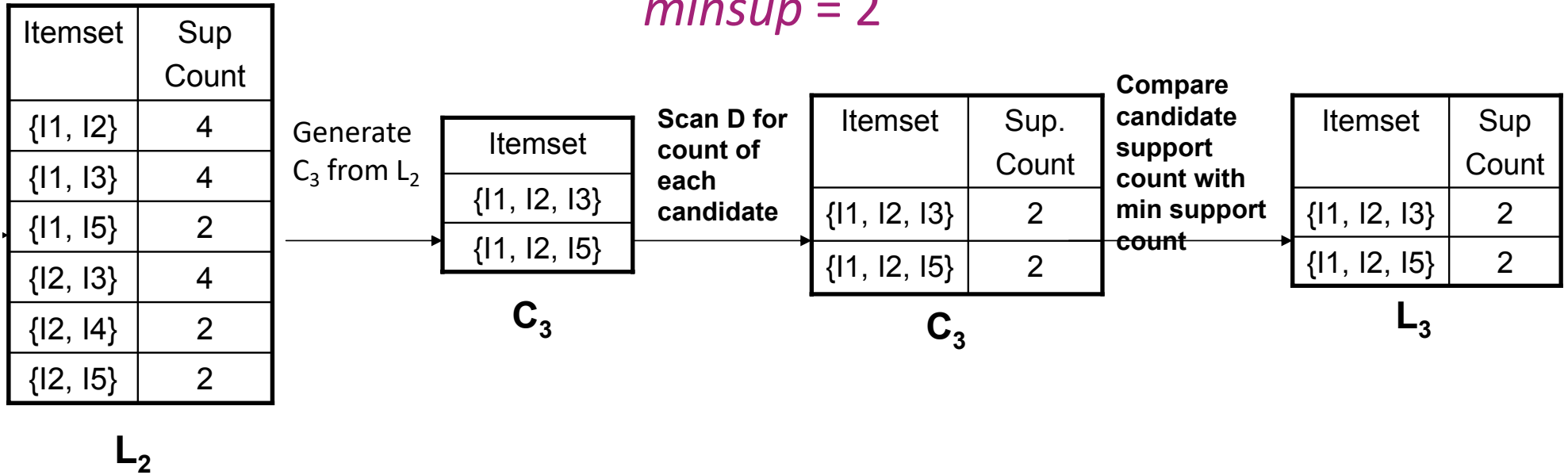
- Generation of set of candidate 3-itemsets, C₃, involves use of Apriori property
- To find C₃, compute L₂ *Join* L₂
- C₃ = { {I1, I2, I3}, {I1 I2, I5}, {I1, I3, I5}, {I2, I3, I4}, {I2, I3, I5}, {I2, I4, I5} }
- Join step complete, prune step used to reduce size of C₃

Step 3: Generate 3-itemset patterns

- Use Apriori property: all subsets of frequent itemsets must also be frequent
- Ex. $\{I_1, I_2, I_3\}$ has 2-itemsets subsets of:
 - $\{I_1, I_2\}, \{I_1, I_3\}, \{I_2, I_3\}$ are all in L_2 ,
 - keep $\{I_1, I_2, I_3\}$ in C_3
- Ex. $\{I_2, I_3, I_5\}$ has 2-itemsets subsets of:
 - $\{I_2, I_3\}, \{I_2, I_5\}, \{I_3, I_5\}$
 - $\{I_3, I_5\}$ not a member of L_2 , thus $\{I_2, I_3, I_5\}$ not in C_3
- Therefore, $C_3 = \{\{I_1, I_2, I_3\}, \{I_1, I_2, I_5\}\}$

Step 3: Generate 3-itemset patterns

minsup = 2



- Generation of set of candidate 3-itemsets, C₃, involves use of Apriori property
- Prune operation used to reduce size of C₃

Step 4: Generate 4-itemset patterns

- Algorithm uses L_3 *Join* L_3 to generate 4-itemsets, C_4 .
 - Join results in $\{ \{I1, I2, I3, I5\} \}$
 - itemset is pruned since $\{ \{I2, I3, I5\} \}$ is not frequent
- Algorithm terminates, having found all frequent items

$$\mathcal{F} = \{L_1, L_2, L_3\} = \\ \{1, 2, 3, 4, 5, 12, 13, 15, 23, 24, 25, 123, 125\}$$

- Still left to do
 - generate association rules from itemsets
 - improve efficiency

Frequent Itemset Mining with Vertical Data

Association Analysis: ECLAT

Frequent Mining with Vertical Data

- Vertical format
 - for each item store a list of transaction IDs (tids)
- tid-list: list of tids for itemsets
 - $t(AB) = \{T_{11}, T_{25}, \dots\}$
- Derive frequent patterns based on vertical intersections
 - $t(X) = t(Y)$: X and Y always happen together
 - $t(X) \subset t(Y)$: transaction having X always has Y

ECLAT – Equivalence Class Transformation

- For each item, store a list of transaction ids (tids)

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				

↓
TID-list

ECLAT

- Determine support of any k-itemset by intersection

A		B	=	AB
1		1		1
4		2		5
5	\cap	5		7
6		7		8
7		8		
8		10		
9				

- Use **diffset** to accelerate mining
 - only keep track of difference of tids
 - $\text{Diffset}(AB, A) = \{ 4, 6, 9 \}$, $\text{Diffset}(AB, B) = \{ 2, 10 \}$

ECLAT

- 3 traversal approaches for itemsets
 - top-down, bottom-up, and hybrid
- Advantages: very fast support counting
- Disadvantages: intermediate tid-lists may become too large for memory
- References:
 - ECLAT – Zaki et al., KDD 1997
 - Mining closed patterns with vertical format: CHARM – Zaki & Hsiao, SDM 2002

Frequent Pattern Tree Methods: FPGrowth

Association Analysis: FPGrowth

Algorithms for Mining Frequent Patterns

- Bottlenecks of Apriori
 - breadth-first (i.e., level-wise) search
 - candidate generation and test
 - may generate huge number of candidates
- FPGrowth Approach (Han, Pei, Yin SIGMOD, 2000)
 - depth-first search
 - avoid explicit candidate generation
- Main Idea – grow long patterns from short ones using local frequent items only

Construct FP-tree

- Compress a large database into a compact, **Frequent-Pattern tree (FP-tree)** structure
 - highly condensed, but complete for frequent pattern mining
 - helps avoid costly database scans
- Develop an efficient, FP-tree based frequent pattern mining method
 - divide and conquer methodology: decompose mining tasks into smaller ones
 - avoid candidate generation: sub-database test only

Construct FP-tree: Overview

<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_support = 3

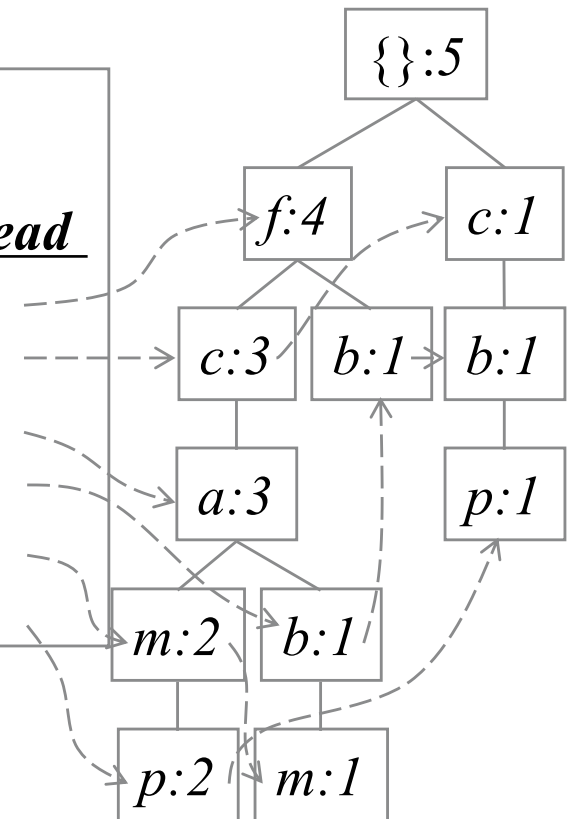
1. Scan DB once, find frequent 1-itemset (single item pattern)
2. Sort frequent items in frequency descending order, f-list
3. Scan DB again, construct FP-tree

Header Table

Item frequency head

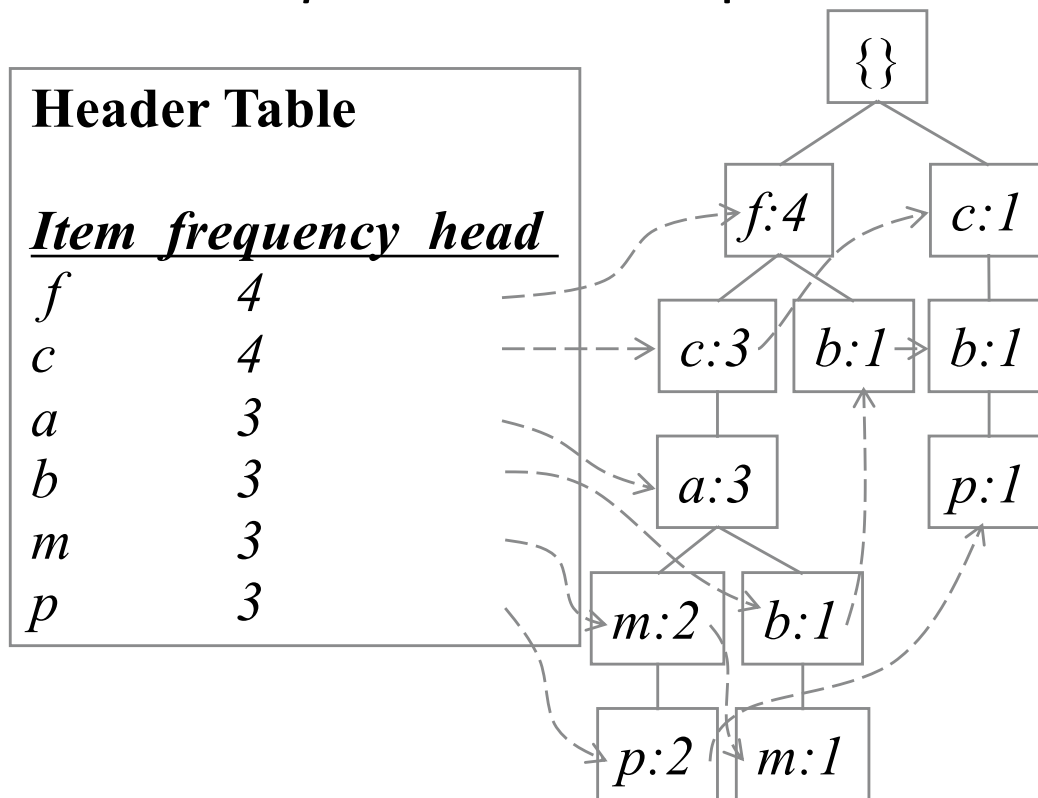
<i>f</i>	4
<i>c</i>	4
<i>a</i>	3
<i>b</i>	3
<i>m</i>	3
<i>p</i>	3

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



Find Patterns using FP-Tree

- Starting at the frequent item header table in the FP-tree
- Traverse the FP-tree by following the link of each frequent item p
- Accumulate all of the transformed prefix paths of item p to form p 's conditional pattern base



Conditional pattern bases

item *cond. pattern base*

c $f:3$

a $fc:3$

b $fca:1, f:1, c:1$

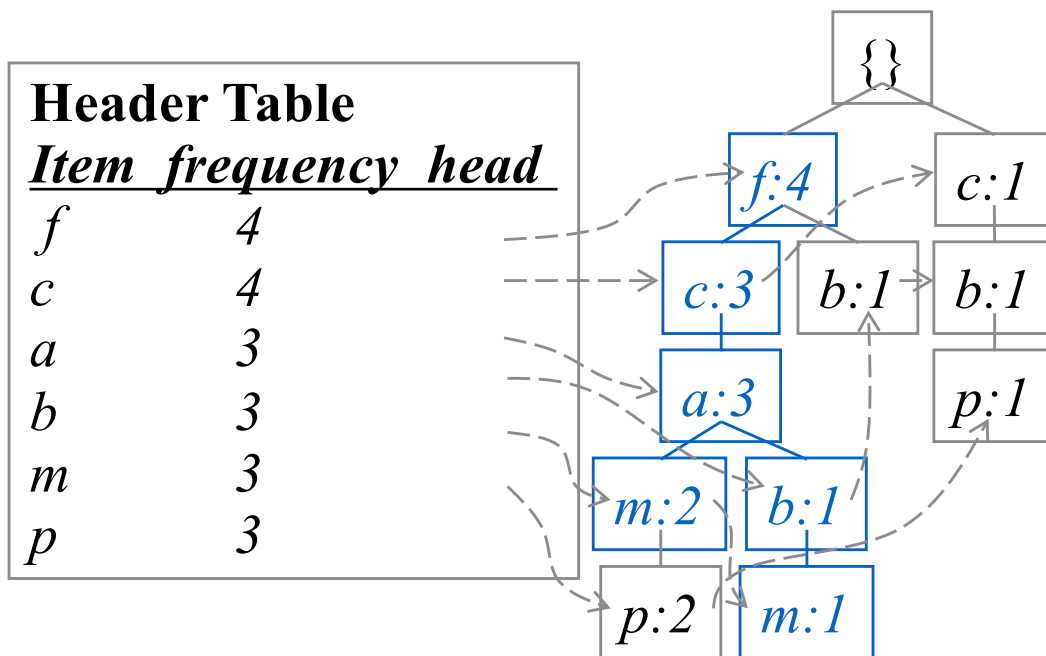
m $fca:2, fcab:1$

p $fcam:2, cb:1$

From Conditional Pattern-base to Conditional FP-Tree

For each pattern-base

- Accumulate the count for each item in the base
- Construct the FP-tree for the frequent items of the pattern base



m-conditional pattern base:
fca:2, fcab:1



{}

|

f:3

|

c:3

|

a:3



All frequent
patterns relate to *m*
m,
fm, cm, am,
fcm, fam, cam,
fcam

m-conditional FP-tree

FPGrowth: Example

Part 1: Construct the FP-tree

FPGrowth Example

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

- Consider database D
- Let $minsup = 2$
- First scan is same as Apriori to derive 1-itemsets and their support counts
- Set of frequent items is sorted in order of descending support count
- $L = \{I2:7, I1:6, I3:6, I4:2, I5:2\}$

Construct FP-tree

- Create root of tree, labeled “null”
- Scan D a 2nd time (first scan was to create 1-itemsets and L)
- For each transaction, items are processed in L order (sorted order)
- Branch created for each transaction with items having their support count separated by colon
- Whenever same node is encountered in another transaction just increment support count of common node or prefix
- To facilitate tree traversal, an item header table is built so that each item points to its occurrences in the tree via a chain of node-links
- The problem of mining frequent patterns in D is transformed to mining the FP-tree

Construct FP-Tree: Start

Start, root = null

{ }:



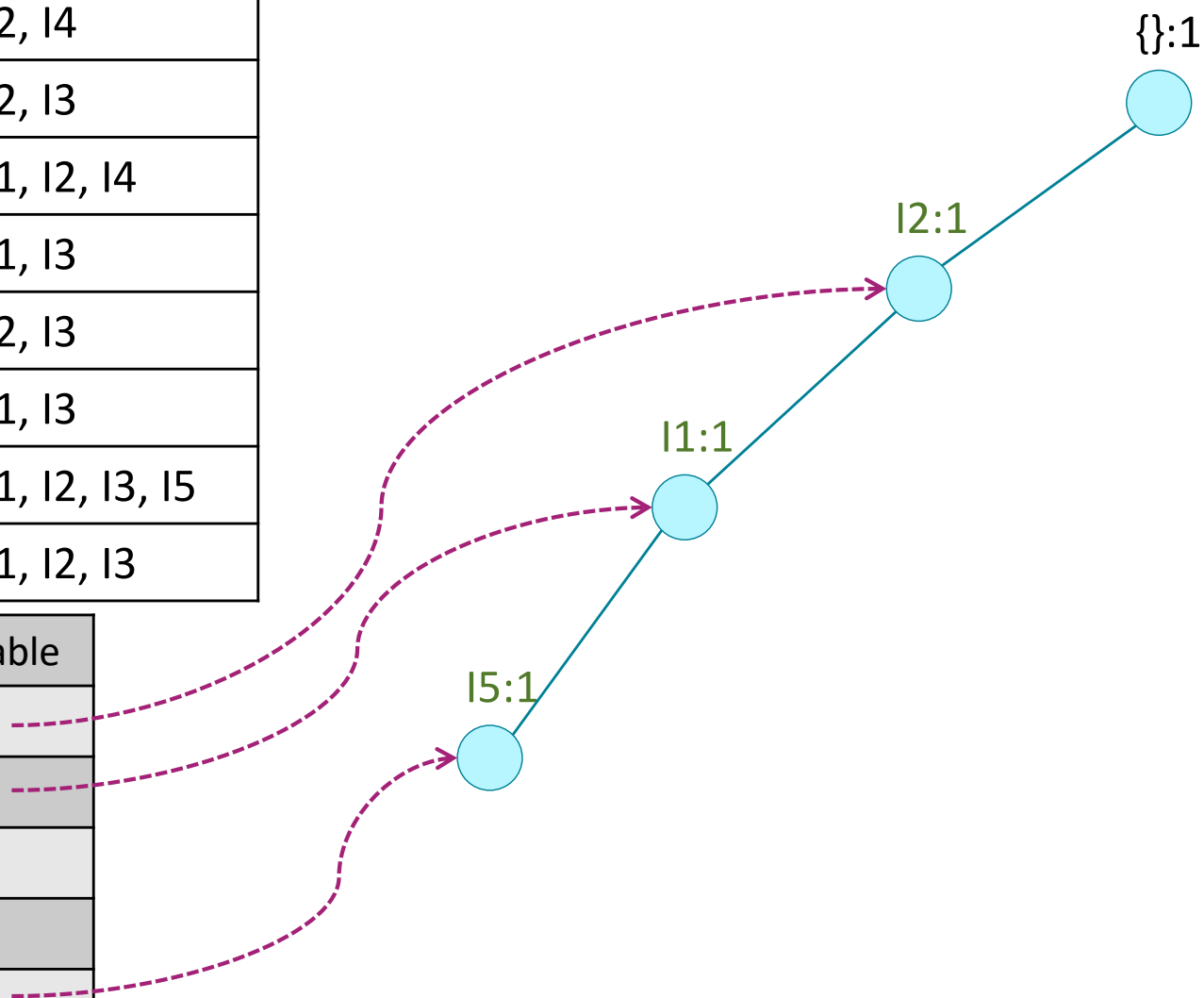
TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

Header Table	
I2	
I1	
I3	
I4	
I5	

Construct FP-Tree: Trans. #1

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

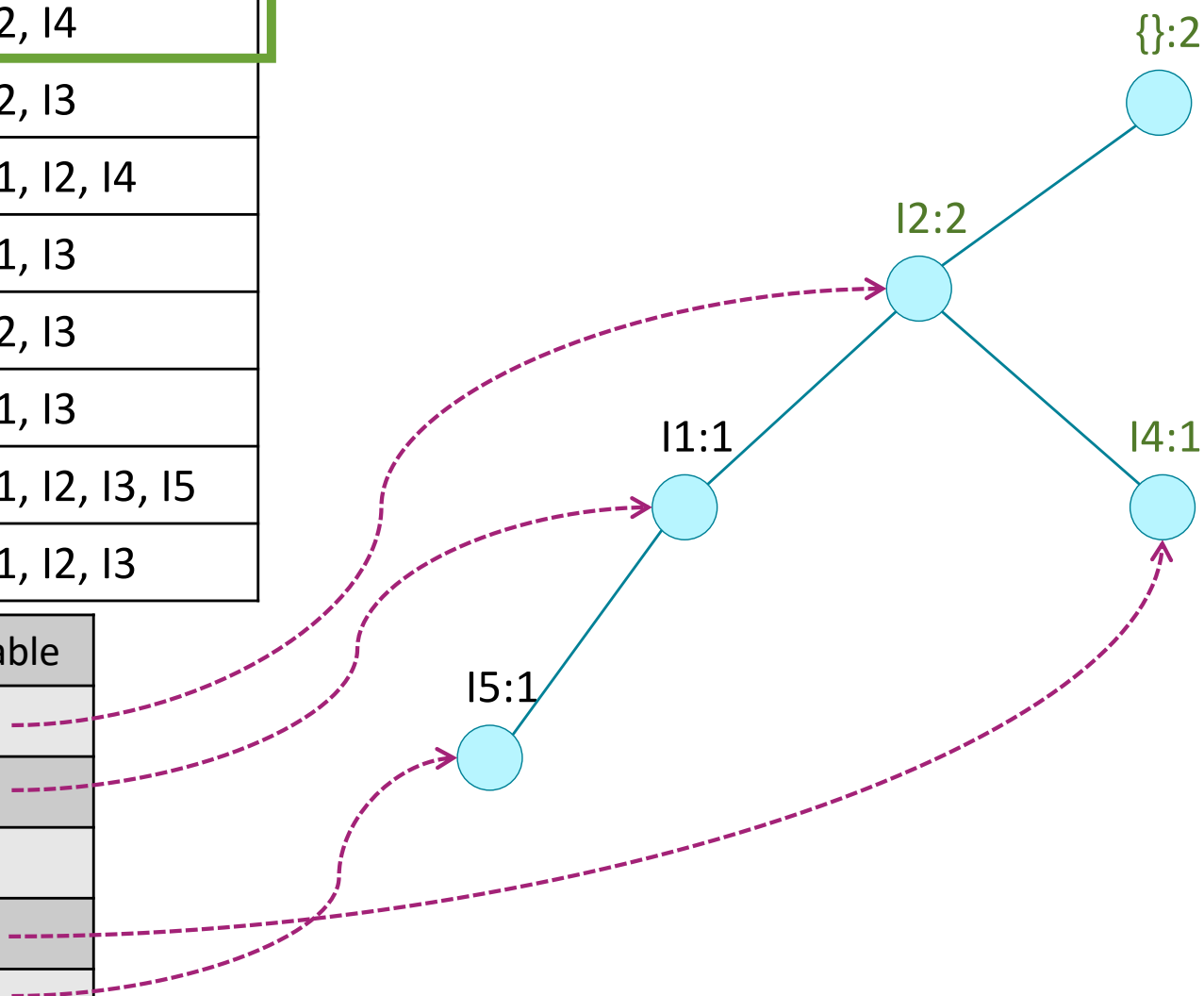
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #2

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

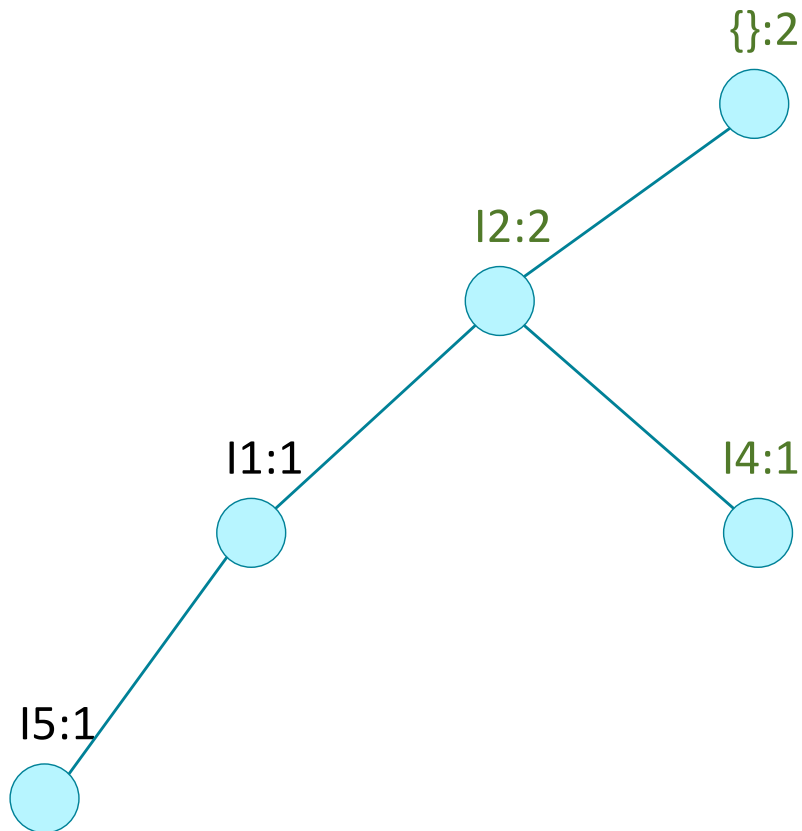
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #2

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

Remove header links to better see building of tree

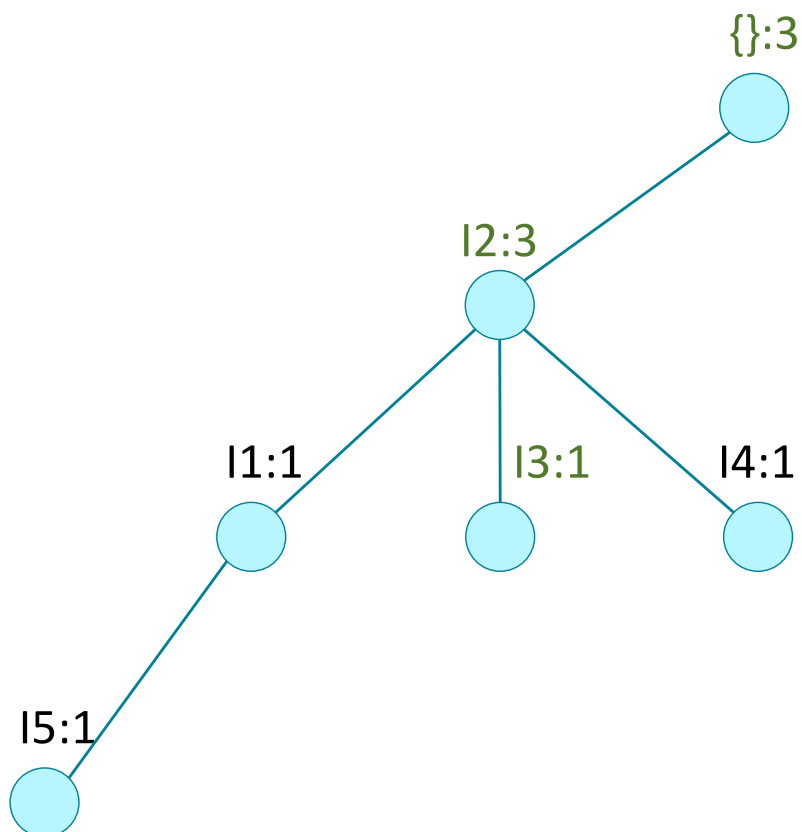


Header Table	
I2	
I1	
I3	
I4	
I5	

Construct FP-Tree: Trans. #3

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

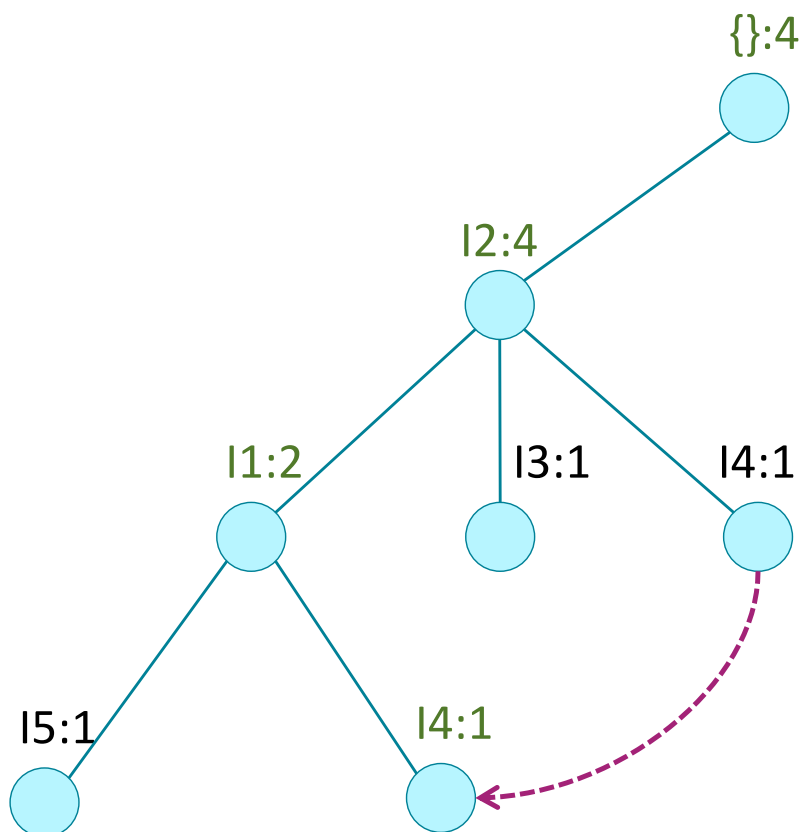
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #4

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

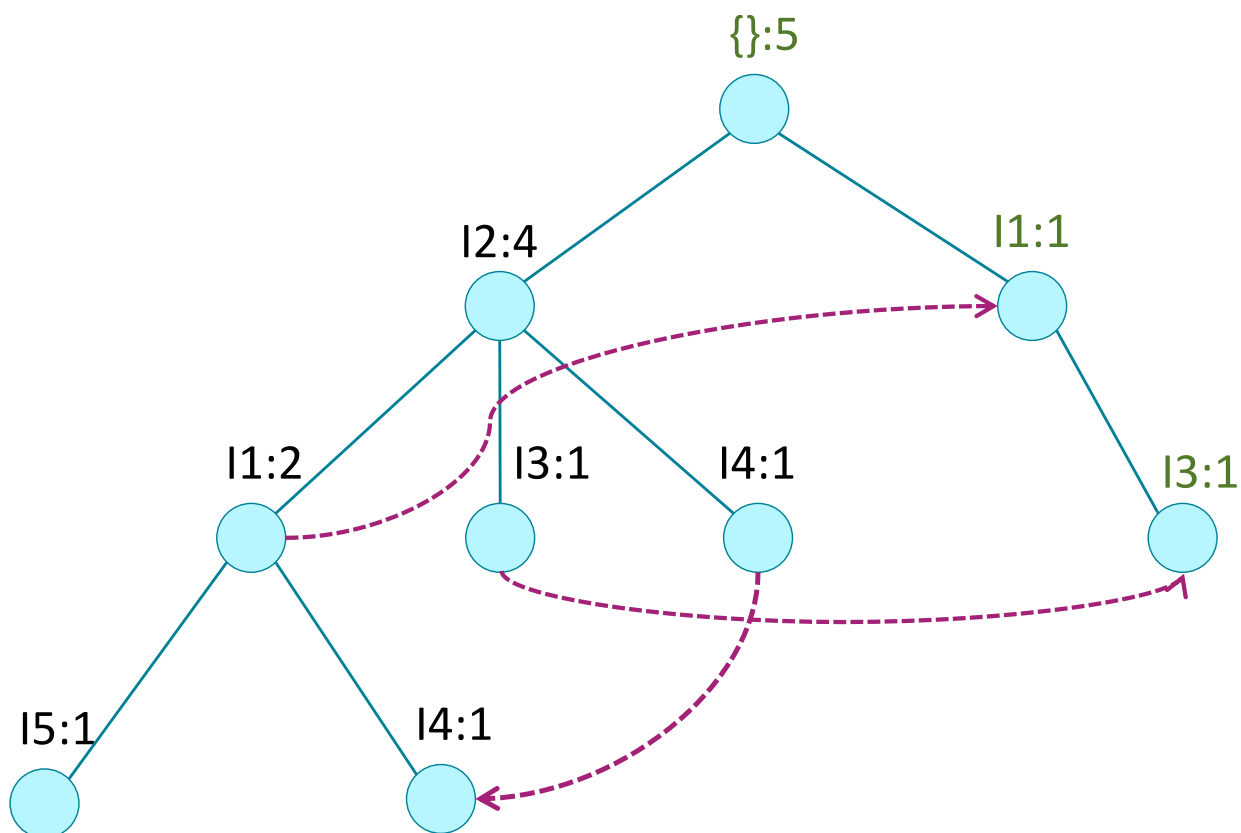
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #5

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

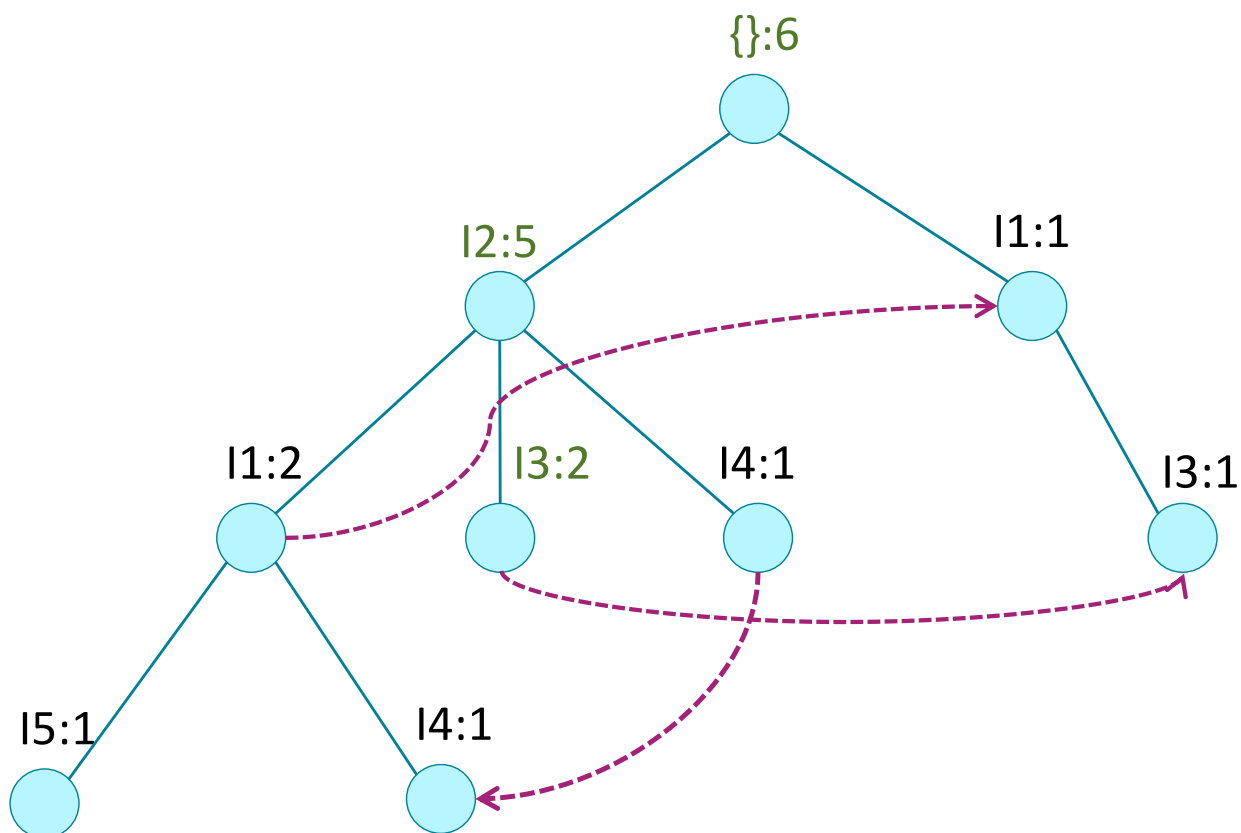
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #6

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

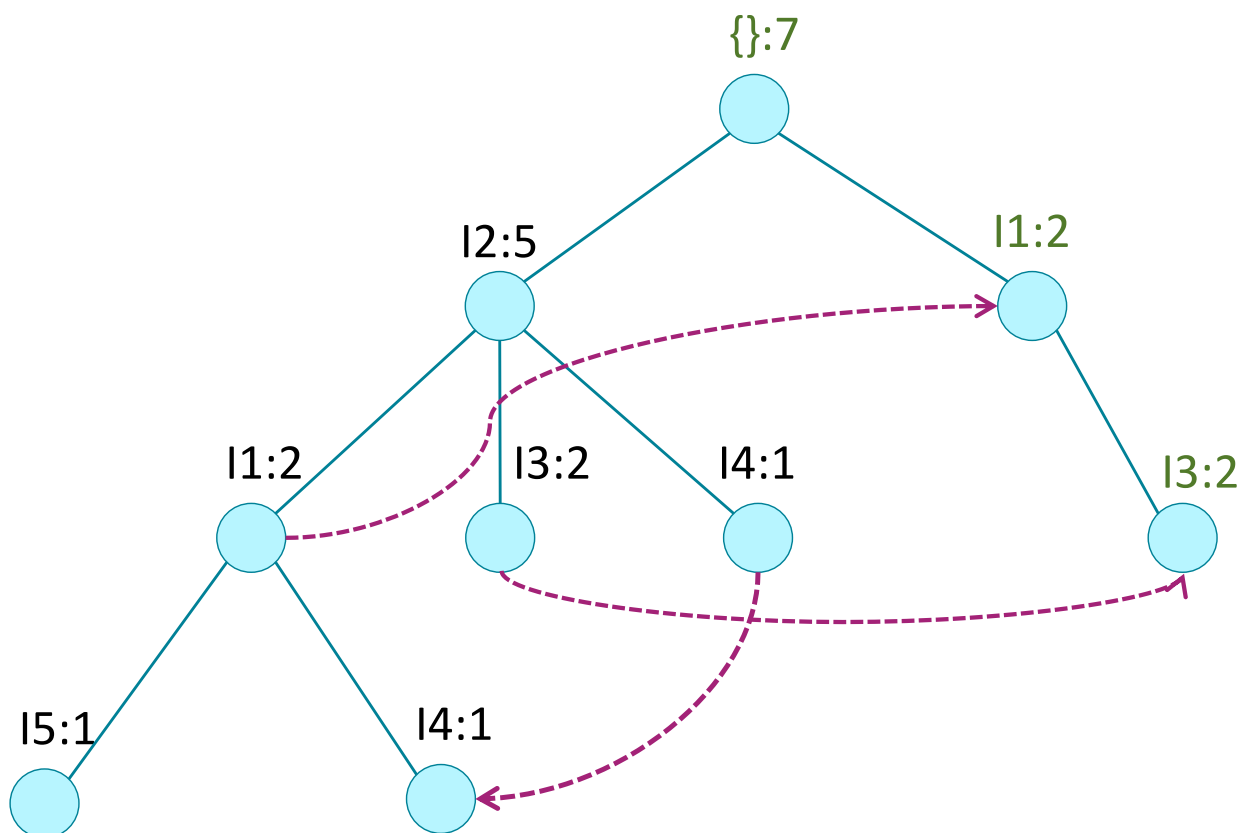
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #7

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

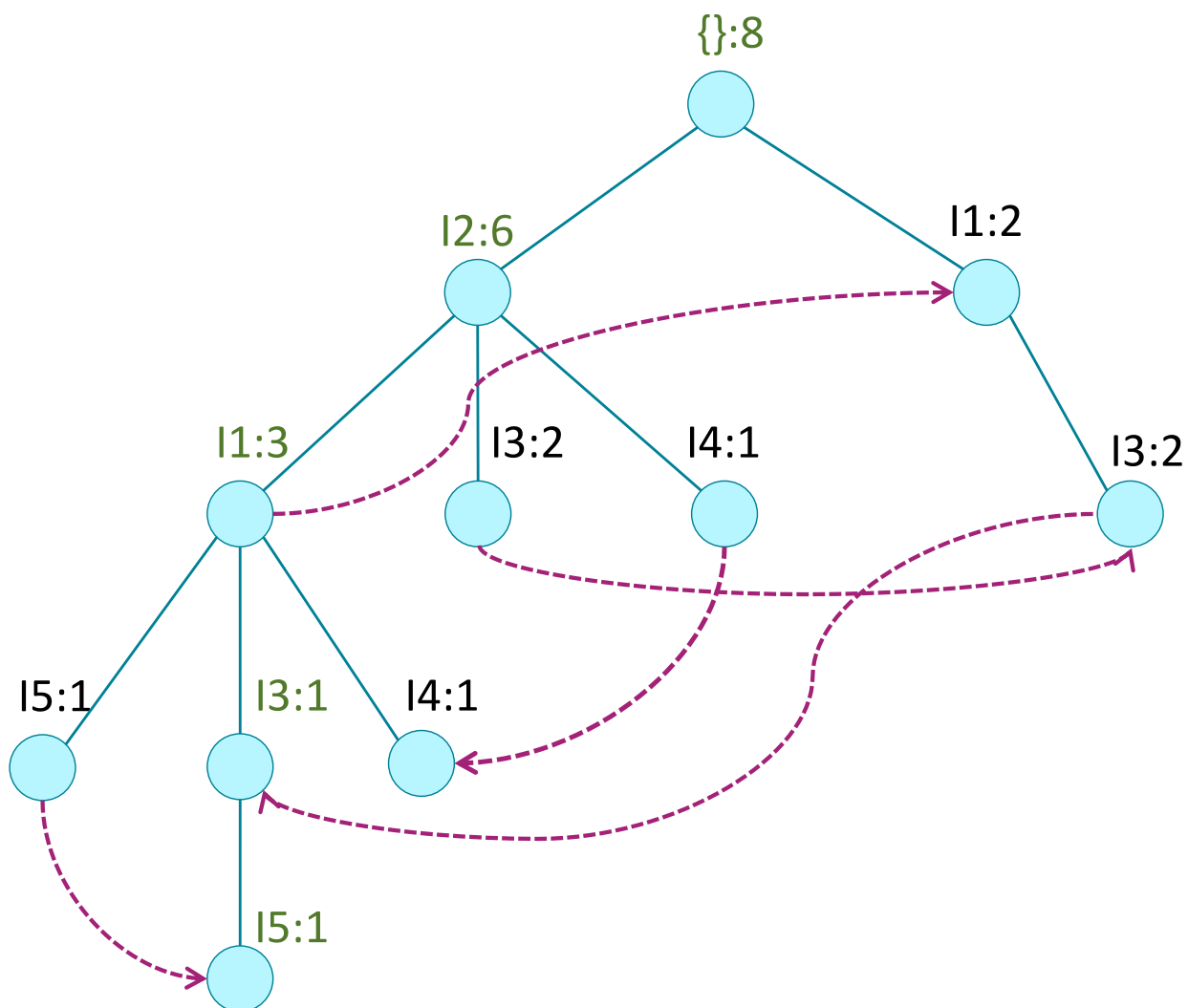
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #8

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

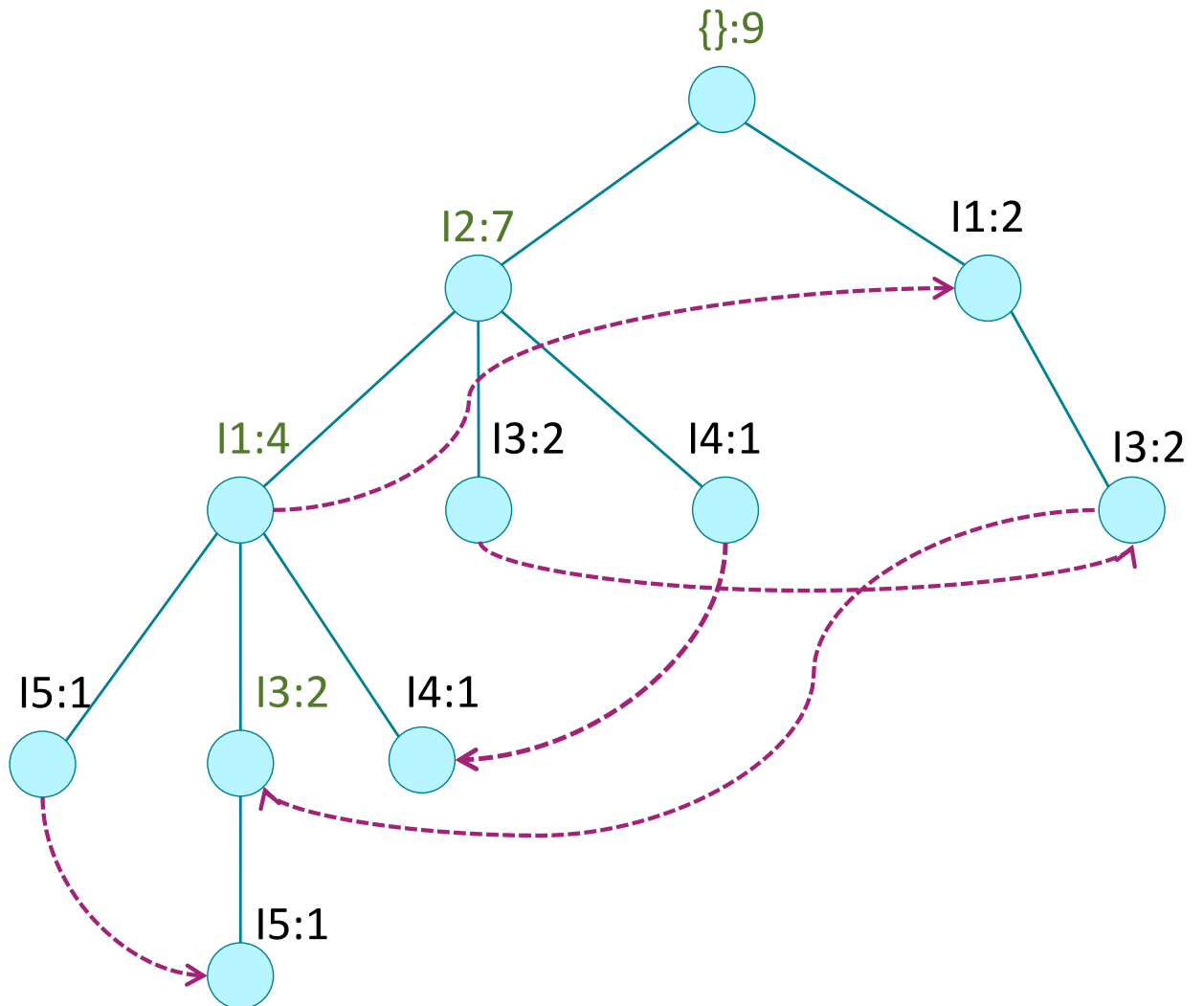
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Trans. #9

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

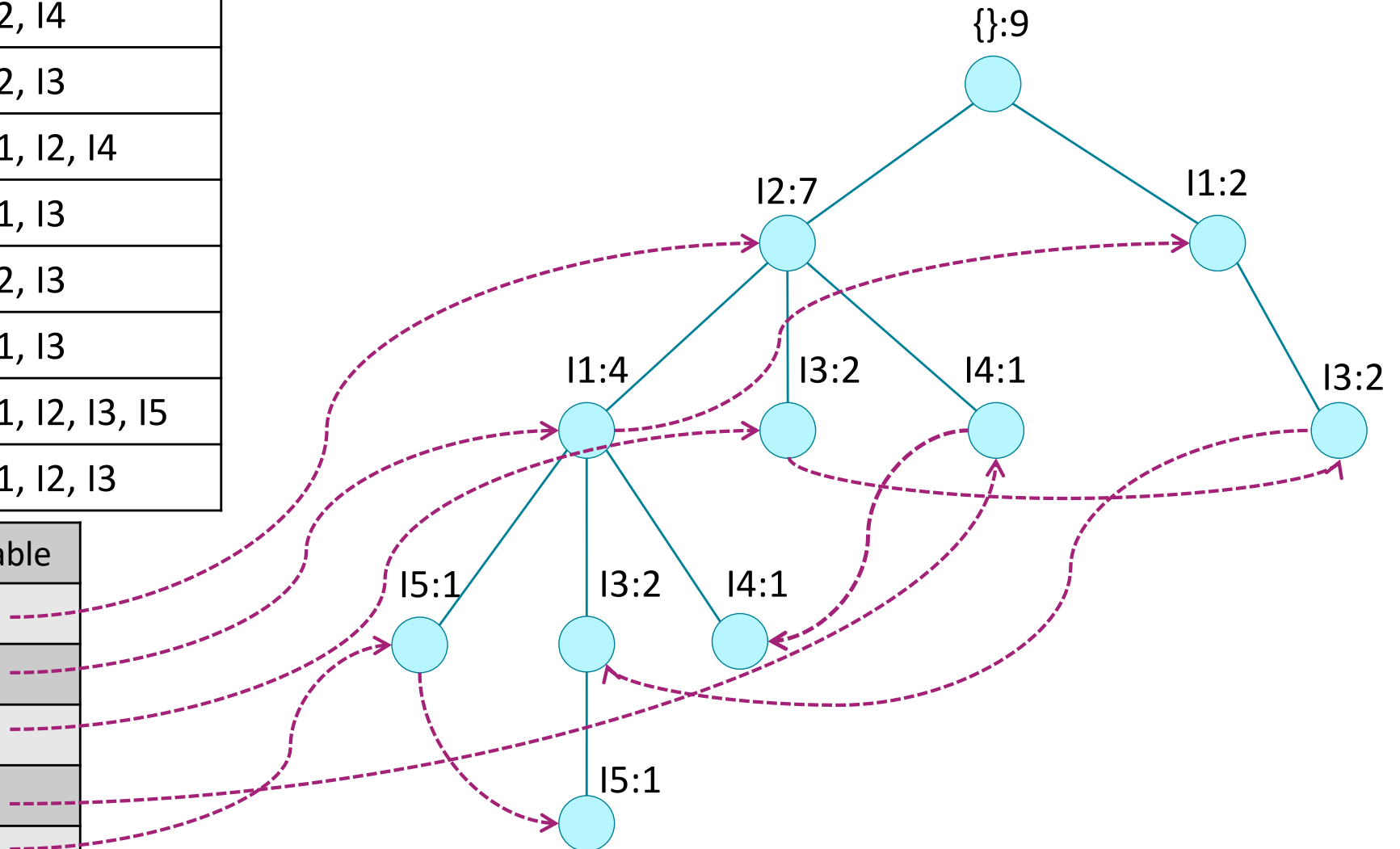
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct FP-Tree: Complete

TID	Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

Header Table	
I2	
I1	
I3	
I4	
I5	



FPGrowth: Example

Part 2: Mine the FP-tree

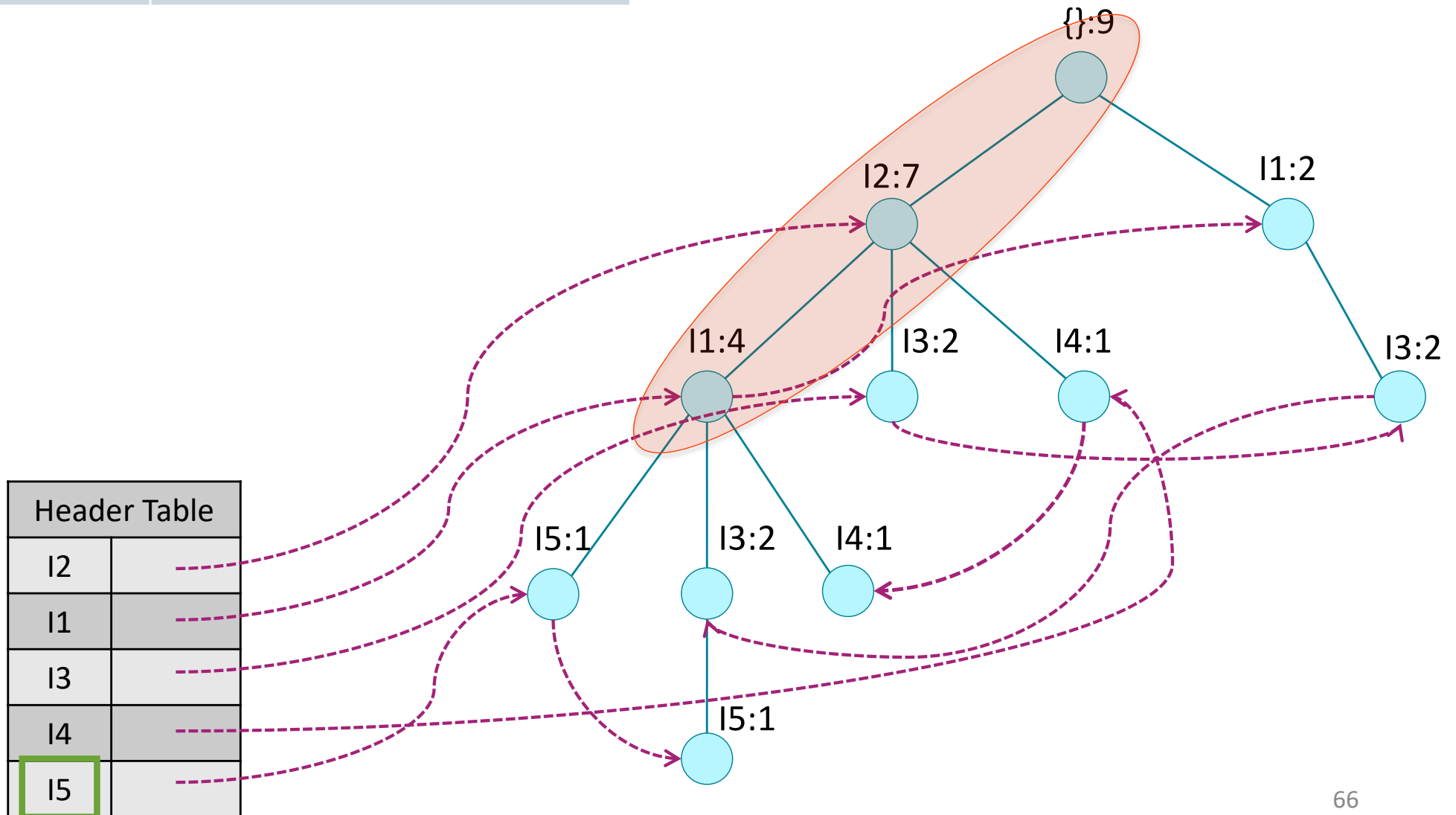
Mine FP-tree: Create Conditional Pattern Base

Steps

1. Start from each frequent length 1-pattern i (as an initial suffix pattern) in increasing order of support
2. Construct its conditional pattern base which consists of the set of prefix paths in the FP-tree co-occurring with suffix pattern (path from i in FP-tree to root)
3. Then, construct its conditional FP-tree & perform mining on such a tree
4. The pattern growth is achieved by concatenation of the suffix pattern with the frequent patterns generated from a conditional FP-tree
5. The union of all frequent patterns (generated by step 4) gives the required frequent itemset

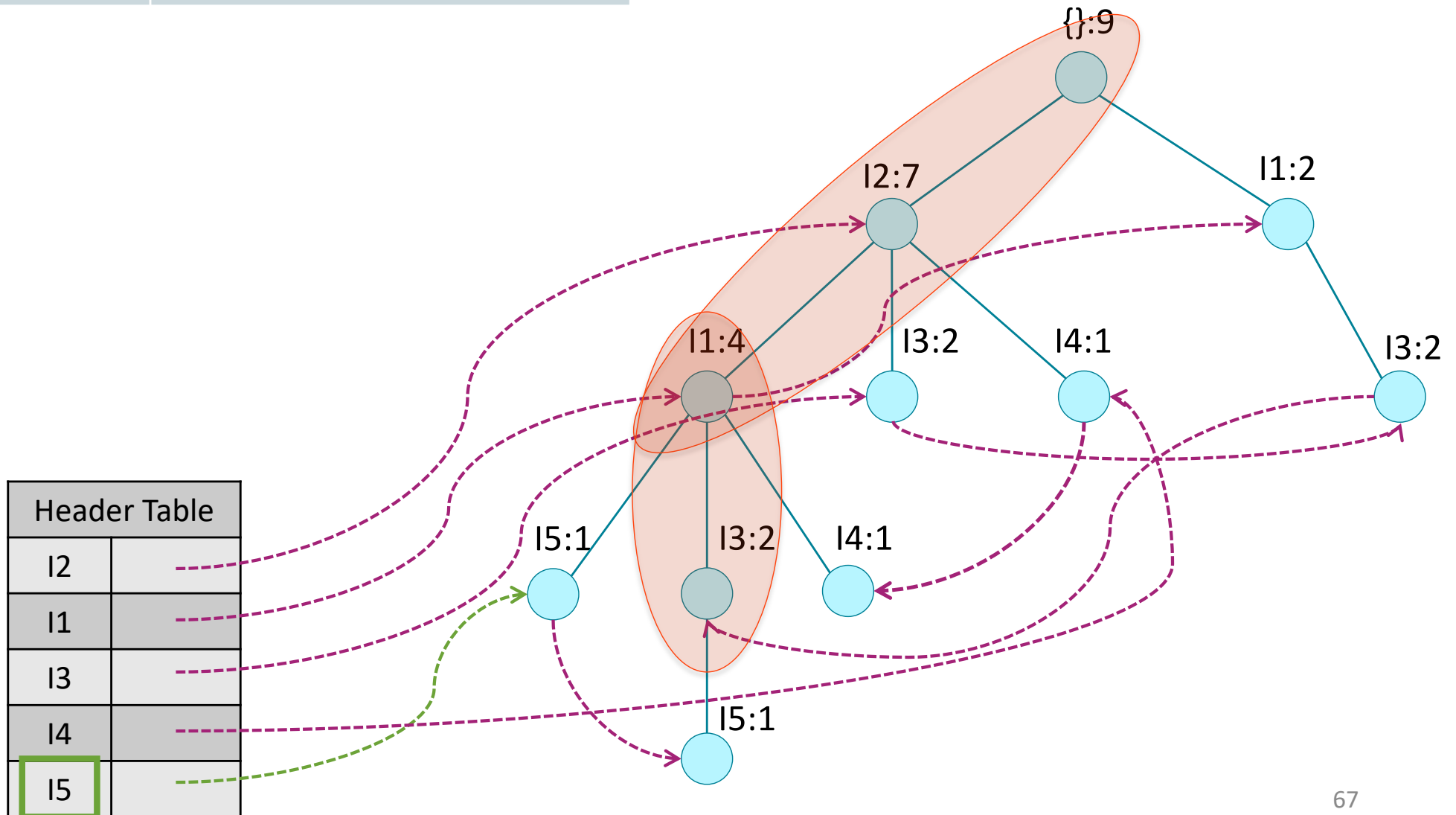
Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}}



Construct Conditional Pattern Base

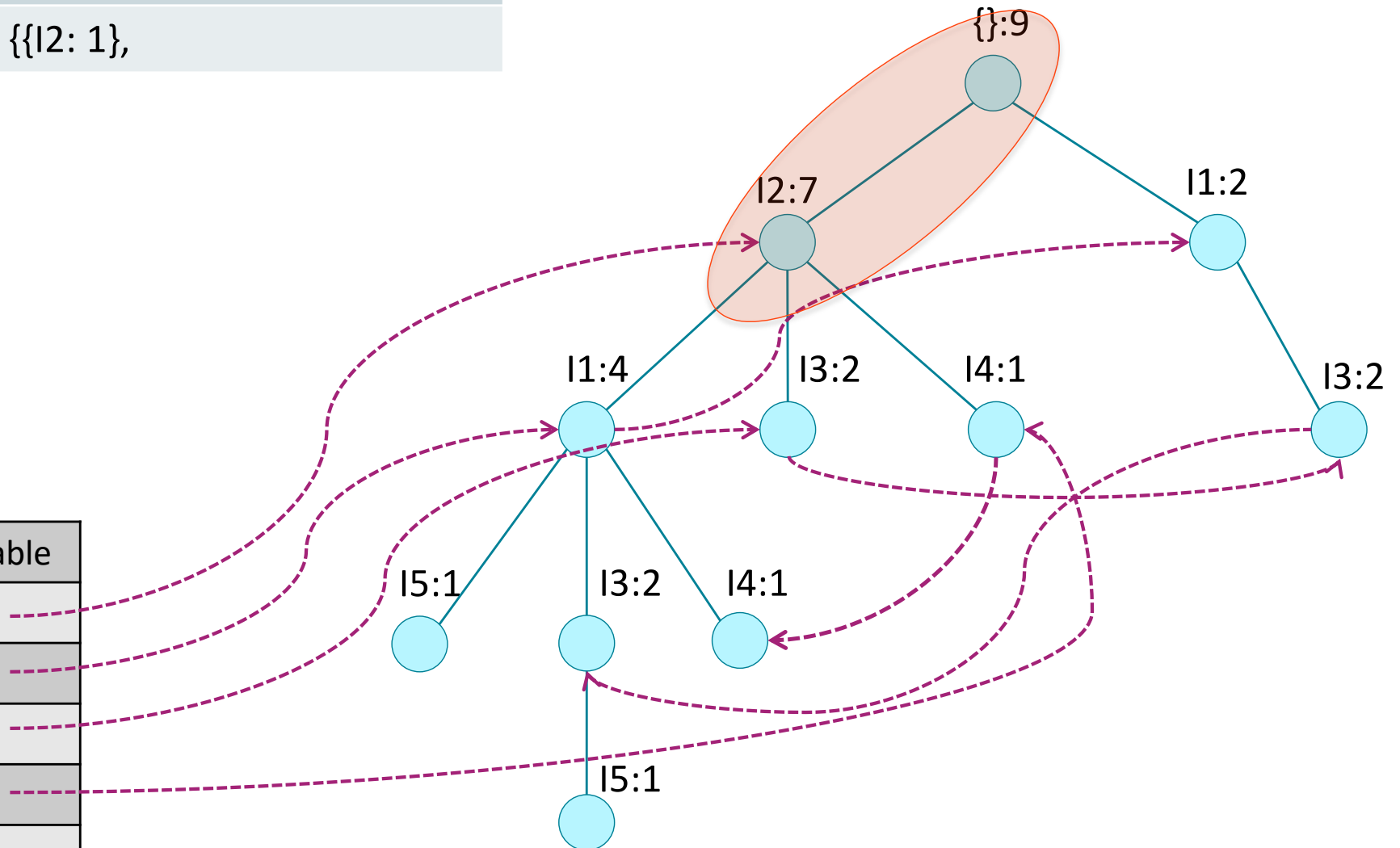
Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}



Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}
I4	{{I2: 1},

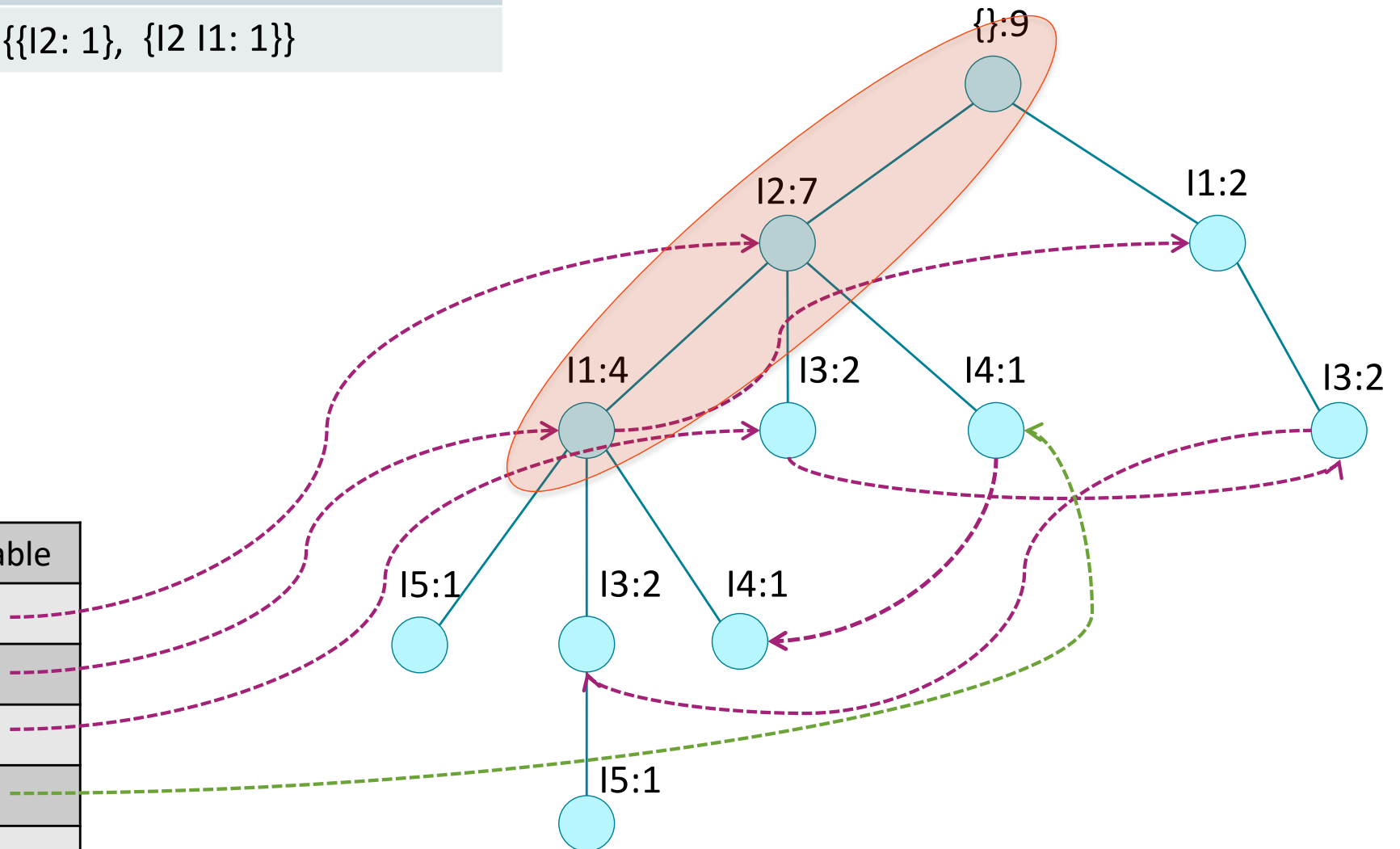
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}
I4	{{I2: 1}, {I2 I1: 1}}

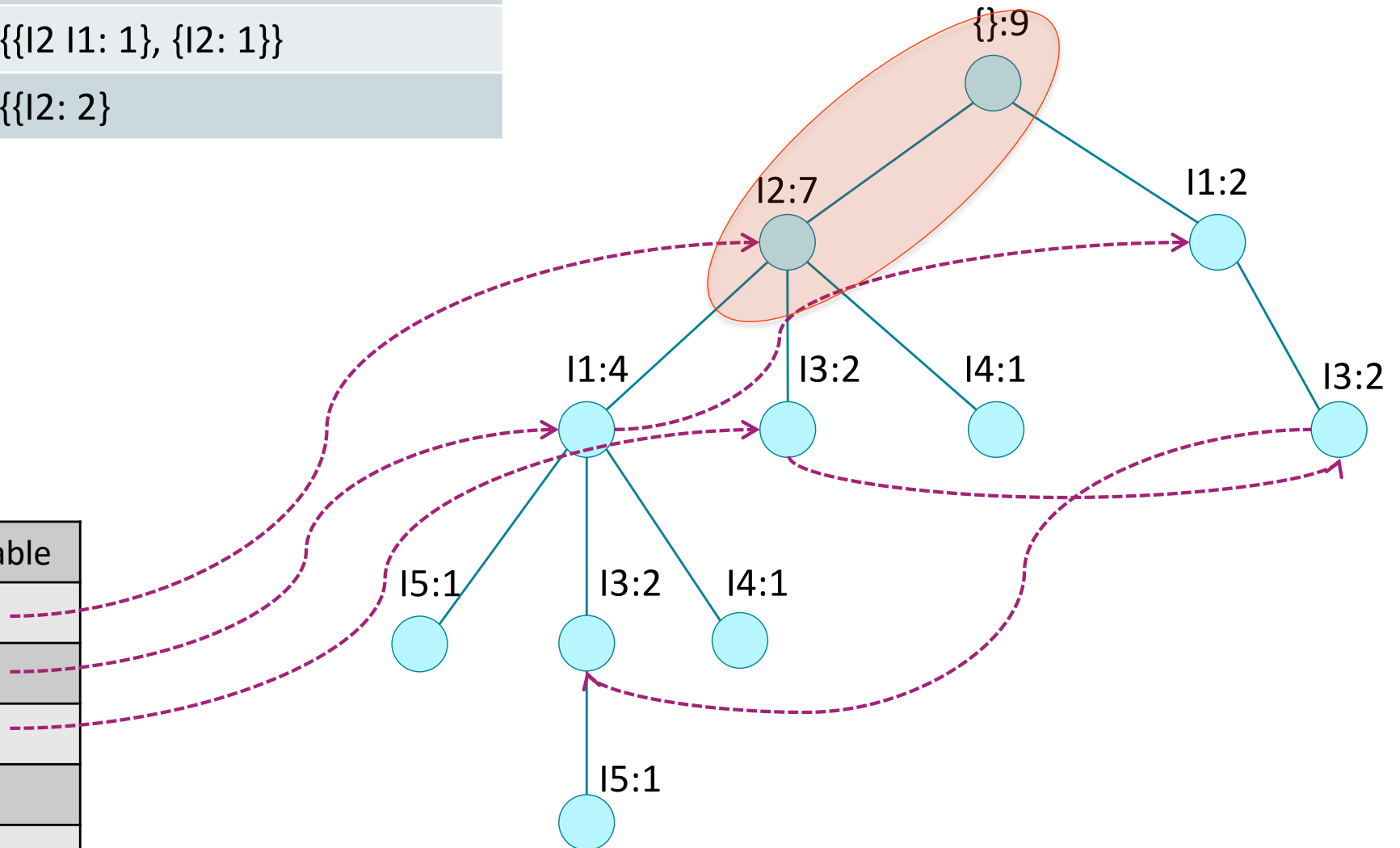
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}
I4	{{I2 I1: 1}, {I2: 1}}
I3	{{I2: 2}}

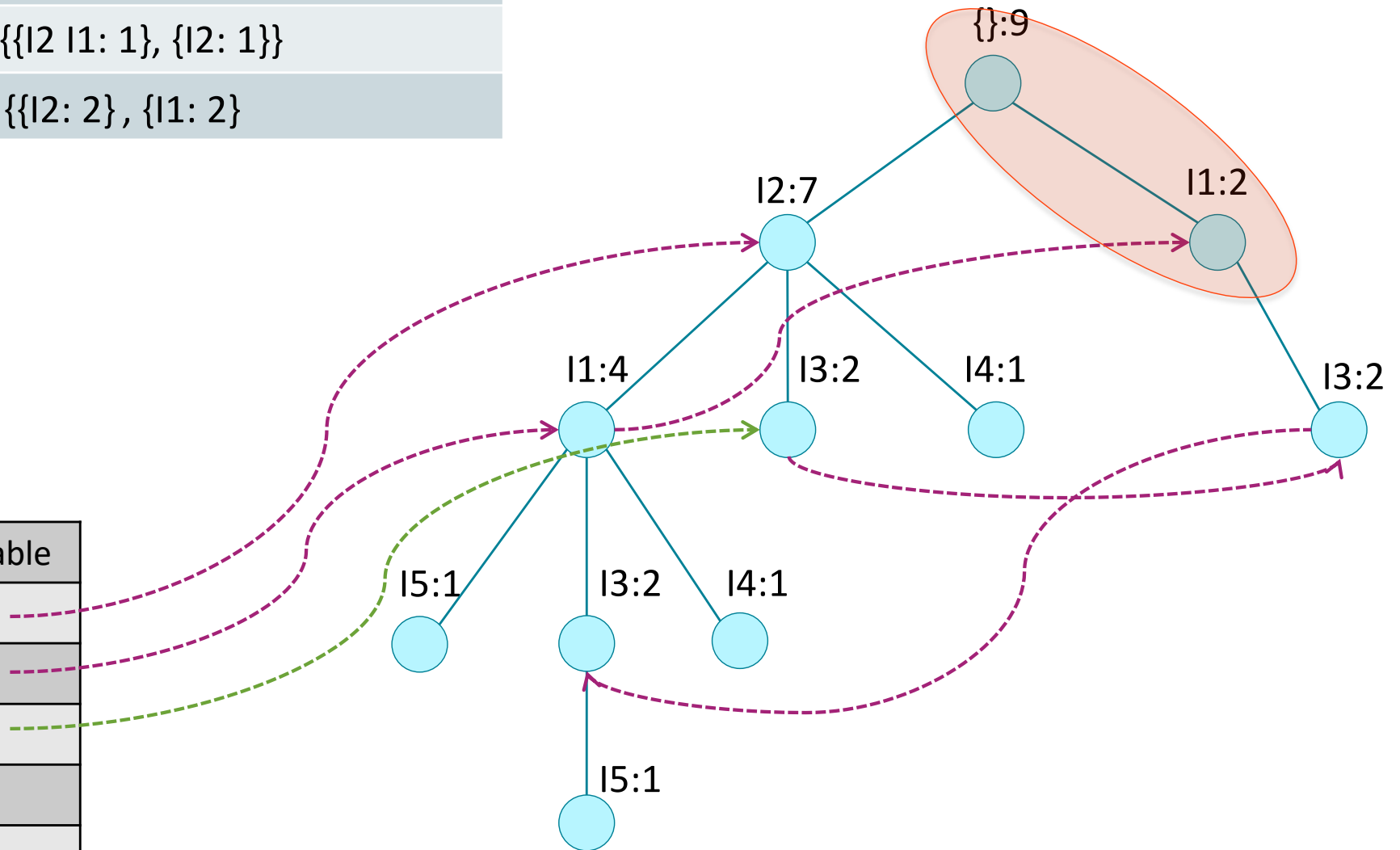
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}
I4	{{I2 I1: 1}, {I2: 1}}
I3	{{I2: 2}, {I1: 2}}

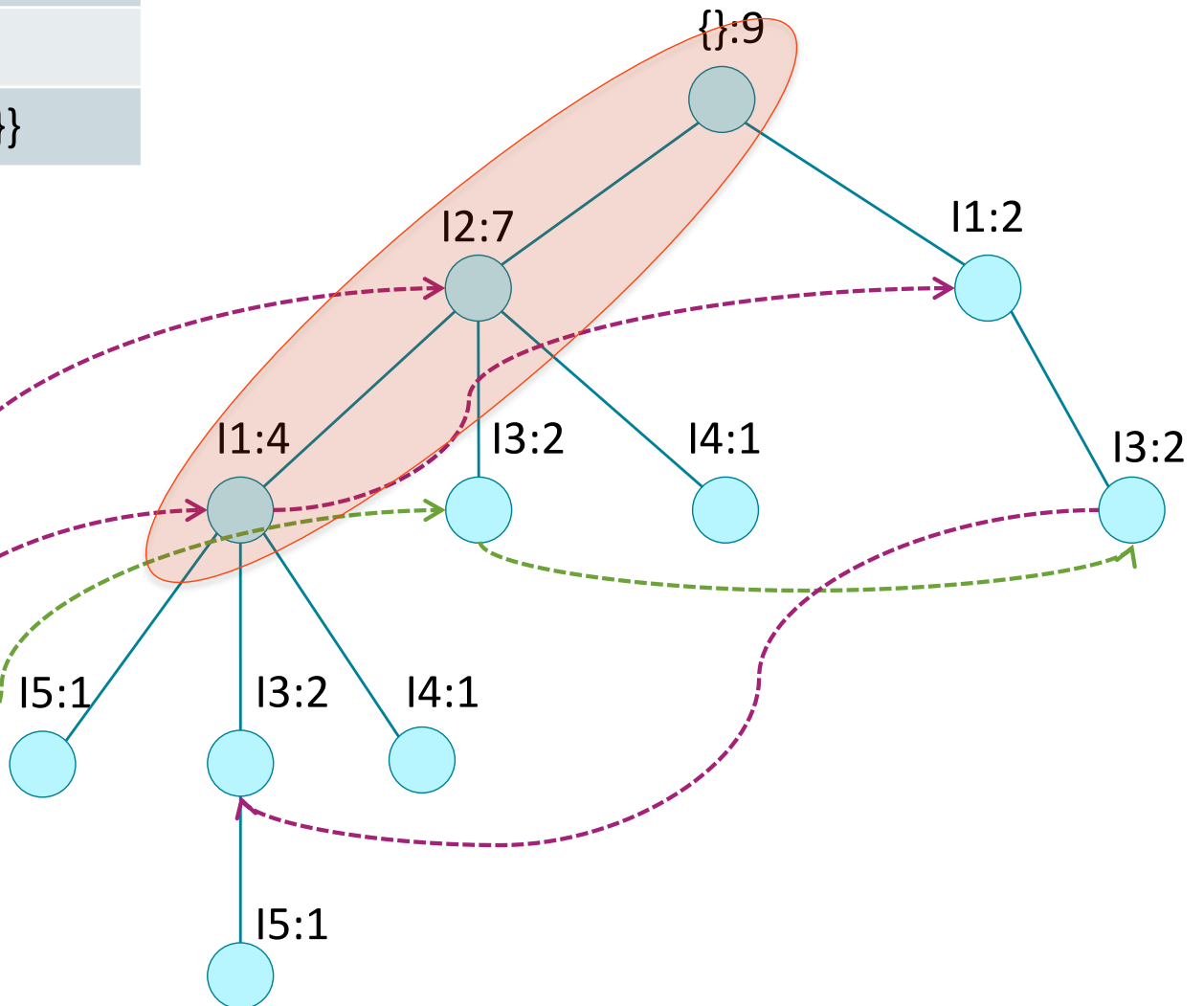
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}
I4	{{I2 I1: 1}, {I2: 1}}
I3	{{I2: 2}, {I1: 2}, {I1 I2:2}}

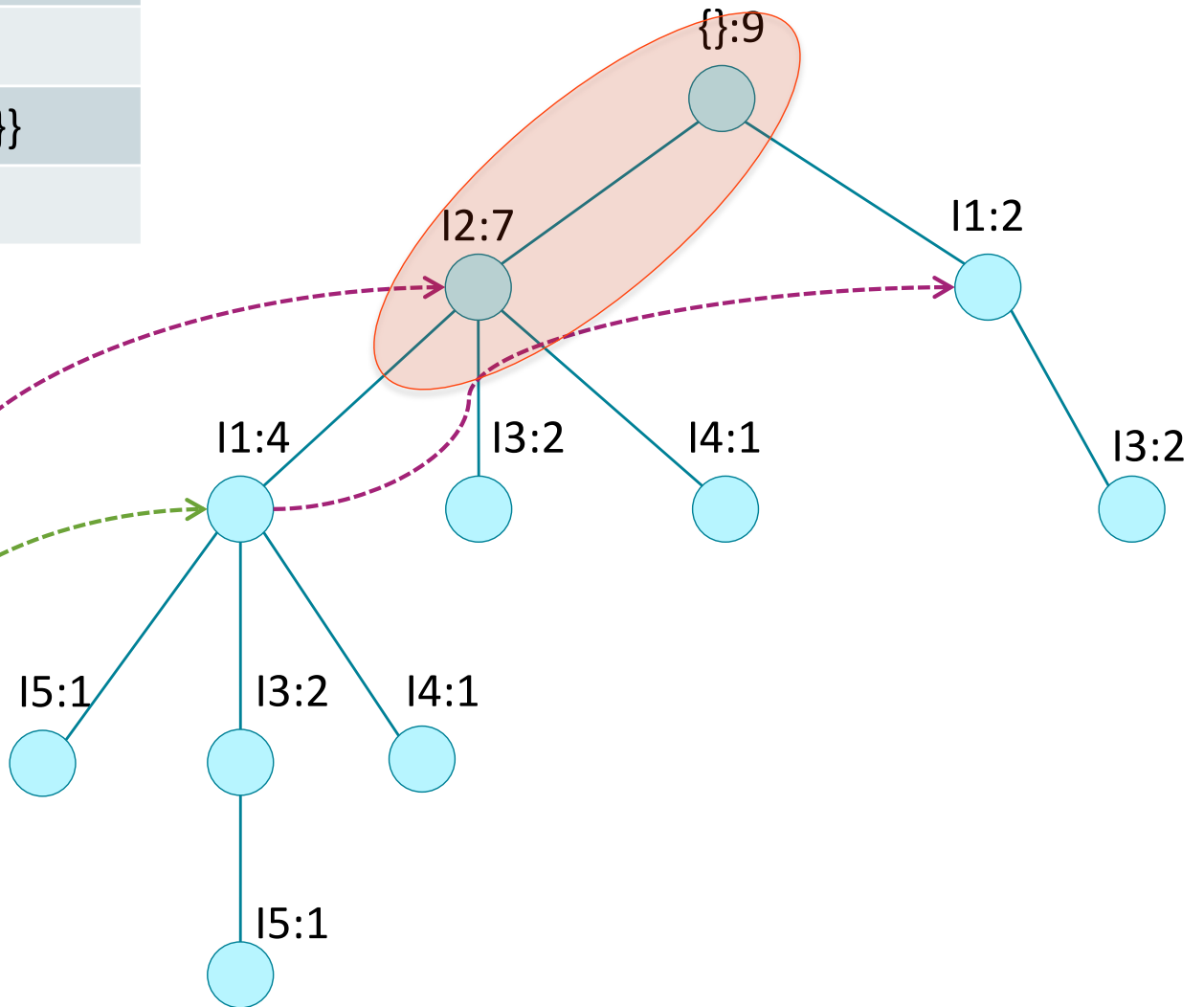
Header Table	
I2	
I1	
I3	
I4	
I5	



Construct Conditional Pattern Base

Item	Cond. Pattern
I5	{{I2 I1: 1}, {I2 I1 I3: 1}}
I4	{{I2 I1: 1}, {I2: 1}}
I3	{{I2 I1: 2}, {I2: 2}, {I1: 2}}
I1	{{I2: 4}}

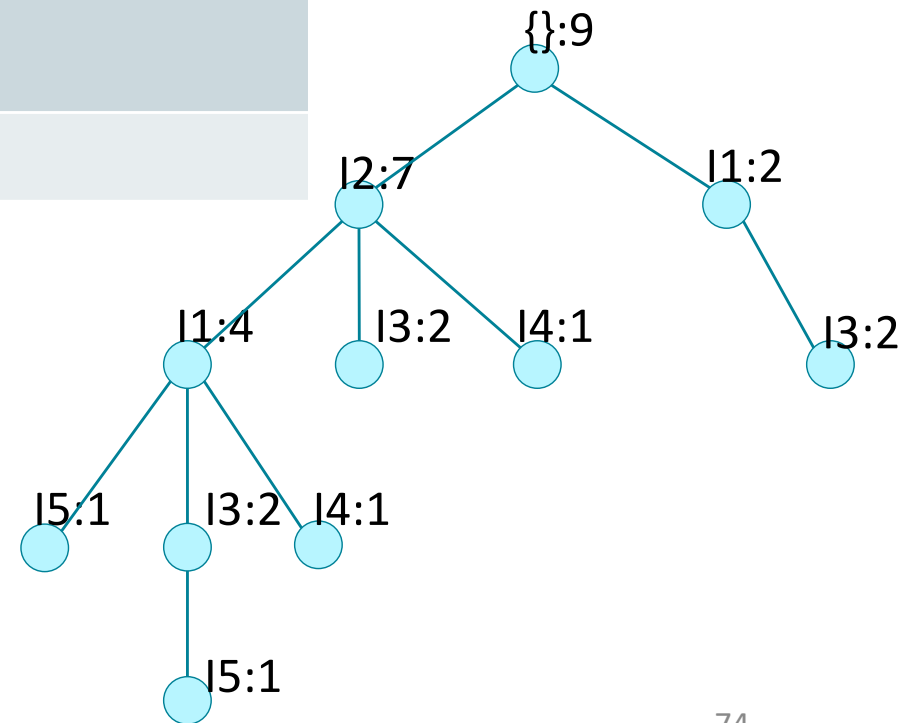
Header Table	
I2	
I1	
I3	
I4	
I5	



Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}		
I4	{(I2 I1: 1), (I2: 1)}		
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}		
I1	{(I2: 4)}		

- Create Cond. FP-tree using conditional patterns
- Frequent pattern with each suffix is generated by considering all possible combinations of the item and FP-tree

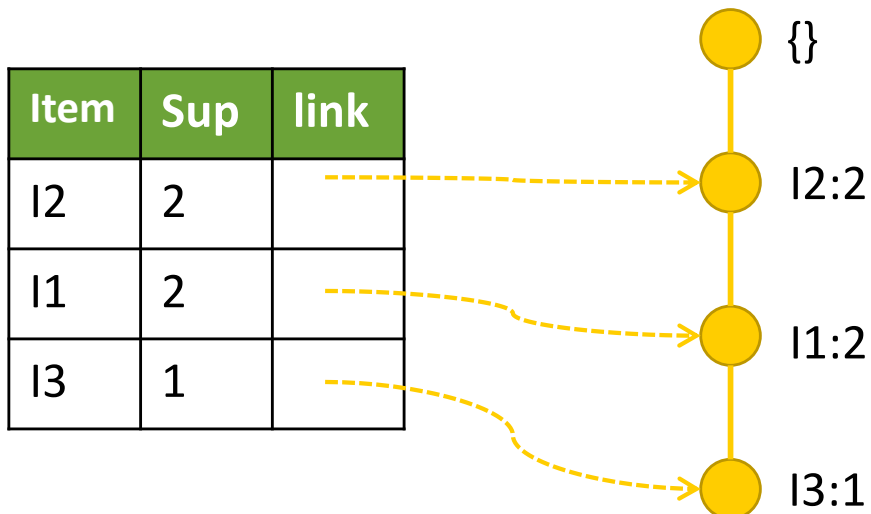


FPGrowth: Example

Part 3 – Create Conditional FP-tree and generate patterns

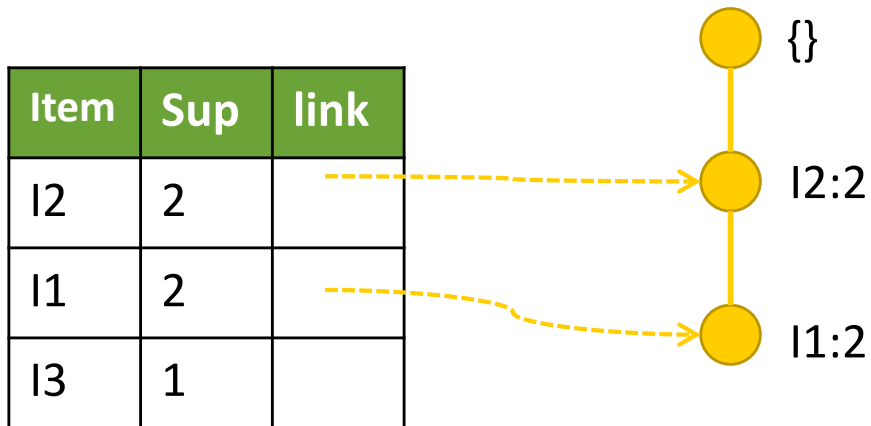
Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}		
I4	{(I2 I1: 1), (I2: 1)}		
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}		
I1	{(I2: 4)}		



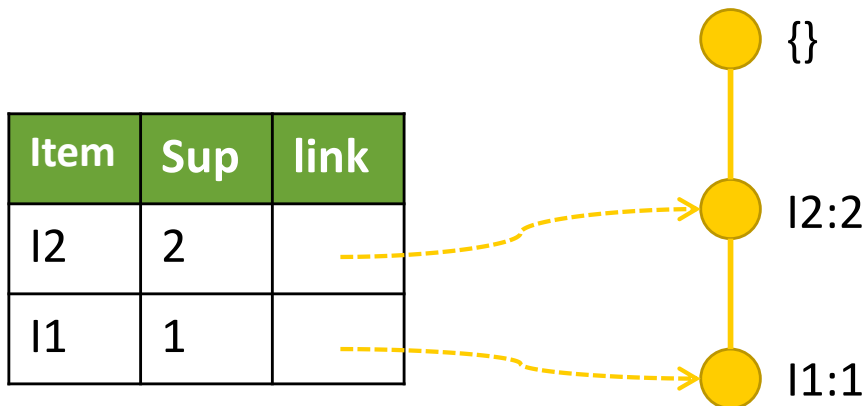
Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}	<I2: 2, I1: 2>	{I2 I5: 2}, {I1 I5: 2}, {I2 I1 I5: 2}
I4	{(I2 I1: 1), (I2: 1)}		
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}		
I1	{(I2: 4)}		



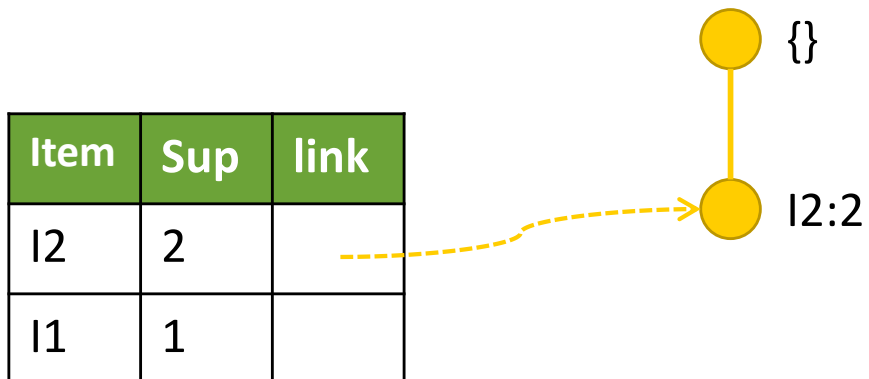
Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}	<I2:2, I1:2>	{I2 I5: 2}, {I1 I5: 2}, {I2 I1 I5: 2}
I4	{(I2 I1: 1), (I2: 1)}		
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}		
I1	{(I2: 4)}		



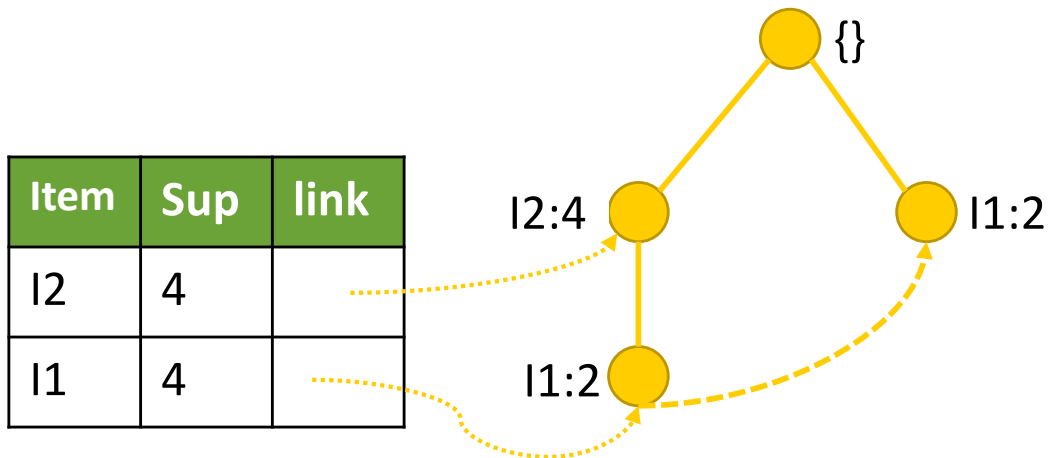
Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}	<I2:2, I1:2>	{I2 I5: 2}, {I1 I5: 2}, {I2 I1 I5: 2}
I4	{(I2 I1: 1), (I2: 1)}	<I2: 2>	{I2 I4: 2}
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}		
I1	{(I2: 4)}		



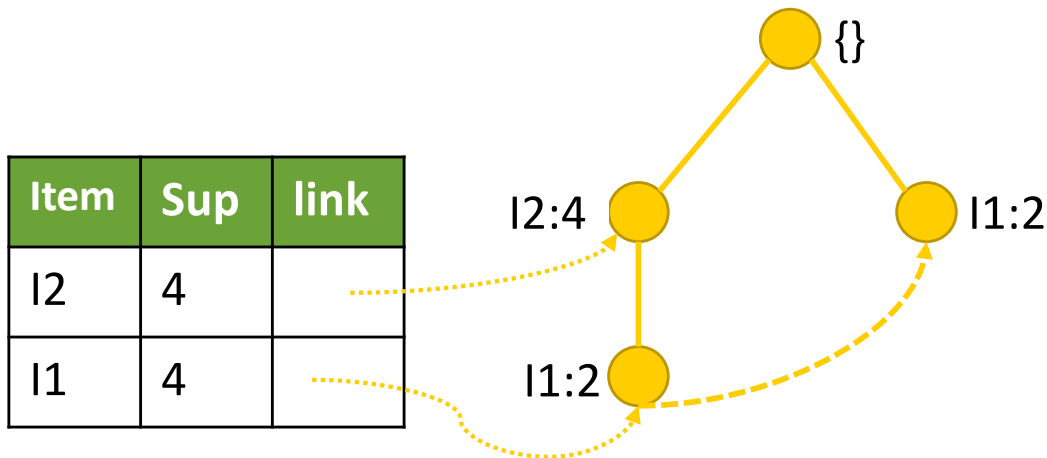
Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}	<I2:2, I1:2>	{I2 I5: 2}, {I1 I5: 2}, {I2 I1 I5: 2}
I4	{(I2 I1: 1), (I2: 1)}	<I2: 2>	{I2 I4: 2}
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}	<I2: 4, I1: 2>, <I1: 2>	
I1	{(I2: 4)}		



Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}	<I2:2, I1:2>	{I2 I5: 2}, {I1 I5: 2}, {I2 I1 I5: 2}
I4	{(I2 I1: 1), (I2: 1)}	<I2: 2>	{I2 I4: 2}
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}	<I2: 4, I1: 2>, <I1: 2>	{I2 I1 I3: 2}, {I1 I3: 4}, {I2 I3: 4}}
I1	{(I2: 4)}		



Cond. Pattern Base -> Create Cond. FP-Tree

Item	Cond. Pattern	Cond. FP-tree	Frequent Pattern
I5	{(I2 I1: 1), (I2 I1 I3: 1)}	<I2:2, I1:2>	{I2 I5: 2}, {I1 I5: 2}, {I2 I1 I5: 2}
I4	{(I2 I1: 1), (I2: 1)}	<I2: 2>	{I2 I4: 2}
I3	{(I2 I1: 2), (I2: 2), (I1: 2)}	<I2: 4, I1:2>, <I1: 2>	{I2 I3: 4}, {I1 I3: 4}, {I2 I1 I3: 2}
I1	{(I2: 4)}	<I2: 4>	{I2, I1: 4}

Frequent Itemsets Identified:

{ 1, 2, 3, 4, 5, 12, 13, 15, 23, 24, 25, 123, 125 }

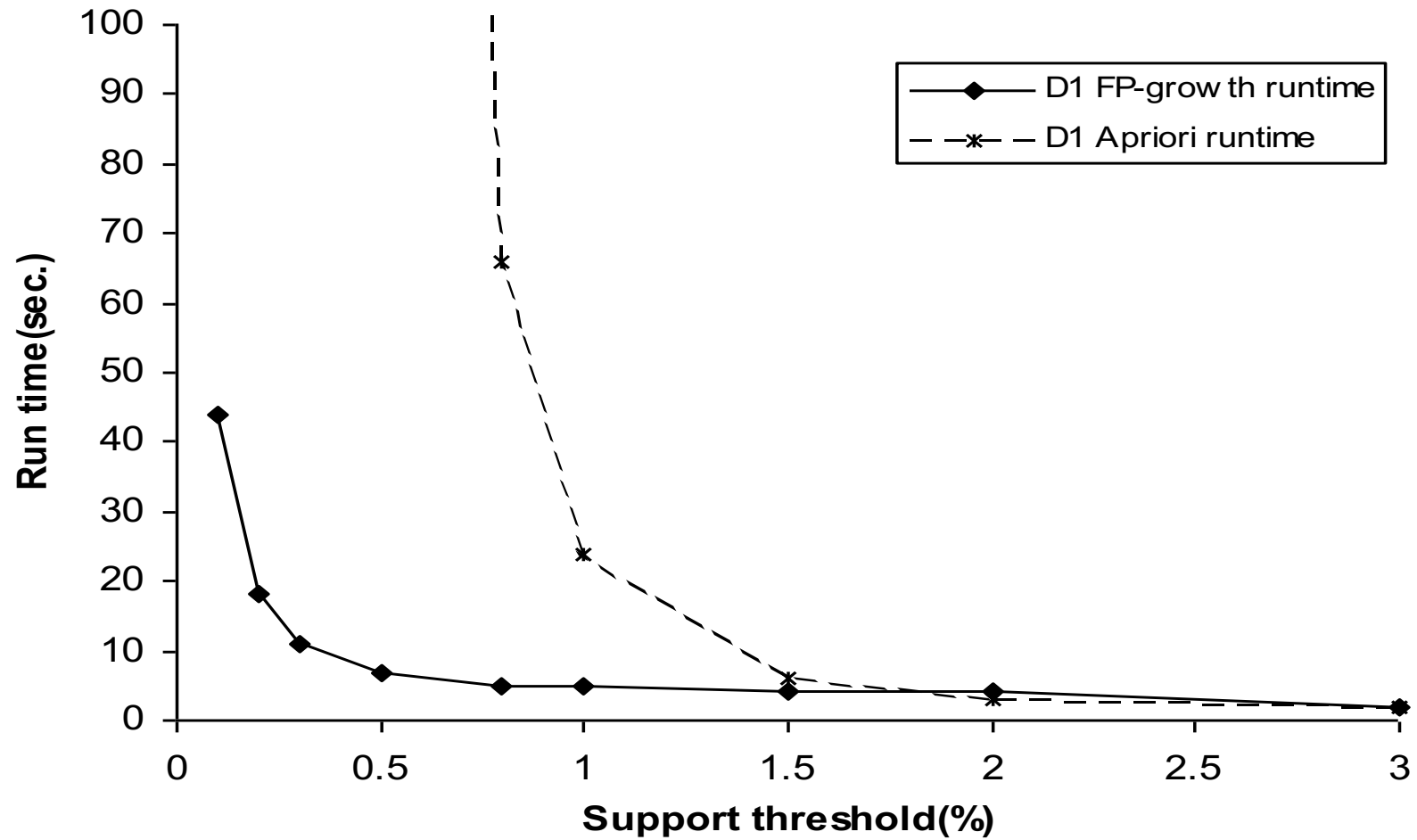
Benefits of FP-tree

- 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: the more frequently occurring, the more likely to be shared
 - Never larger than original database

Benefits of FPGrowth

- Performance study shows
 - FPGrowth is an order of magnitude faster than Apriori, also faster than tree-projection
- Reasoning
 - no candidate generation, no candidate test
 - use compact data structure
 - eliminate repeated database scan
 - basic operation is counting and FP-tree building

FPGrowth vs. Apriori



Other Improvements in Mining

- AFOPT (Liu et al., KDD 2003)
 - A “push-right” method for mining condensed frequent pattern (CFP) trees
- Carpenter (Pan et al., KDD 2003)
 - Mine data sets with small rows but numerous columns
 - Construct a row-enumeration tree for efficient mining
- Fpgrowth+ (Grahne and Zhu, FIMI 2003)
 - Efficiently using prefix trees, open-source implementation
 - ICDM 2003
- TD-Close (Liu et al., SDM 2006)

Additional Improvements in Mining

- Mining closed frequent itemsets and max-patterns
 - CLOSET (DMKD'00), FPclose, and FPMax (Grahne & Zhu, Fimi'03)
- Mining sequential patterns
 - PrefixSpan (ICDE'01), CloSpan (SDM'03), BIDE (ICDE'04)
- Mining graph patterns
 - gSpan (ICDM'02), CloseGraph (KDD'03)
- Constraint-based mining of frequent patterns
 - Convertible constraints (ICDE'01), gPrune (PAKDD'03)
- Computing iceberg data cubes with complex measures
 - H-tree, H-cubing, and Star-cubing (SIGMOD'01, VLDB'03)
- Pattern-growth-based Clustering
 - MaPle (Pei, et al., ICDM'03)
- Pattern-Growth-Based Classification
 - Mining frequent and discriminative patterns (Cheng, et al, ICDE'07)