



(\*) PAMPA-11D has been thoroughly evaluated on ~~both~~ real high dimensional datasets. ~~shown in terms of execution time,~~  
Experimental results show the efficiency and effectiveness of PAMPA-11D  
in ~~existing~~ performing the frequent-<sup>use</sup> mining with good load balancing  
on ~~multiple~~ <sup>multiple</sup> processors.

Nevertheless, many ~~preconditions~~ in security domains such as biometrics

applications / ~~have~~ <sup>have</sup> a page number

or networking. often requires to deal with this type of data. For instance, most gene expression datasets are characterized by a large number of items (related to tens of thousands of genes) and a few records (one transaction per patient or tissue). Many applications in computer vision deal with high-dimensional data, such as face recognition. Some state-of-the-art studies have built this type of large datasets measuring the occupancy of different car lanes, car-

transition then describes the occupancy rate in a caprot location and in a given timestamp [4]. In the networking domain, instead, the heterogeneity environment provides many different datasets characterized by high-dimensionality data, such as URL, reputation, advertising, and social network datasets [1, 5].

The work introduces PaMPa-HD, a parallel MapReduce-based frequent closed itemset mining algorithm for high-dimensional datasets, based on the Carpenter algorithm. PaMPa-HD outperforms the single-machine Carpenter implementation and the best state-of-the-art distributed approaches in total execution time and memory support threshold. Furthermore, the simple live PaMPa-HD design mentioned takes into account critical design aspects such as load balancing and robustness to memory issues.

The paper is organized as follows: Section 2 introduces the frequent (clearest) trust-mining problem. Section 3 briefly describes the centralization version of Carpenter, and Section 4 presents the proposed Paillier-IND algorithm. Section 5 describes the experimental evaluations proving the effectiveness of the proposed technique. Section 6 provides a brief review of the state of the art, and Section 7 discusses possible applications of Paillier-IND. Finally, Section 8 introduces future works and conclusions.

To effectively deal with those high-dimensional datasets, novel and distributed approaches are needed.

a large number  
 of events characterised  
 by a variety of  
 features. Transcribed  
 into high  
 TMSI have  
 closely  
 compared  
 to  
 be  
 generated.

tid	ten
1	ab
2	ad
3	ac
4	af
5	ah

(a) Horizontal  
tion of D

$\mathcal{D}$	
tid	items
1	abc, closv
2	ad, chlp, cv
3	acch, oq, v
4	afv
5	ab, d, gl, qst

(a) Horizontal representation of  $\mathcal{D}$

$TT$	
item	tblst
a	1,2,3,4,5
b	1,5
c	1,3
d	2,5
e	2,3
f	4,5
g	5
h	2,3
i	1,2,5
o	1,3
p	2
q	3,5
r	2
s	1,5
t	3,5
v	1,2,3,4

(1\*) Transposed representation of  $D$

$T^*T^*_{(2,3)}$	
item	tidlist
a	4,5
c	-
b	-
v	1

(c)  $T^T_{(2,3)}$ : example of condition-transposed table

## 2. Frequent itemset mining background

Let  $\mathcal{I}$  be a set of items. A transactional dataset  $\mathcal{D}$  consists of a set of transactions  $\{t_1, \dots, t_n\}$  where each transaction  $t_i \in \mathcal{D}$  is a set of items (i.e.,

$t_i \in \mathcal{I}$ , and it is identified by a transaction identifier (*tid*). Figure 1a reports an example of a transaction dataset with 5 transactions. The dataset is used as a running example through the paper.

An itemset  $I$  is defined as a set of items (i.e.,  $I \subseteq \mathcal{I}$ ) and it is characterized by a *tblst* and a support value. The *tblst* of an itemset  $I$ , denoted by  $tblst(I)$ , is defined as the set of *tbls* of the transactions in  $\mathcal{D}$  containing  $I$ , while the support of  $I$  in  $\mathcal{D}$ , denoted by  $sup(I)$ , is defined as the ratio between the number of transactions in  $\mathcal{D}$  containing  $I$  and the total number of transactions in  $\mathcal{D}$  (i.e.,  $tblst(I)/|\mathcal{D}|$ ). For instance, the support of the itemset  $\{mnb\}$  in the running example dataset  $\mathcal{D}$  is 2/5 and its *tblst* is  $\{1, 3\}$ . An itemset  $I$  is considered frequent if its support is greater than a user-provided minimum support threshold  $minsup$ .

Given a transactional dataset  $\mathcal{D}$  and a minimum support threshold  $minsup$ , the frequent itemset mining [6] problem consists in extracting the complete set of frequent itemsets from  $\mathcal{D}$ . In this paper, we focus on a valuable subset of frequent itemsets called frequent closed itemsets [3]. Closed itemsets allow representing the same information of traditional frequent itemsets in a more compact form.

A transactional dataset can also be represented in a vertical format, which is usually a more effective representation of the dataset when the average number of items per transactions is orders of magnitude larger than the number of transactions. In this representation, also called transposed table *TT*, each row consists of an item  $i$  and its list of transactions, i.e.  $TransList(i)$ . Let  $r$  be an arbitrary row of *TT*, *r.indlist* denotes the listset of row  $r$ . Figure 4 reports the transposed representation of the running example reported in Figure 3.

aterei la figura 13

in Figure 1a.

Given a transposed table  $TT$  and a tidlist  $X$ , the conditional transposed table of  $TT$  on the tidlist  $X$ , denoted by  $TT|_X$ , is defined as a transposed table such that: (1) for each row  $r_i \in TT$  such that  $X \subseteq r_i.tidlist$  there exists one tuple  $r'_i \in TT|_X$  and (2)  $r'_i$  contains all tids in  $r_i.tidlist$  whose tid is higher than any tid in  $X$ .

For instance, consider the transposed table  $TT$  reported in Figure 6a. The projection of  $TT$  on the tidlist  $\{2, 3\}$  is the transposed table reported in Figure 1c.

Each transposed table  $TT|_X$  is associated with an itemset composed by the items in  $TT|_X$ . For instance, the itemset associated with  $TT|_{\{2,3\}}$  is  $\{activity\}$  (see Figure 1c).

### 3. The Carpenter algorithm

3. The Carpenter algorithm to perform

As discussed in section 6, the most popular techniques (e.g., Apriori [7] and FPGrowth [8]) adopt the itemset enumeration approach ~~to generate~~ the frequent itemsets. However, itemset enumeration revealed to be ineffective with datasets with a high average number of items per transactions [3]. To tackle this problem, the Carpenter algorithm [3] was proposed. Since the Car-

In this problem, the Carpenter algorithm [3] was proposed. Specifically, Carpenter is a frequent itemset extraction algorithm devised to handle datasets characterized by a relatively small number of transactions but a huge number of items per transaction. To efficiently solve the itemset mining problem, Carpenter adopts an effective depth-first transaction enumeration approach based on the transposed representation of the input dataset. To illustrate the centralized version of Carpenter, we will use the running example dataset



Recall that in

—

5

Carpenter performs a depth first search of the enumeration tree to mine the set of frequent closed itemsets. Referring to the tree in Figure 2, the depth first search would lead to the visit of the nodes in the following order:  $\{1\}$ ,  $\{1,2\}$ ,  $\{1,2,3\}$ ,  $\{1,2,3,4\}$ ,  $\{1,2,3,4,5\}$ ,  $\{1,2,3,5\}$ ,  $\{\dots\}$ . For each node, Carpenter applies a procedure that decides if the itemset associated with that node is a frequent closed itemset or not. Specifically, for each node, Carpenter decides if the itemset associated with the current node is a frequent closed itemset by considering: 1) the tallest  $X$  associated with the node, 2)

Sovereign us

Free the

simtel 28 11 40 40 00

Computer - Poi

eggs - pro.

superhero performs

steps

1) Building

## trans section

## majoration

tree by explosive

2) ~~Adapt first~~

2nd Visiting

the tree  
through a depth

first search. ~~Different primary sales are~~

~~exploited to avoid the cost of free internet~~

FAI UN CHECK X VERIFICARE CHE

Ingravidire  
un po' il farto  
(io non leggo :-))

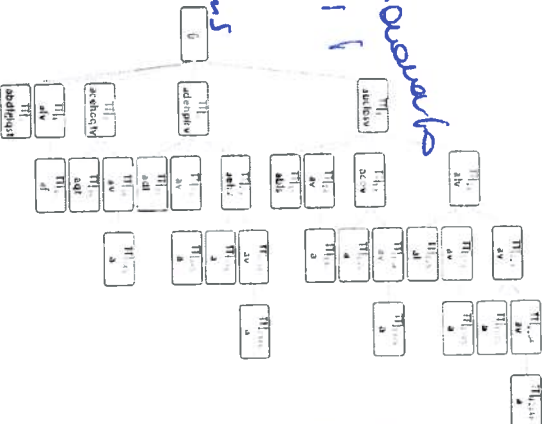


Figure 2: The transaction enumeration tree of the running example dataset in Figure 1a.

For the sake of clarity, no pruning rules are applied to the tree.

the conditional transposed table  $TT_X$ , 3) the set of frequent closed itemsets found up to the current step of the tree search, and 4) the enforced minimum support threshold (*minsup*). Based on the theorems reported in [3], if the itemset  $I$  associated with the current node is a frequent closed itemset then

- $I$  is included in the frequent closed itemset set. Moreover, by exploiting the analysis performed on the current node, part of the remaining search space (i.e., part of the enumeration tree) can be pruned to avoid the analysis of nodes that will never generate new closed itemsets. At this purpose, three pruning rules are applied on the enumeration tree, based on the evaluation performed on the current node and the associated transposed table  $TT[N]$ .
- **Pruning rule 1.** If the size of  $N$ , plus the number of distinct tids in the rows of  $TT[N]$  does not reach the minimum support threshold, the subtree rooted in the current node is pruned.
  - **Pruning rule 2.** If there is any tid  $tid$  that is present in all the tidlists of the rows of  $TT[N]$ ,  $tid$  is deleted from  $TT[N]$ . The number of discarded tids is updated to compute the correct support of the itemset associated with the pruned version of  $TT[N]$ .
  - **Pruning rule 3.** If the itemset associated with the current node has been already encountered during the depth first search, the subtree rooted in the current node is pruned because it can never generate new closed itemsets.

The tree search continues in a depth first fashion moving on the next node of the enumeration tree. More specifically, let  $tid$  be the lowest tid in the tidlists of the current  $TT[N]$ , the next node to explore is the one associated with  $N' = N \cup \{tid\}$ .

Among the three rules mentioned above, pruning rule 3 assumes a global knowledge of the enumeration tree explored in a depth first manner. This

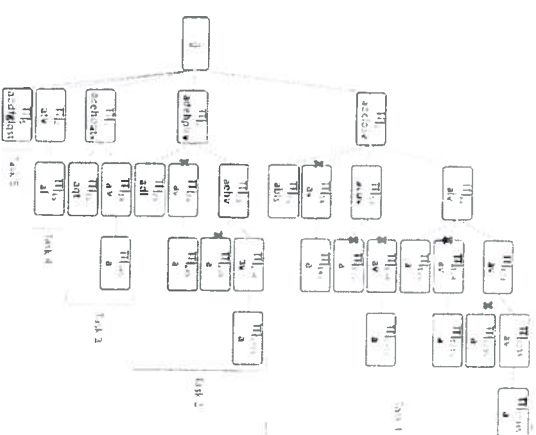


Figure 3: Running, for example, each node expands a branch of the tree independently. Pruning rule 1 and 2 are not applied. The pruning rule 3 is applied only within the same task: the red crosses on the edges represent pruned nodes due to local pruning, rule 1, e.g., the one on node (2,4) represents the pruning of node (2,4).

as detailed in section 4, is very challenging in a distributed environment that adopts a shared-nothing architecture, like the ones we address in this work.

#### 4. The PaPiPe-HD algorithm

Given the complete enumeration tree (see Figure 2), the centralized Carpenter algorithm extracts the whole set of closed itemsets by performing a depth first search (DFS) of the tree. Carpenter also prunes part of the search space by applying the three pruning rules illustrated above. The PaPiPe-HD algorithm proposed in this paper splits the depth first search process in a set of (partially) independent sub-processes, that autonomously evaluate sub-trees of the search space. Specifically, the whole problem can be split by assigning each subtree rooted in  $TT[X]$ , where  $X$  is a single transaction id in the initial dataset, to an independent sub-process. Each sub-process applies the centralized version of Carpenter on its conditional transposed table  $TT[X]$  and extracts a subset of the final closed itemsets. The subsets of closed itemsets mined by each sub-process are merged to compute the whole closed itemset result. Since the sub-processes are independent, they can be executed in parallel by means of a distributed computing platform, e.g., Hadoop. Figure 3 shows the application of the proposed approach on the running example. Specifically, five independent sub-processes are executed in the case of the running example, one for each row (transaction) of the original dataset.

Partitioning the enumeration tree in sub-trees allows processing bigger enumeration trees with respect to the centralized version. However, this approach does not allow fully exploiting pruning rule 3 because each sub-process works independently and is not aware of the partial results (i.e. closed itemsets) already extracted by the other sub-processes. Hence, each sub-process can only prune part of its own search space by exploiting its

"local" closed itemset list, while it cannot exploit the closed itemsets already mined by the other sub-processes. For instance, Task T2 in Figure 3 extracts the closed itemset *ac* associated with node  $TT[1,2,3]$ . However, the same closed itemset is also mined by T1 while evaluating node  $TT[1,2,3]$ . In the centralized version of Carpenter, the duplicate version of *ac* associated with node  $TT[1,2,3]$  is not generated because  $TT[1,2,3]$  follows  $TT[1,2,3]$  in the depth first search, i.e., the tasks are serialized and not parallel. Since pruning rule 3 has a high impact of the reduction of the search space, as detailed in Section 5, its inapplicability leads to a negative impact on the execution time of the distributed algorithm as described so far. To address this issue, we share partial results among the sub-processes. Each independent sub-process analyzes only a part of the search subspace; then, when a maximum number of visited node is reached, the partial results are synchronized through a synchronization phase. Of course, the exploration of the tree finishes also when the subspace has been completely explored. Specifically, the sync phase shares the partial results (i.e., nodes of the tree still to be analyzed and found closed itemsets) globally applying pruning rule 3. The pruning strategy consists of two phases. In the first one, all the transposed tables and the already found closed itemsets are analyzed. The transposed tables and the closed itemsets related to the same itemset are grouped together in a bucket. For instance, in our running example, each element of the bucket  $B_{ac}$  can be

- a frequent closed itemset *ac* extracted during the subtree exploration of the node  $TT_{1,1}$ ,
- a transposed table associated to the itemset *ac* among the ones that still have to be expanded (nodes  $TT_{1,2,3}$  and  $TT_{2,3,4}$ ).