

# Providing a new approach to improve in speed and longevity of the Internet of Things based on RFID

Arman Kavooosi Ghafi

Department of Computer  
Software, Central Islamshahr  
Branch, Islamic Azad  
University, Iran

Behnaz Farzi

Department of Computer,  
Buinzahra Branch, Islamic Azad  
University, Buinzahra, Iran

Helena Kojooyan Jafari

Department of Computer  
Software, Central Islamshahr  
Branch, Islamic Azad  
University, Iran

---

**Abstract:** Identification Technology Using Radio Frequency tags (RFID) is a very advanced technology that is fairly named the greatest revolution after the Internet. Internet of Things is based on this technology and, it will be rapidly prevailed. A set of constraints that lie ahead is the major challenges of development and application of RFID networks. One of the most fundamental concerns is tag Readers optimal deployment in large-scale RFID network planning (RNP), which leads to optimal performance and increase in lifetime and speed of network. With considering coverage, signal interference and load balance as optimization targets and determination of optimum, the establishment issue of tag reader is converted to compound multi-objective optimization problem. In this article, in order to find the answers of the problem, the particle swarm algorithm (PSO) Combination with multi-objective optimization, based on Pareto's theory MOPSO that is able to solve the problem with more than one objective, was used. Simulation results show that the algorithm MOPSO compared to the optimization algorithms coverage, signal interference and load balancing has been effective. Therefore, with optimal deployment of tag readers, overall performance of system is improved.

**Keywords:** RFID, MOPSO, Internet of things, lifetime, tag readers.

---

## 1. INTRODUCTION

In large-scale RFID network systems, fast and efficient collection of information stored on the tags has a great importance, regarding the usage of applications, but this problem is not yet been fully investigated, and the comprehensive solution for it has not been presented yet. Due to limitations of domain and coverage range of tag readers, the tag readers should be in a compact format, in order to completely coverage a certain area, and to use minimum number of tag reader that is impossible. There are some problems, which are imposed on the deployment of RFID network that can be cited as an example: to determine the number of deployed readers, seating and setting parameters for each reader that the issue of planning network (RFID RNP) in RFID with large scale at an optimal level, has seriously challenged. The RNP problem is a large-scale nonlinear optimization problem with many variables, constraints and objectives proven that the RNP is a NP-Hard problem. Although the major factors in RNP is deploying reader for maximum coverage. Read speed of information and the RFID network longevity is one of the major concerns that must be considered. Because if the reader and tags are considered identical, assuming a uniform statistical distribution of tags in a network environment, Obviously, balancing the load between the readers, for example scope and balanced coverage volume, can lead to increased network speed and longevity. Accordingly in this article, with regard to coverage, interference signals and load balancing as a optimization goal, first of all we explain the inferential model of multi-objective optimization, and then using MOPSO improved algorithms, will be clarified, in order to have an increasing rate in the speed, and life length of internet of things based on RFID network.

## 2. Introduction of RFID

RFID is the concept for automatic identification of an object using radio signals -based storage and retrieval of remote data (Clauberg et al., 2013). Generally, RFID technology is helping of the following equipment to implement:

1. Tag
2. Tag Reader
3. Antenna - Signal Booster
4. Information Management Software
5. Database

Radio waves carry information between sender and receiver information devices. To the sender piece of information, tag and to receiver pieces of information, called a reader or tag reader. Labels often lie on the object. If the tags are classified according to their energy sources that are being used, we have three main types of them. Active, passive and semi-passive Tags. Active and passive tags have many differences, But it can be noted that active tags, receives Their energy needs from their mobile battery, While passive tags itself has not been a source of energy and for tripping should use the energy of electromagnetic waves emitted by the tag reader, and their range and scope of read is less than active tags (Finer et al., 2011). Passive tags are low cost and long lifetimes and also small dimensions. Semi-passive tags is also another type of label that in addition to using its internal battery, can also use wave energy emitted by the tag reader. Antenna is used to transmit radio signals between the tag reader and tag that is used both for Tags and tag reader. There is a Data management software in order to process the collected data. This special software is usually on a local server that allows

the exchanged Data with the tag reader, to be collected and processed, and stored in a database, and also be restored if needed. RFID technology can be an alternative for the barcodes. In fact, RFID is more than just a bar code, because it has an automatic system scanner. This two technologies have major differences. The two main differences can be elaborated in the following ways; RFID technology is capable of carrying large volumes of information, and also it does not need to have line of sight for data collection and communication. (Hvlmyvyst and Stephenson, 2006).

### 3. Introduction of the Internet of Things

The report prepared by the analysts of policy and strategies unit of ITU, has a look at the next step in a continuous and always ship (ITU, 2005). Based on this the new technologies, such as RFID, smart calculation promising, world of networked, and interconnected equipment. At that time, everything from car cyclic to brush come into communications area that is the announcer of rising in a new era, and will lead current Internet (which contains data and people) to the Internet of Things.

## 4. Method

### 4.1 Optimization of Tag readers' deployment

In this section, the first theoretical topics of issue, the formula array, and then the proposed algorithm, based on PSO will be presented.

## 5. Problem formulation

In large-scale RFID applications systems, three essential limitations should be considered comprehensively; increase overall performance and efficiency, prevention of data loss of labels, or tags which undoubtedly is unacceptable. In this field of study, Coverage is the main goal and the other two are directly related to this problem.

### A) The maximum coverage

In the traditional method of modeling coverage rates, Cover Rate is obtained by dividing the number of covered tags toward the total number of tags. If we assume that T is a set of tags, based on coverage area, R be a set of tag readers, based by method or algorithm which is used in environment, E be the minimum threshold power of reliable receiving of ship between reader and tag,  $E_{R,T}$  is most optimal power of tag, to get information of T member, tag by r tag reader that is a member of R collections. When we want to define the coverage area for R series tag readers, only the tags that the amount of their  $E_{R,T}$  is higher than the minimum threshold of (E), and its similar amount is for any other of tag reader, in this case, it has better conditions than other tag readers which are taken into account in its range. Tags coverage area R can be applied as the number of tags that can be expressed in the following equation:

$$C_T(R) = \{T \in TS \mid E_{R,T} \geq E, E_{R,T} \geq E, E_{R',T} (\forall R' \in RS, R' \neq R)\}$$

The optimization objective of coverage is the deployment of finite number of tag readers for full coverage area which is given below.

$$TS = \bigcup_{R \in RS} C_T(R), \text{with minimize of } |RS|$$

Therefore, the coverage rate can be defined as follows (equation (3)), that the goal is maximization of it:

$$\text{Max } f_{CR} = \frac{\sum_{R \in RS} C_T(R)}{|TS|}$$

Where  $\sum_{R \in RS} C_T(R)$  is the number of active tags used in the work area.  $|TS|$  is the number of tags widespread in the environment. If  $f_{CR} = 1$ , then the system has reached its optimal coverage.

### B) reduce signal interference

Reduce signal interference can be considered as equivalent to a reduction of undesirable interaction of multiple tag readers. If interference or overlapping areas will increase by several readers as a result signal interference level will go up. So by reducing interference or overlapping areas, signal interference can also be reduced. Of course the number of tags based on overlapping regions has a significant impact on network performance. Assuming this to be done in this area of tag. But because of the presence or absence of tags, there is not prognosis in these areas. So for optimum director should be on reducing weaving or overlapping efforts. For each tag in a particular scene, if the power received by the other Tag readers be less than the best power  $E_{R,T}$ , but larger than the minimum threshold of E, so  $E_{R,T} \geq E_{R',T} \geq E$ , then there is the possibility of signal interference for the tag. The level of signal interference for tags can be defined by  $\varphi(T)$  and as follows:

$$\varphi(T) = \sum (E_{R',T} - E), \text{ where } E_{R,T} \geq E_{R',T} \geq E$$

Thus, for each tag, if the value of E corresponding to the its overlying Tag readers be more than amount of E related to all other of Tag readers-based in environment, simply put, when

$$\max\left(\frac{E_{R,T}}{E_{R,T} + \varphi(T)}\right) \text{ is the maximum, can be sure that}$$

Signal interference with other tags does not happen in practice. With regard to the proposed, the goal function of reduce signal interference can be described as follows by equation (5):

$$\text{max } f_{SI} = \frac{\sum_{T \in TS} \left(\frac{E_{R,T}}{E_{R,T} + \varphi(T)}\right)}{|TS|}$$

So when  $f_{SI} = 1$  and to reach its maximum value, the system reaches its optimum level in terms of reducing signal interference.

**C) Increasing the lifetime of the network with load balancing**

In line with the objectives mentioned above, set of observed tags should be covered properly and balanced by based tag readers. The ship between tag and its covering Tag readers via RF signal is carried out, used energy or power needed to any exchange of information between tag and tag reader with  $C_{R,T}$  is shown. It is clear that the total communication cost or power consumption for each R tag reader will be as the sum of the total cost of communications coverage area, and defined as  $C_R = \sum_{R=\gamma(T)} C_{R,T}$  which

$R = \gamma(T)$  indicates that the T tag is covered by R tag reader. The aim of load balancing is finding an optimal distribution of tags to minimize the total power consumption ( $C = \sum_{R \in RS} C_R$ ) or is reduction of the cost and power consumption, which result in increasing of the network life time. To simplify our analysis, we put the amount of energy consumed for each exchange of information between tag and reader equal to one (Units). A simple system for energy consumption by equation (6) expressed as the number of tags by the reader R is covered

$$\text{Min } C = \sum_{R \in RS} C_R = \sum_{R \in RS} n_R$$

The number of assigned tags are used to serve as a solution for load balancing, which is trying to minimize

$$\text{Min } \prod_{i=1}^{|RS|} \left( \frac{1}{n_{R(i)}} \right) \text{ and } n_{R(i)} \text{ is the number of tags}$$

covered by the (i) tag reader. If load balancing is the only goal of system optimization, this solution will have the best impacts and results. But, because there are a lot of variables that are fine. Its use does not seem useful for multi-objective programs. In this paper, according to information entropy theory, and the corresponding  $S = -\sum P_i \log p_i$

Replace The P value with  $\frac{n_{R(i)}}{|TS|}$  between and define load

balancing function with equation (7), where  $|RS|$  represents the number of tags, and  $|TS|$  represents the number of tag readers. To simplify calculations we rewrite equation (7) into equation (8), where  $\ln |RS|$  is used for ease of use and normalization. Once, it means that all tags are just balanced in covering readers. When  $f_{LB} = 1$ , it means that all tags are just balanced in covering readers.

$$\text{Max } f_{LB} = -\sum_{i=1}^{|RS|} \frac{n_{R(i)}}{|TS|} \ln \left( \frac{n_{R(i)}}{|TS|} \right)$$

$$\text{Max } f_{LB} = -\sum_{i=1}^{|RS|} \frac{\frac{n_{R(i)}}{|TS|} \ln \left( \frac{n_{R(i)}}{|TS|} \right)}{\ln(|RS|)}$$

**6. The proposed solution based on MOPSO**

As described in the previous section, in order to fulfill the formula functions of the three goals of the establishment of readers on the RFID systems at large scale, in order to identify and maximum coverage, and reduce interference and consequently increase network speed, as well as covering the balance to reduce adjusted cost and power consumption, and extend the life of the network, can be considered as the combination of optimized multi-objective optimization. PSO algorithm is a popular multi-dimensional optimization algorithm and is a proper technique for this purpose, the application of simple, high quality solutions and results, and efficient computing and high speed of convergence are PSO algorithm's strengths. PSO algorithm used at the position of each particle as a deployment solution is code reader, assuming  $m = |RS|$  is the number of readers, the position of the particle i define as  $x_i = \{x_{i,1}, x_{i,2}, \dots, x_{i,2m}, x_{i,2m}\}$  in which each pair number represents the coordinates of a reader, so long  $|x_i| = 2m$ , mathematically speaking,

$$(x_{i,2k-1}, x_{i,2k}) \in x_i, (k = 1..m)$$

That reflects the position of coordinate's reader k in the i-th particle. Particle i fitness function (fitness (Xi)), is criterion for evaluating the position or solution of Xi in performance environments of network with m reader. In the proposed method the three objectives of functions that were formulated in the previous sections Instead of fitting function, we use than Pareto theory and we choose a set of optimal answer as the Pareto set.

$$\left\{ \begin{array}{l} v_i^{(t+1)} = wv_i^t + c_1 \gamma (pb_i^{(t)} - x_i^{(t)}) + c_2 R (gb^{(t)} - x_i^{(t)}) \\ x_i^{(t+1)} = x_i^{(t)} + v_i^{(t+1)} \end{array} \right\}$$

In The PSO algorithm, each particle by a path with concept of timing of the allocation of resources be coded. The main idea of PSO is finding the best time of movement of the entire particle with the energy assessment, and implementation cost and time of implementation. Each particle contains a status indicating timing solutions and a velocity vector that represents the direction and amount of motion of the particle

respectively. Position and speed of movement of the particles of the wishes expressed by the following expressions: Where  $V_i^t$  represents the speed and  $P_i^t$  represents the position of the particle  $i$  at repetition  $t$  and respectively  $gb$  and  $pb$  are the best position or solution, and local (individual) and best overall (general).  $r$  and  $R$  are probabilistic random values in the range  $[0,1]$  and  $w$  and  $c1$  and  $c2$  are weighting parameters.

To select the best local and global position, When more than one purpose in mind, the selection becomes a major challenge that To solve the problem we use of the concept of Pareto's theory to compare and choose the best timing or position of the particle in terms of the number of targets, . The process and the proposed algorithm is given below

## 7. PSO multi-objective algorithm

Step (1): we put value  $t=0$  for the first iteration. All  $k$  particles in the Stset community using the greedy algorithm are initialized.

Step (2): for each particle calculations using the equations (previous section) and the description of paragraph (2) of this section do and produce new particles  $K$  wishes. Total new and old particle forms  $2k$  of community of Tset temporary memory.

Step (3): Using the theory of Pareto and formulated three objective function, we have found the optimized particles or the position of local optimization ( $Pbest$  s) and it named with  $POset$  and we show the number of member of optimal set with  $Npo$ . If  $Npo < K$  we go to step (4), otherwise the step (5) are going.

Step (4): we select the number of  $K-Npo$  particles from particles that are not optimal Pareto to add to the collection  $POset$  to complete a set number and give up to a number  $K$  to form set of particle  $ind$  The next step  $St+1set$  .

Step (5): we will select randomly  $K$  number of chosen Pareto optimal set of  $POset$  to create next step complex particle swarm  $St+1set$  set.

Step 6: Check the termination condition of the algorithm, consider that obtained conditions is favorable or the number of iterations has been completed or not. In case of termination condition go to step (8), otherwise the step (7) are going.

Step (7): The amount of  $t = t + 1$  up to date and to step (2) are going.

Step (8): algorithm Completed and we refer to Pareto optimal set of  $POset$  to decide and choose the nationwide optimal solution of  $Gbest$ .

## 8. Results

In this section will present the results of simulation using the proposed establishment of the reader MOPSO (multi-objective particle swarm algorithm), to analyze and compare the results thereof. These simulation was conducted by using MATLAB in 2011.

## 9. Evaluation criteria

Since the aim of this paper is to determine the location of readers on the RFID network to increase the speed and reduce the cost and increase longevity. And in this regard raised three main objectives that include increased maximum coverage in order to read the maximum number of tags and reliability of the network, Reduce signal interference in areas of overlap to reduce energy waste And thus extend the life of network and increase speed of network with correct and reliable connections Load balancing, or balance between reader and number and cost of tag readers Resulting in increased longevity and speed of the network because of the network secure connections. Thus for the performance evaluation of MOPSO proposed optimized algorithm , The secure relevance rate or lack of signal interference that is calculated by the number of tags identified by the total number of tags available in the network environment or the coverage, as the inverse probability function of the of relevance rate And to measure of load balancing of the logarithmic covered tags  $rat$  by any Tag reader to the total number tag readres In comparison to the number of tags, tag readers, reader coverage and algorithm parameters related to PSO (population, weighting coefficients, etc.) and multi-objective optimization (sink size and the number of objective functions, etc.) is used.

## 10. Simulation parameters

MATLAB programming languages is used to test and validate the proposed algorithm. PC specifications  $ci5$  4 GB Intel processor and operating system Windows 7 was selected. The number of repeat elections was simulated 500 times, to simulate different scenarios were used to assess. The proposed algorithm with different values for the number of tags (200, 400, 600 and 800) was tested. The size of the network environment was varies according to the number of tags and was calculated from

$$\text{equation } R^2 = [0, \sqrt[2]{|TS|}] * [0, \sqrt[2]{|TS|}]. \text{For}$$

example, for the number of tags to 200, a circumference of  $14 * 14$  was chosen to do so for the number of tags 400 tense environment,  $20 * 20$ , for the number of tags 600, dimensions network environment  $24 * 24$  and to the end to the number of tags 800 twisted aspects of the environment  $28 * 28$  were selected. Also experiment with different values of tag readers (3,5,10, ..., 30) and evaluation of their role in the RFID system performance took place. All evaluations by three indicators related to the three objectives formulated, were tested. Other settings applied to simulate is given in the table (below).

**Table 1. Simulation parameters**

## 11. The simulation results

### 11.1 Overall assessment of the proposed algorithm

In the proposed algorithm in the first step 200 tags into an area with the size of 14 in 14 randomly and evenly were distributed. Below figure shows the distribution of tags. Each tag has been shown in blue as small hollow circle. Tag readers with red + sign marked that the locations of them determined by algorithm-based recommendations MOPSO been. we repeat it for different values of tag and tag readers so that the results is given in the below picture .

Value	The parameter
200,400,600,800	The number of tags
3-30	The number of tag readers
14,20,24,28	Dimensions of Network
2	The minimum threshold power
10	The maximum power of the tag readers
200	Initial population of PSO algorithm
100	Sink size for MOPSO
1	C1 for PSO (personal learning coefficient)
2	(C2 for PSO (Global learning coefficient)
0.1	Alpha coefficient
2	Beta coefficient
0.1	Mutation rate
Value	The parameter
200,400,600,800	The number of tags
3-30	The number of tag readers
14,20,24,28	Dimensions of Network
2	The minimum threshold power
10	The maximum power of the tag readers
200	Initial population of PSO algorithm
100	Sink size for MOPSO
1	C1 for PSO (personal learning coefficient)
2	(C2 for PSO (Global learning coefficient)
0.1	Alpha coefficient
2	Beta coefficient
0.1	Mutation rate

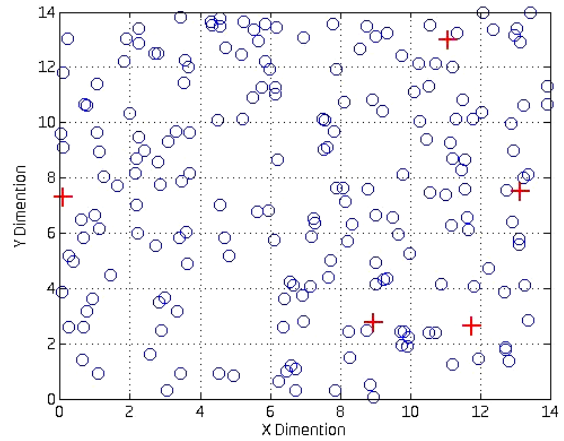


Figure 1: Distribution of 6 readers for random distribution network with 200 tags

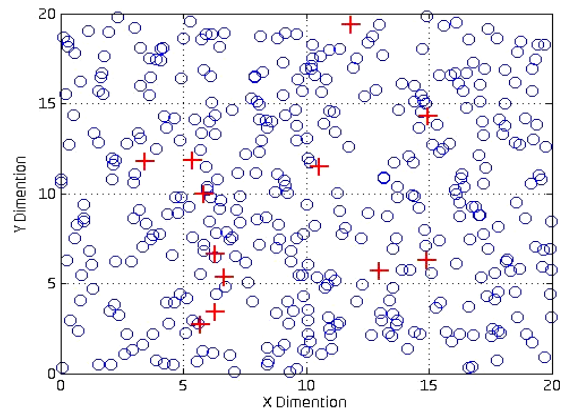


Figure 2: Distribution of 12 readers for random distribution of network with 400 tags

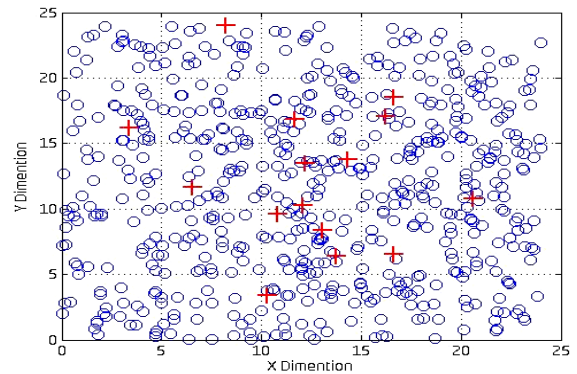


Figure 3: Distribution 15 tag reader for network with 600 tags by random distribution

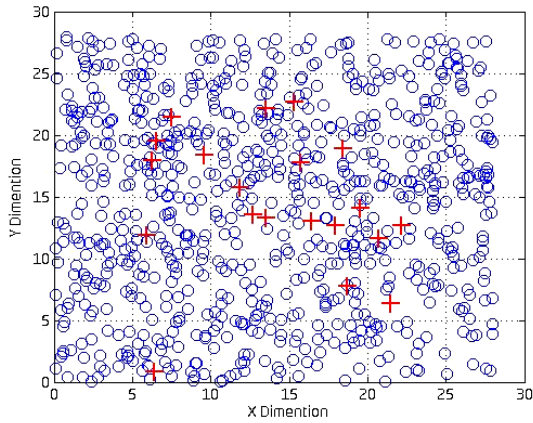


Figure 4: Distribution of 20 Tag reader for random distribution of network with 800 tags

### 12. Efficiency assessment on the cover tags

Since the basic objective of the proposed algorithm, is maximum coverage RFID tags by a reader network, at this stage of evaluation first the coverage of tags with different amounts of tag readers for a 200 number fixed number tag reader within the network environment Investigated discussed that the results is given in the table (below) and figure (below)

Table 2: The coverage rate of 200 tags by a different number for Tgkhvan

Reader	Coverage
3	0.8
5	0.865
10	0.94
12	0.965
15	0.955
17	0.97
20	0.955
23	0.96
25	0.985
27	0.96
30	0.99

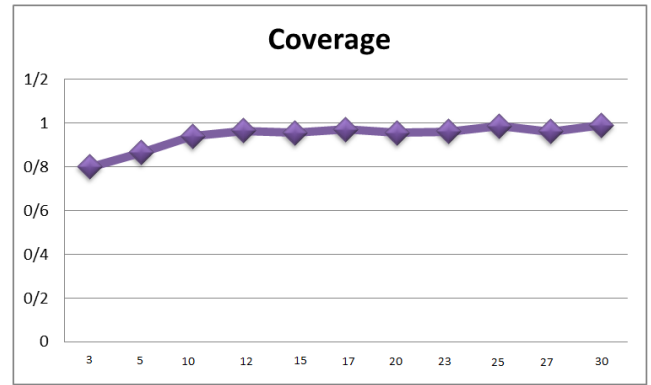


Figure 5: tags Coverage rate (tag number: 200) than the number of readers

In the figure (above), the horizontal axis and vertical axis represents the number of Tag reader and coverage rates in the range [01] shows. The table (above) and shape (above), it is observed that when the number of Tag reader that are less than 12 have coverage below 90%. With the increasing of number of Tag readers, amount of coverage increased, but Tag reader increase from 10 to 30 does not make a significant difference in the rate of coverage, and it is inferred that to cover the 200 tag, in the RFID network environment with dimensions of 14 \* 14, the number of 10 to 12 Tag reader will suffice.

### 13. Efficiency assessment to reduce signal interference and increase the speed and reliability of the grid

The second purpose of this article was reduction of the maximum signal coverage, to evaluation the proposed method of image values derived from the formula used to signal interference. Therefore, the results obtained, show the rate of signal interference for different numbers Tag reader for a fixed number 200 indicates the number of tags in a network environment. The table (below) and form (below) shows the results.

Table 3: Rates of non-interference, tag of number of for Tag reader of 200

Reader	~Signal interference
3	0.99558
5	0.99517
10	0.99315
12	0.99262
15	0.99178
17	0.99081
20	0.99009

23	0.98917
25	0.98827
27	0.98826
30	0.98733

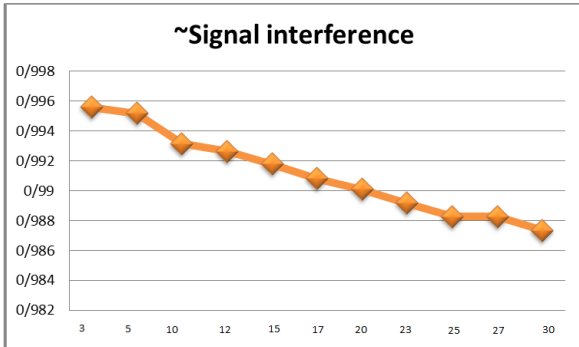


Figure 6: The amount of signal interference on the cover tags (tag of number of: 200) than the of number of readers

In the figure (above), the horizontal axis and vertical axis represents the number of Tag reader and load rates in [01] shows. Of the table (above) and shape (above), it is observed that when the number of Tag reader are less than 12, the number of readers, rate of balance of the load increases significantly, but this amounts between 12 to 30 the difference is subtle number and the corresponding fluctuations could be due to random distribution is different in each simulation. The results also found that 11 to 12 pieces for readers on an area the size of 14 \* 14 and the most optimal tag of number of that are 200.

#### 14. Efficiency assessment in the three goals compilation

Figure (below) show The results of proposed algorithm to achieve three objectives defined in the multi-objective optimization Pareto collection , in the form (below) is observed that when of number of tags is 12, the convergence caused in the three objectives (coverage, non-interference and load balance, and the optimum mode occurs, thus increasing more than this size is not necessary for the reader, and cause increasing cost. Therefore the best mode is intended for environments with tag of number of 200 number, is the 12 reader.

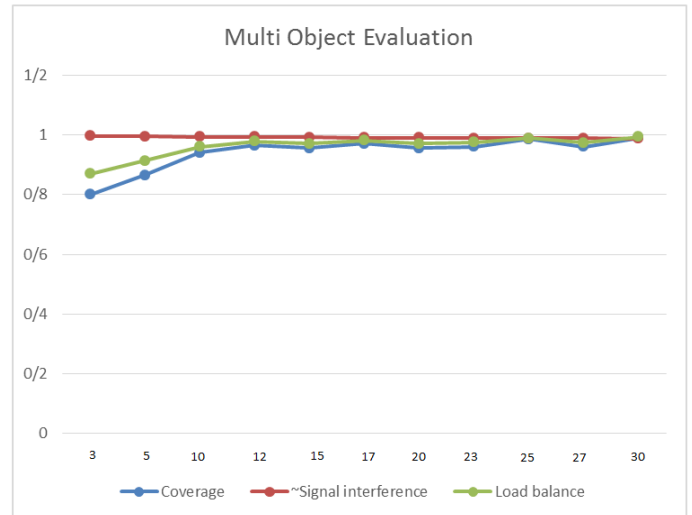


Figure 7: The three goals of efficiency rate (tag of number of: 200) than the of number of readers

#### 15. Evaluate and compare the proposed algorithm with previous approaches

One of the important evaluations, is comparing the proposed algorithm with the previous method. In this evaluation of Ming Tao's, et al method in (2015) SA-PSO, Keyvan Kashkul Nejhad et al (2011) AFS and Mr. Lian Bo et al (2014) ABC for are chosen to be compared with the proposed method respectively. In the form (below), coverage rate are given for comparison algorithms. As can be seen in the figure, the proposed algorithm even in small numbers compared to the same Tag readers that are able to provide better coverage. When the number of readers is more than 10 the algorithms, the performance is similar, and difference is not tangible. And the proposed algorithm has been opened with a slight difference above the rest.

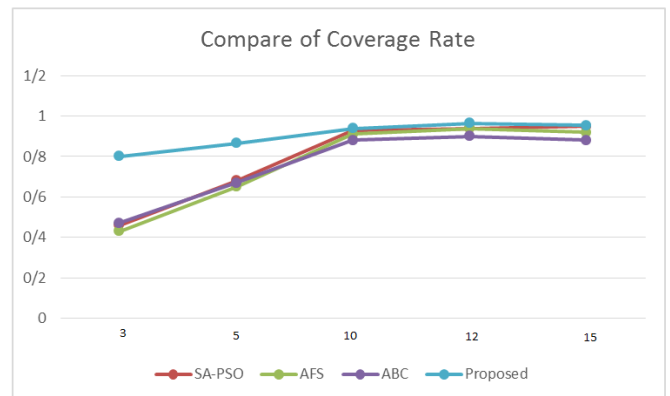


Figure 8: The algorithm coverage rates comparison to some readers

In the following comparison, the proposed algorithm in terms of solving the problem rate and reduce signal interference compared with the three previous algorithm. Figure (below) shows the results, in the figure ABC algorithms had the weakest performance, and AFS algorithm is in better shape. AFS performance and SA-PSO algorithm, and the proposed algorithm in the range of Tag reader number between 3 and 12 almost identical same performance and non-signal

interference is of the top of 99%, but when the number of readers increases the rate of non-interference in the AFS algorithm is reduced and by the end of this slope is proportional to the number Tag reader of decline. The SA-PSO algorithm until that the number of tag readers are less than 23 acted feet the feet of the proposed algorithm in good condition and with the passage of 23 has been accelerated its decline and its slope is steeper. All algorithms by increasing tag readers their rate of non-interference is reduced, due to the increasing of levels of overlap and read a tag (which is located in the overlap region), At the same time is by two or more readers. But Tag reader improper placement can cause an increase in overlapping regions is resulting in increased interference. Increase the rate of non-interference evaluation index. According to the multi-objective algorithm and using Pareto multi-objective optimization theory could take several targets in conjunction And the use of indicators covering and interfering signals simultaneously, reduce signal interference and increase the rate of non-interference evaluation index.

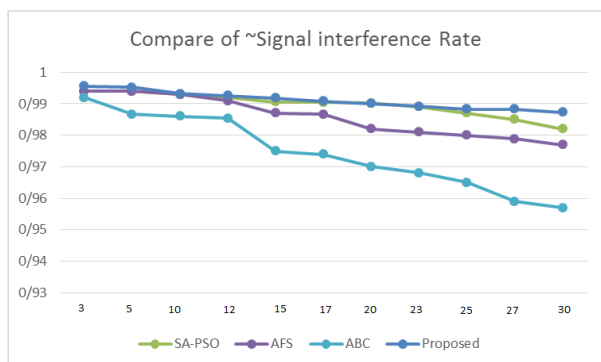


Figure 9: Comparison of the rate of signal interference algorithm to some readers

In the latest assessment, the rate of compliance with the proposed load balancing algorithm was tested in comparison with the three other algorithms. The third objective was to assess the use of the evaluation function. The figure (below) shows the results.

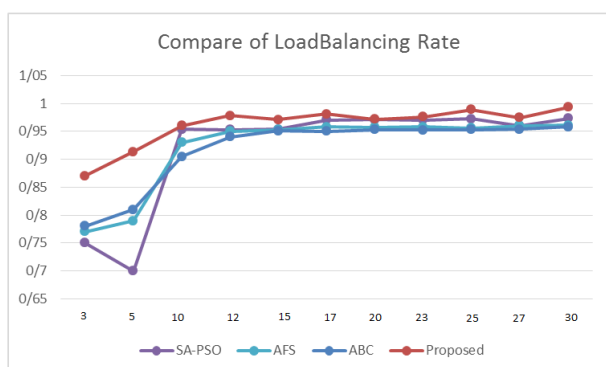


Figure 10: Comparison of the rate of load balancing algorithms to some readers

In figure (above), it is observed that when the number is less than 10 tag reader load balancing algorithms are weak in, in

this range SA-PSO algorithm is the weakest mode. And the proposed algorithm is best achieved. ABC algorithm and AFS have same performance. But when the number of tag readers are more than 10 the Performance and smoothing algorithms are on a constant scope and are stable. Compliance rate balance of the proposed between algorithms is better than others, and between 97% and 99% each. Here, too, even at low doses cause of better proposed performance of settled algorithm is, using of multiple objective functions and power PSO algorithm at the same time.

## 16. Conclusion

In this section, the simulation results of the proposed algorithm to determine the optimal positioning and deployment of tag readers based on MOPSO algorithm in, the Internet of Things based on RFID network was presented. Then according to the diagram and simulation results are presented in tables, proved that the proposed algorithm has better performance in all aspects, especially when the number of readers is less than the proposed algorithm is better able to list and achieved good results. It was observed that excessive reader effecting improvement in three indicators and the extra cost is not raised. In addition excessive reader in a network environment with fixed dimensions increases more signals are overlapping areas. In total, according to the results, the efficiency of the proposed algorithm and method of applying meta-heuristic algorithm And the definition of beneficial triple objective functions and using Pareto multi-objective optimization theory concepts to determine the suitability solution or the particles proved Also revealed that the optimum MOPSO algorithm, has good feature set optimized for deployment reader prone position and has a network with more speed and long life.

## 17. Research findings

In this paper, a new algorithm for optimal deployment of tag readers for RFID networks, in order to extend the life and network speed are presented. The algorithm is actually a mapping which is based on three parameters: maximum coverage, reduce signal interference and increase load balancing. With the use of Particle Swarm (PSO) and multi-objective optimization Pareto's theory, that is formed by the combination is referred to as MOPSO, has tried to reduce overall cost, increase speed and increase the lifetime of the network, due to multi-objective fitness function. in PSO algorithm design frame talent place as the solution and particulate components were used as pair-pair And using the concepts of Pareto theory and multi-objective optimization from three angles and with creation a balance between them was achieved.

## 18. Innovation research

In this paper, try to pay in analysis of the RFID network information in Internet of Things. New aspects of the project is using analytical solutions in information networks of the Internet of Things in the event that That in the vast geographical expanse of massive amounts of data and data feeds used in data networks and also reduce the lifetime of the



network. This paper aims to increase lifespan and maximum RFID tag detection in the world of the Internet of Things Taking advantage of a meta-heuristic methods to determine the position of the reader. The most important innovation in this article is using meta-heuristic algorithm in combination with multi-objective optimization PSO as defined MOPSO and efficient functions for these three purposes, which greatly increase the network coverage, reduce signal interference, increase the rate of load balancing and the thereby reducing the total cost of network speed and network lifetime and good connections.

## 19. Suggestions

As future work that can be done to develop a method presented in this paper, the following suggestions are presented:

- Use the proposed algorithm in wireless ad-hoc networks
- Use MOPSO algorithm considering deleted tag readers.
- Use Multi Objective Genetic Algorithm (MOGA) in the present case and compared with the method proposed in this paper.
- Applying another useful parameters for determining the fitness function.
- Use the proposed approach to integrated WSN and RFID network.
- Use other methods and hybrid approaches to determine the optimal number and position of their own tag locations (for example, using neural networks)

## 20. REFERENCES

- [1] Clauberg, R., RFID and Sensor Networks, RFID Workshop, University of St. Gallen, Switzerland, Sept. 27, 2013.
- [2] Fine, C., Klym, N., Trossen, D., and Tavshikar, M., The Evolution of RFID Networks: The Potential for Disruptive Innovation, MIT Center for eBusiness, CeB Research Brief, Vol. VIII, No. 1, 2011.
- [3] Holmqvist, M., and Stefansson, G., Mobile RFID - A Case from Volvo on Innovation in SCM, System Sciences, HICSS '06. Proceedings of the 39th Annual Hawaii International Conference on, vol 6, IEEE, 2006.
- [4] Ming Tao, ShuqiangHuang, YangLi, MinYan, YuyuZhou " SA-PSO based optimizing reader deployment in large-scale RFID Systems " journal homepage: [www.elsevier.com/locate/jnca](http://www.elsevier.com/locate/jnca).
- [5] Keyvan Kashkouli Nejad, Xiaohong Jiang, Michitaka Kameyama. "High Performance Tag Singulation for Memory-less RFID Systems" This full text paper was peer reviewed at the direction of IEEE Communications Society subject matter experts for publication in the IEEE ICC 2011 proceedings.

[6] Lianbo Maa, b, Kunyuan Hua, Yunlong Zhua, Hanning Chena, "Cooperative artificial bee colony algorithm for multi-objective RFID network planning" Volume 42, June 2014, Pages 143–162.

# A Comprehensive Survey on Privacy Preserving Big Data Mining

S.Srijayanthi

Department of Computer Science and Engineering  
R.M.K Engineering College, India

R.Sethukkarasi

Department of Computer Science and Engineering  
R.M.K Engineering College, India

**Abstract:** In recent years, privacy preservation of large scale datasets in big data applications such as physical, biological and biomedical sciences is becoming one of the major concerned issues for mining useful information from sensitive data. Preservation of privacy in data mining has ascended as an absolute prerequisite for exchanging confidential information in terms of data analysis, validation, and publishing. Privacy-Preserving Data Mining (PPDM) aids to mine information and reveals patterns from large dataset protecting private and sensitive data from being exposed. With the advent of varied technologies in data collection, storage and processing, numerous privacy preservation techniques have been developed. In this paper, we provide a review of the state-of-the-art methods for privacy preservation

**Keywords:** big data; confidential; data mining; privacy preservation; sensitive

## 1. INTRODUCTION

Privacy preservation in data mining has emerged as an unconditional prerequisite for exchanging private information in data analytics as internet phishing posed an intense menace on propagation of sensitive information over the web. Despite thriving of Big data provides potential values in healthcare, business analytics, government surveillance, and so on, substantial caution is essential in balancing the data utility and privacy preservation in the big data collection, storage and processing. Failing to protect privacy is immoral as it causes monetary loss and stern reputation impairment. Most methods for privacy computations use some form of data transformation to provide privacy preservation which reduce the granularity of representation resulting in some loss of effectiveness of data management or mining algorithms. This is the natural trade-off between information loss and privacy.

A data set is viewed as a file with  $n$  records, where each record contains  $m$  attributes. The attributes can be classified as [16] identifiers, quasi-identifiers, confidential outcome attributes and non-confidential outcome attributes. There are several approaches implemented for privacy preserving data mining. They are classified based on the following dimensions [44]:

- (i) Data modification
- (ii) Data or rule hiding
- (iii) Privacy preservation
- (iv) Data mining algorithm
- (v) Data distribution
- (vi)

## 2. LITERATURE SURVEY

### 2.1 Data Modification

Data modification techniques modify the original values of a database and the transformed database is made available for mining.

The basic idea of value-based perturbation approach is to add random noise to the data values. The technique [3] proposed is based on random noise addition and is as follows: Consider  $n$  original data  $A_1, A_2, \dots, A_n$  of one-dimensional distribution following the same independent and identical distribution (i.i.d). The  $n$  random variables  $B_1, B_2, \dots, B_n$  are generated to hide  $X_i$  data values. Distributed data is generated as  $W_1, W_2, \dots, W_n$  where  $W_i = A_i + B_i$ . According to the perturbed dataset  $W_1, W_2, \dots, W_n$  and a reconstruction procedure based on Bayes rule, the density function will be estimated by Eq. (1)

$$f_{X(a)} = \frac{1}{n} \sum_{i=1}^n \frac{\int_{-\infty}^a f_Y(w_i - a) f_X(a) \, da}{\int_{-\infty}^{+\infty} f_Y(w_i - z) f_X(z) \, dz} \quad (1)$$

The reconstruction procedure is improved by Expectation Maximization (EM) algorithm. This method is able to retain privacy while accessing the information implicit in the original attributes. It is more effective in terms of information loss. The authors proved that the EM algorithm converges to the maximum likelihood estimate of the original distribution based on the perturbed data.

A randomized Response (RR) technique was developed in the statistics community for the purpose of protecting surveyee's privacy and was first introduced by Warner [53]. Two models were proposed to solve the survey problem: Related-Question Model in which each respondent is asked two related questions, the answers to which are opposite to each other and Unrelated-Question Model in which two unrelated questions are asked with one probability for one of

the questions is known. The Multivariate Randomized Response (MRR) technique [18] was proposed for multiple-attribute data set. The method consists of two parts: the first part is the multivariate data disguising technique used for data collection; the second part is the modified ID3 decision tree building algorithm used for building a classifier from the disguised data. The framework [9] conducts a multivariate regression analysis to generate predicted probabilities for the sensitive item. They showed to use the sensitive attitude inferred from the multivariate regression analysis as a predictor for an outcome regression model.

The condensation based technique [4] was proposed to generate pseudo-data from clustered groups of k-records. Principal component analysis of the behaviour of the records within a group is used in the generation of pseudo-data. The use of pseudo-data provides an additional layer of protection. Also, the aggregate behaviour of the data is preserved thus useful for a variety of data mining problems. The technique [5] constructed groups of non-homogeneous size from the data, such that it is guaranteed that each record lies in a group whose size is at least equal to its anonymity level. Then, pseudo-data were generated from each group to create a synthetic data set with the same aggregate distribution as the original data. Aggarwal [1] has proposed a method for anonymization of string data that creates clusters from the different strings, and then generates synthetic data which has the same aggregate properties as the individual clusters. Since each cluster contains at least k-records, the anonymized data is guaranteed to at least satisfy the definitions of k-anonymity.

The main idea of random rotation perturbation technique is that the original dataset with d columns and N records represented as  $X_{d \times N}$ . The rotation perturbation of the dataset X will be defined as  $g(X) = RX$ . Where  $R_{d \times d}$  is a random rotation orthonormal matrix. A key feature of rotation transformation is preserving the Euclidean distance, inner product in a multi-dimensional space. The optimal algorithm [11] perturbs all columns together. The authors defined an efficient multi-column privacy measure for evaluating the privacy quality of any rotation perturbation. The level of difficulty for the estimation of the original data is by variance of the difference. Let  $r_{ij}$  represent the *element*(i, j) in the matrix R, and  $c_{ij}$  be the *element*(i, j) in the covariance matrix of X. The Variance of Difference ith column is computed by Eq. (2)

$$\text{Cov} \left( \begin{matrix} X \\ X \end{matrix} \right)_{(i,i)} = \sum_{j,k} r_{ij} r_{ik} c_{kj} - 2 \sum_{j} r_{ij} c_{ij} + c_{ii} \quad (2)$$

Geometric

data perturbation consists of a sequence of random geometric transformations, including multiplicative transformation (R), translation transformation ( $\Psi$ ), and distance perturbation ( $\Delta$ ). Authors Chen and Liu [12] have developed three protocols: (i) simple protocol to transmit

encrypted perturbed data to the service provider. (ii) negotiation protocol enables multi-round voting to reach an agreed perturbation. (iii) The space adaptation protocol provides a better balance between scalability, flexibility of data distribution, and the overall satisfaction level of privacy guarantee. The authors Chen and Liu [13] proposed a multi-column privacy evaluation model and designed a unified privacy metric to address the problems. The authors analysed the resilience of the rotation perturbation approach against three types of inference attacks: naive-inference attacks, ICA-based attacks, and distance-inference attacks. The authors constructed a randomized optimization algorithm to efficiently find a good geometric perturbation that is resilient to the attacks.

Random projection projects a set of data points from a high dimensional space to a randomly chosen lower dimensional subspace. The basic idea of random projection arises from the Johnson-Lindenstrauss Lemma. The authors Kargupta et al. [23] proposed spectral filtering technique that can estimate values of individual data-points from the perturbed dataset and thus can be used to reconstruct the distribution of actual data as well. Signal-to-Noise Ratio (SNR) quantifies the relative amount of noise added to actual data to perturb it and is given by Eq. (3)

$$SNR = \frac{\text{Variance of Actual Data}}{\text{Noise Variance}} \quad (3)$$

The authors Liu et al. [32] showed that the projection can preserve the inner product, which is directly related to several distance-related metrics, by conducting row wise and column-wise projection of the sample data. The authors Li et al. [29] expanded scope of additive perturbation based PPDM to multi-level trust (MLT). The method allows data owners to generate differently perturbed copies of its data for different trust levels on demand, offering maximum flexibility to data owners. The key challenge lies in circumventing from combining copies at different trust levels to jointly reconstruct the original data. This is addressed by properly correlating perturbation across copies at different trust levels.

The Singular value decomposition SVD technique is used to distort portions of the datasets. The SVD of the data matrix A is given by the Eq. (4)

$$A = U \Sigma V^T \quad (4)$$

where A be a sparse matrix of dimension  $n \times m$  representing the original dataset. U is  $n \times n$  orthogonal matrix and  $V^T$  is  $m \times m$  orthogonal matrix. A transformed matrix with a much lower dimension is defined by Eq. (5)

$$A_k = \bigcup_k \sum_k V_k T \quad (5)$$

The proposed data distortion method [48], sparsified SVD, is better than SVD. Entries smaller than a certain threshold in are set to zero after reducing the rank of the SVD matrices. This operation is called as dropping operation. The distorted data matrix  $\bar{A}_k$  is written as Eq. (6)

$$\bar{A}_k = \bar{U}_k \Sigma_k \bar{V}_k^T \quad (6)$$

$\bar{A}_k$  is twice distorted in the sparsified SVD method and thus it is harder to reconstruct the entries in A. The computation of SVD for large scale dataset matrices is expensive which can be substantially reduced by employing clustered SVD strategies.

In the proposed algorithm [21], attributes are grouped according to their distance difference similarity by clustering the data set using decision tree classification. The algorithm packetizes the attributes in each group and for each group it creates an equivalence class following the unique attribute-distinct diversity anonymization model. The weights given to attributes improve clustering and give the ability to control the generalization's depth.

In Non-negative matrix factorization (NNMF) technique is a vector space method to obtain a representation of data using non-negative constraints. Considering  $n \times m$  nonnegative matrix dataset A with  $A_{ij} \geq 0$  and a pre-specified positive integer  $k \leq \min\{n, m\}$ , nonnegative matrix factorization finds two non-negative matrixes  $W \in R^{n \times k}$  with  $W_{ij} \geq 0$  and  $H \in R^{k \times m}$  with  $H_{ij} \geq 0$ , such that  $A \approx WH$  and the objective function given by Eq. (7) is minimized

$$f(W, H) = \frac{1}{2} \|A - WH\|_F^2 \quad (7)$$

where  $\|A - WH\|_F$  the Frobenius norm. Matrices W and H have desirable properties in data mining applications. The work in [51] contributed least-square compression form of original datasets and iterative methods to solve the least-square optimization problem.

Each release of the data must be such that every combination of values of quasi-identifiers can be indistinguishably matched to at least k respondents. In k-anonymity techniques [39], the granularity of representation of the pseudo-identifiers is minimized by generalization and suppression. k-Optimize algorithm [8] assumes ordering among the quasi-identifiers. The Incogito method [25] has been proposed for computing k-minimal generalization with the use of bottom-up aggregation along domain generalization hierarchies. [19] starts with a general solution, and then specializes some attributes of the current

solution so as to increase the information, but reducing the anonymity. The reduction in anonymity is always controlled, so that k-anonymity is never violated.

A hybrid approach [49] proposed combined Top Down Specialization (TDS) and BUG (Bottom Up Generalization) together for efficient sub-tree anonymization over big data. The approach automatically determined which component to be used to conduct the anonymization when a data set was given, by comparing the user specified k-anonymity parameter with a threshold derived from the dataset. Both components TDS and BUG are developed based on Map Reduce (MR) to gain high scalability by exploiting powerful computation capability of cloud.

In TDS [50] scalable approach proposed, a data set is anonymized by performing specialization operations. In the first phase, data sets are partitioned and anonymized in parallel, producing intermediate results. The intermediate results are merged and further anonymized to produce k-anonymous data sets in the second phase. The goodness of a candidate specialization is measured by a search metric, Information Gain per Privacy Loss (IGPL).

BUG approach of anonymization is an iterative process starting from the lowest anonymization level. The goodness of a candidate generalization is measured by a search metric, Information Loss per Privacy Gain (ILPG).

ILPG of generalization is given by the Eq. (8)

$$ILPG(gen) = IL(gen) / (PG(gen) + 1) \quad (8)$$

Information loss is given by Eq. (9)

$$IL(gen) = \sum_{c \in Child(q)} \left( \frac{R_c}{Rq} \right) I(R_c) - I(Rq) \quad (9)$$

The privacy gain is given by Eq. (10)

$$PG(gen) = A_{p(gen)} - A_{c(gen)} \quad (10)$$

The extended k-Anonymity [47], ( $\alpha$ , k)-Anonymity, combines two principles: (i) each equivalence class must have size at least k (ii) at most  $\alpha$  percent of its tuples can have the same sensitive value. The authors presented an optimal global-recoding ( $\alpha$ , k)-anonymization algorithm and a scalable local-recoding technique that shows less data distortion.

The k-Anonymity technique is vulnerable to many kinds of attacks if the background knowledge is known. Such kinds of attacks include are homogeneity attack and background knowledge attack. The l-diversity technique proposed not only maintains the minimum group size of k, but also focuses on maintaining the diversity of the sensitive attributes. Therefore, the l-diversity model [24] for privacy is defined as a group of indistinguishable tuples are l-diverse

if they contain at least  $l$  “well-represented” values for the sensitive attributes. Liu and Wang [31] proposed an extension of  $l$ -diversity using full-subtree generalization and suppression techniques. It is stated that the confidence of the adversary in inferring a target’s sensitive information is

bounded by the percentage  $conf(S_i | QI_j)$  of the records that contain the same value  $S_i$  in the equivalence class  $j$ . Authors limit this bound by guaranteeing  $conf(S_i | QI_j) \leq \theta_i$ , where parameter  $\theta_i$  is a given privacy threshold in the interval  $[0, 1]$ . A dynamically created structure, Cut enumeration tree, enumerates all possible generalizations of  $QI$  attributes according to the generalization level and information loss of each candidate solution.

To prevent skewness attack, the authors [26] proposed a privacy model, called  $t$ -closeness, which states that an equivalence class is said to have  $t$ -closeness if the distance between the distribution of a sensitive attribute in this class and the distribution of the attribute in the whole table is no more than a threshold  $t$ . It severely affects the data utility as it needs the distribution of sensitive values to be the same in all equivalence classes. The authors used Earth Mover Distance (EMD) to measure the closeness between two sensitive values which does not prevent attribute linkage on numerical sensitive attributes. Table 1. shows the comparison of data perturbation techniques.

Table 1. Comparison of perturbation techniques

Criteria	perturbation		
	Value Based	Data Mining Task	Dimensional Reduction
Privacy Loss	Average	Low	Very Low
Information Loss	Low	Very Low	Very Low
Modify Mining Algorithms	Yes	No	No
Data Dimension	Single	Multi	Multi

## 2.2 Association Rule Hiding

Association rule hiding algorithms can be divided into three classes namely: (i) Heuristic (ii) Border-based approaches (iii) Exact approaches

The heuristic approaches sanitize a set of transactions from the database to hide the sensitive knowledge. It is efficient and scalable. In several circumstances they suffer from undesirable side-effects that lead them to suboptimal solutions. The authors Atallah et al. [7] proposed a greedy iterative search algorithm to hide sensitive association rules through the reduction in the support of their generating

itemsets. The limitation is the loss of support for a large itemset, as long as it remains frequent in the sanitized outcome. Verykios et al. [45] proposed two heuristic algorithms. The first algorithm hides the item having the maximum support from the minimum length transaction. The second algorithm sorts the generating itemsets with respect to their size and support. The algorithm removes the items from the corresponding transactions in a round-robin fashion, until the support of the sensitive itemset drops below the minimum support threshold. Amiri [6] proposed three effective, multiple rule hiding heuristics approach: (i) Aggregate approach (ii) Disaggregate approach (iii) Hybrid approach.

DSRRC (Decrease Support of Right hand side item of Rule Clusters) algorithm [34] clusters the sensitive rules based on certain criteria in order to hide as many as possible rules at one time. One shortcoming of this algorithm is that it cannot hide association rules with multiple items in antecedent and consequent. The authors Domadiya and Rao [15] introduced a heuristic based algorithm called Modified Decrease Support of RHS item of Rule Clusters (MDSRRC) to secure the delicate association rules using multiple items in consequent (RHS) and antecedent (LHS). This algorithm successfully addressed the drawbacks of rule hiding DSRRC algorithm.

Saygin et al [40] proposed the usage of unknowns instead of altering 1’s to 0’s and vice versa to hide association rules. The two schemes [46] proposed include unknowns and aimed at the hiding of predictive association rules. The algorithms proposed require a reduced number of database scans and exhibit an efficient pruning strategy. The first scheme decreases the confidence of a rule by increasing the support of the itemset in its LHS. The second approach reduces the confidence of the rule by decreasing the support of the itemset in its RHS.

The border based approach modifies the original borders in the lattice of the frequent and the infrequent patterns in the dataset. The sensitive knowledge is hidden by enforcing the revised borders in the sanitized database.

The Sun and Yu [41] proposed a scheme which first computes the positive and the negative borders in the lattice of all itemsets. A weight is assigned to each element of the expected positive border which is dynamically computed as a function of the current support. The algorithm deletes the candidate item that will have the minimal impact on the positive border. The authors Moustakides and Verykios [35] proposed an algorithm to remove all the sensitive itemsets belonging to the revised negative border. Among all minimum border itemsets, the one with the highest support is selected. This max-min itemset determines the item through which the hiding of the sensitive itemset will incur.

The exact approach considers the hiding process as a constraint satisfaction problem solved using integer or linear programming. Exact approaches are efficient than heuristic schemes, at a high computational cost. They formulate the sanitization process as a constraint satisfaction problem. The scheme [33] consists of an exact and a heuristic part in which the exact part formulates a Constraint Satisfaction Problem (CSP) with the objective of identifying the minimum number of transactions that need to be sanitized. An integer programming solver is then used to identify the best solution. An approach [20] uses the itemsets belonging in the revised positive and negative borders to identify the candidate itemsets for sanitization. It obtains efficient solution of the CSP, by using binary integer programming.

### 2.3 Privacy Preservation

It refers to the privacy preservation technique used for the selective modification of the data. The cryptographic methods tend to compute functions over inputs provided by multiple recipients without actually sharing the inputs with one another. The challenge is to conduct such a computation while preserving the privacy of the inputs. [14] presented four secure multiparty computation based methods that can support privacy preserving data mining. The methods described include, the secure sum, the secure set union, the secure size of set intersection, and the scalar product. The authors Kantarcioglu and Clifton [22] addressed the problem of secure mining of association rules over horizontally partitioned data based on the assumption that each party first encrypts its own itemsets using commutative encryption, then the ready encrypted itemsets of every other party. A secure comparison takes place between the final and initiating parties to determine if the final result is greater than the threshold plus the random value. Based on cryptographic techniques Chakravorty et al. [10], replaced the personal/quasi- identifiers of collected sensor data with hashed values before storing them into a de-identified storage. In the Dong and Chen [17] proposed an efficient secure dot product protocol based on the Goldwasser–Micali Encryption and Oblivious Bloom Intersection for privacy preserving association rule mining. The protocol is faster as it employs mostly cheap cryptographic operations, e.g. hashing and modular multiplication. The authors Wang et al. [52] proposed a privacy-preserving public auditing system for data storage security in cloud computing utilizing the homomorphic linear authenticator and random masking to guarantee that the third party would not learn any information stored on the cloud server during the auditing process.

The reconstruction based methods first randomize the original data to hide the sensitive data and then reconstruct the interesting patterns based on the statistical features without knowing true values. The work [3] addresses the problem of building a decision tree classifier from training data in which the values of individual records have been

perturbed. By using the reconstructed distributions, they are able to build classifiers whose accuracy is comparable to the accuracy of classifiers built with the original data. The approach [2] is based on an Expectation Maximization (EM) algorithm for distribution reconstruction which converges to the maximum likelihood estimate of the original distribution on the perturbed data.

### 2.4 Data Mining Algorithm

To extract useful information from big data without breaching the privacy, privacy preserving data mining techniques have been developed to identify patterns and trends from data. These techniques can be broadly grouped into clustering, classification and association rule based techniques.

The authors Zhou et al. [54] proposed a parallel k-means clustering algorithm by using three functions of MapReduce. First, the Map tasks calculate the closest distance for data points from every initial centroid of clusters. Next, the combiner calculates a partial sum of values. The Reduce tasks compute the centroids by dividing the partial sum of samples in to the number of samples assigned to a similar cluster. The Mapper processes each data object and called several times which increases the problem in handling large data sets.

In Incremental k-means Algorithm (IKM) [36], the Mapper loads data segment, and executes the IKM on the loaded data segment. The Reducer receives the intermediate results and executes the IKM again to obtain the clustering results. This approach provides an approximate solution and does not provide exact clustering results. Li et al. [27] focused on concurrently running k-means processes based on MapReduce with multiple initial center groups. Its main objective is to avoid serial execution of k-means and more focus on initial centroids. In this approach, the hopeless k-means process attempts are abandoned, which speeds up the future iterations. However, because of using MapReduce, it still lacks the ability to cache data between iterations for improving performance.

Classification is a technique of identifying, to which predefined group a new input data belongs. Classification algorithm is designed to process data in two ways [42]. It either classifies the data by themselves or forward the input data to another classifier. It is computationally efficient particularly when handling large and complex data.

The algorithm [3] in which the original data are altered by adding random offsets was proposed. Bayesian formula is used to derive the density function of the original data. A random forest is built based on Mahout RF Partial implementation to classify imbalanced big data. The algorithm calculates the leaves weights for each tree. Then, the leaf weight is the accumulated weight divided by the

number of instances classified and then the algorithm combines the outputs from each mapper. For each instance in all classes, the accumulated weight is divided by the number of trees involved in the classification.

A global SVM classification model [43] was constructed based on gram matrix computation to securely compute the kernel matrix from the distributed data. The algorithm, [30] Privacy-Preserving SVM Classifier PPSVC approximates the decision function of the Gaussian kernel SVM classifier without compromising the sensitive attribute values possessed by support vectors. The PPSVC is robust against adversarial attacks and the accuracy is comparable to the original SVM classifier. Quantum based support vector machine [38] for big data classification minimizing the computational complexity and the required training data was proposed.

## 2.5 Distributed Privacy Preservation

The key goal in most distributed methods for privacy-preserving data mining is to allow computation of useful aggregate statistics over the entire data set without compromising the privacy.

The BOPPID (Boosting – based Privacy Preserving Integration of Distributed data) algorithm [28] in which each participant has different set of records with both common features and local unique attributes. AdaBoost algorithm was employed to build an ensemble classifier. By sharing the local models with each other, all the participators can build their individual integrated model without direct access to the datasets. To prevent “negative impact” during integration, the models from the other participators whose data distribution is very different from the data distribution of this participator are excluded. The proposed method overcomes the need of third-party and reduces the communication cost. An algorithm [37] was proposed for differentially private data release for vertically partitioned data. The two-party differentially private data release algorithm anonymized the raw data by sequence of specialization and added noise. The proposed distributed exponential mechanism takes candidate and score pairs as inputs. Candidates are selected based on their score functions. The score is determined using Max utility function given by Eq. 11

$$Max(D, v) = \sum_{c \in child(v)} \left( \max |D_c^{cls}| \right) \quad (11)$$

## 3. CONCLUSION

The privacy preservation for data analysis is a challenging research issue due to increasingly larger volumes of data sets, thereby requiring intensive investigation. Each privacy preserving technique has its own importance. Data encryption and anonymization are widely adopted ways to combat privacy breach. However, encryption is not suitable for data that are

processed and shared. Anonymizing big data and managing anonymized data sets are still challenges for traditional anonymization approaches. Privacy-preserving data mining is emerged for to two vital needs: data analysis in order to deliver better services and ensuring the privacy rights of the data owners. Substantial efforts have been accomplished to address these needs. In this paper, an overview of the recent approaches for privacy preservation was presented. The privacy guarantees, advantages and disadvantages and possible enhancement of each approach were stated.

## 4. REFERENCES

- [1] C. Aggarwal, “On randomization, public information and the curse of dimensionality”, In: IEEE 23rd International Conf. on Data Engineering, 2007.
- [2] D. Agrawal, and C. Aggarwal. “On the design and quantification of privacy preserving data mining algorithms”, In: Proc. of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems, 2001.
- [3] R. Agrawal, and R. Srikant, “Privacy-preserving data mining”, ACM Sigmod Record. Vol. 29, No. 2, 2000.
- [4] C. Aggarwal and S. Yu Philip. “A condensation approach to privacy preserving data mining”, In: International Conf. on Extending Database Technology. Springer Berlin Heidelberg, 2004.
- [5] C. Aggarwal and P. S. Yu, “On Variable Constraints in Privacy Preserving Data Mining”, SDM. 2005.
- [6] A. Amiri, “Dare to share: Protecting sensitive knowledge with data sanitization”, Decision Support Systems 43(1), pp.181-191, 2007.
- [7] M. Atallah, E. Bertino, A. Elmagarmid, M. Ibrahim, and V. S. Verykios, “Disclosure limitation of sensitive rules” Knowledge and Data Engineering Exchange, 1999.(KDEX'99), Proc. 1999 Workshop on, pp. 45-52, IEEE, 1999.
- [8] R.J. Bayardo and R. Agrawal, “Data privacy through optimal k-anonymization” In: 21st International Conf. on Data Engineering (ICDE'05), pp.217-228, IEEE, 2005.
- [9] G. Blair, K. Imai, and YY. Zhou. “Design and Analysis of the Randomized Response Technique”, Journal of the American Statistical Association 110(511), pp. 1304-1319, 2015.
- [10] A. Chakravorty, T. Wlodarczyk and C. Rong, “Privacy Preserving Data Analytics for Smart Homes” In Security and Privacy Workshops (SPW), pp. 23-27. IEEE, 2013.
- [11] K. Chen and L. Liu, “Privacy preserving data classification with rotation perturbation”, In: Fifth IEEE International Conf. on Data Mining (ICDM'05), pp. 4-pp. IEEE, 2005.
- [12] K. Chen and L. Liu, “Privacy-preserving multiparty collaborative mining with geometric data perturbation”, IEEE Transactions on Parallel and Distributed Systems 20(12), pp. 1764-1776, 2009.
- [13] K. Chen, and L. Liu, “Geometric data perturbation for privacy preserving outsourced data mining”, Knowledge and Information Systems 29(3), pp. 657-695, 2011.
- [14] C. Clifton, M. Kantarcioglu, J. Vaidya, X. Lin, and MY. Zhu, “Tools for privacy preserving distributed data mining”, ACM Sigkdd Explorations Newsletter 4, no. 2, pp.28-34, 2002.
- [15] N. H. Domadiya, U. P. Rao, “Hiding sensitive association rules to maintain privacy and data quality in database”, In: Advance Computing Conf. (IACC), 2013 IEEE 3rd International, pp.1306-1310, 2013.

- [16] J. Domingo-Ferrer, “A survey of inference control methods for privacy-preserving data mining”, In: Privacy-preserving data mining, pp. 53-80. Springer US, 2008.
- [17] C. Dong and L. Chen, “A fast secure dot product protocol with application to privacy preserving association rule mining”, In: Pacific-Asia Conference on Knowledge Discovery and Data Mining, pp. 606-617, Springer International Publishing, 2014.
- [18] W. Du and Z. Zhan, “Using randomized response techniques for privacy-preserving data mining”, In: Proc. of the ninth ACM SIGKDD international conf. on Knowledge discovery and data mining, pp. 505-510, ACM, 2003.
- [19] BCM. Fung, K. Wang, and PS. Yu. “Top-down specialization for information and privacy preservation”, In: 21st International Conf. on Data Engineering (ICDE'05), pp.205-216, IEEE, 2005.
- [20] A. Gkoulalas-Divanis and V S. Verykios. “An integer programming approach for frequent itemset hiding”, In: Proc. of the 15th ACM international conf. on Information and knowledge management, pp. 748-757. ACM, 2006.
- [21] P. Jain, N. Pathak, P. Tapashetti, and A. S. Umesh. “Privacy preserving processing of data decision tree based on sample selection and Singular Value Decomposition”, In: 9th International Conf. on Information Assurance and Security (IAS), 2013, pp. 91-95, IEEE, 2013.
- [22] M. Kantarcioglu, and C. Clifton, “Privacy-preserving distributed mining of association rules on horizontally partitioned data”, IEEE transactions on knowledge and data engineering 16, no. 9, pp. 1026-1037, 2004.
- [23] H. Kargupta, Hillol, S. Datta, Q. Wang, and K. Sivakumar. “On the privacy preserving properties of random data perturbation techniques.” In: Third IEEE International Conference on Data Mining, 2003. ICDM 2003., pp. 99-106, IEEE, 2003.
- [24] A. Machanavajjhala, D. Kifer, J. Gehrke and M. Venkatasubramanian, “l-diversity: Privacy beyond k-anonymity”, ACM Transactions on Knowledge Discovery from Data (TKDD), 1(1), 3, 2007.
- [25] K. LeFevre, DJ. DeWitt, and R. Ramakrishnan. “Incognito: Efficient full-domain k-anonymity”, In: Proc. of the 2005 ACM SIGMOD international conf. on Management of data, pp. 49-60, ACM, 2005.
- [26] N. Li, T. Li, and S. Venkatasubramanian, “t-closeness: Privacy beyond k-anonymity and l-diversity”, In: IEEE 23rd International Conf. on Data Engineering, pp. 106-115, IEEE, 2007.
- [27] C. Li, Y. Zhang, M. Jiao, and Ge Yu. “Mux-Kmeans: multiplex kmeans for clustering large-scale data set”, In: Proc. of the 5th ACM workshop on Scientific cloud computing, pp. 25-32, ACM, 2014.
- [28] Y. Li, C. Bai, and CK. Reddy. “A distributed ensemble approach for mining healthcare data under privacy constraints”, Information sciences 330, pp. 245-259, 2016.
- [29] Y. Li, M. Chen, Q. Li, and W. Zhang. “Enabling multilevel trust in privacy preserving data mining”, IEEE Transactions on Knowledge and Data Engineering 24, no. 9, pp. 1598-1612, 2012
- [30] KP. Lin, and MS. Chen. “On the design and analysis of the privacy-preserving SVM classifier”, IEEE Transactions on Knowledge and Data Engineering 23, no. 11 pp. 1704-1717, 2011.
- [31] Liu, Junqiang, and Ke Wang. “On optimal anonymization for l+-diversity”, In: IEEE 26th International Conf. on Data Engineering (ICDE 2010), IEEE, 2010.
- [32] K. Liu, H. Kargupta, and J. Ryan. “Random projection-based multiplicative data perturbation for privacy preserving distributed data mining” IEEE Transactions on knowledge and Data Engineering 18(1), pp. 92-106, 2006.
- [33] S. Menon, S. Sarkar, and S. Mukherjee. “Maximizing accuracy of shared databases when concealing sensitive patterns”, Information Systems Research 16(3), pp. 256-270, 2005
- [34] Modi, Chirag N., Udai Pratap Rao, and Dhiren R. Patel. “Maintaining privacy and data quality in privacy preserving association rule mining”, In: International Conf. on Computing Communication and Networking Technologies (ICCCNT), 2010.
- [35] GV. Moustakides and VS. Verykios, “A MaxMin approach for hiding frequent itemsets”, Data & Knowledge Engineering 65(1), pp. 75-89, 2008.
- [36] DT. Pham, SS. Dimov, and CD. Nguyen, “An incremental K-means algorithm”, In: Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science 218(7), pp. 783-795, 2004.
- [37] N. Mohammed, D. Alhadidi, BCM Fung, and M. Debbabi, “Secure two-party differentially private data release for vertically partitioned data”, IEEE Transactions on Dependable and Secure Computing 11, no. 1, pp. 59-71, 2014.
- [38] P. Rebertrost, M. Mohseni, and S. Lloyd, “Quantum support vector machine for big feature and big data classification”, arXiv preprint arXiv:1307.0471, 2013.
- [39] P. Samarati, “Protecting respondents identities in microdata release”, IEEE transactions on Knowledge and Data Engineering, 13(6), pp.1010-1027, 2001.
- [40] Y. Sayginl, VS. Verykios, and C. Clifton. “Using unknowns to prevent discovery of association rules”, Acm Sigmod Record 30, no. 4, pp. 45-54, 2001.
- [41] X. Sun, and PS. Yu. “A border-based approach for hiding sensitive frequent itemsets”, In: Fifth IEEE International Conf. on Data Mining (ICDM'05), pp. 8-pp. IEEE, 2005.
- [42] C. Tekin and M. van der Schaar. “Distributed online big data classification using context information”, In: 51st Annual Allerton Conference on Communication, Control, and Computing (Allerton), 2013, pp. 1435-1442, IEEE, 2013.
- [43] J. Vaidya, H. Yu, and X. Jiang. “Privacy-preserving SVM classification”, Knowledge and Information Systems 14, no.2, pp. 161-178, 2008.
- [44] VS. Verykios, E. Bertino, IN. Fovino, LP. Provenza, Y. Saygin, and Y. Theodoridis, “State-of-the-art in privacy preserving data mining”, ACM Sigmod Record 33, no. 1, pp. 50-57, 2004.
- [45] VS. Verykios, AK. Elmagarmid, E. Bertino, Y. Saygin, and E. Dasseni. “Association rule hiding”, IEEE Transactions on knowledge and data engineering 16, no. 4, pp. 434-447, 2004.
- [46] SL. Wang and A. Jafari, “Using unknowns for hiding sensitive predictive association rules”, In: IRI-2005 IEEE International Conf. on Information Reuse and Integration, Conference, 2005, pp. 223-228. IEEE, 2005.
- [47] RCW. Wong, J. Li, AWC. Fu, and K. Wang. “( $\alpha$ ,  $k$ )-anonymity: an enhanced k-anonymity model for privacy preserving data publishing”, In: Proc. of the 12th ACM SIGKDD international conf. on Knowledge discovery and data mining, pp. 754-759. ACM, 2006.



- [48]S. Xu,, J. Zhang, D. Han and J. Wang,. “Singular value decomposition based data distortion strategy for privacy protection”, Knowledge and Information Systems, 10(3). pp. 383-397, 2006.
- [49]X. Zhang, C. Liu, S. Nepal, C. Yang, W. Dou and J. Chen, “A hybrid approach for scalable sub-tree anonymization over big data using MapReduce on cloud”, Journal of Computer and System Sciences, 80(5). pp. 1008-1020, 2014.
- [50]X. Zhang, LT. Yang, C. Liu and J. Chen, “A scalable two-phase top-down specialization approach for data anonymization using mapreduce on cloud”. IEEE Transactions on Parallel and Distributed Systems 25(2), pp. 363-373, 2014.
- [51]J. Wang, W. Zhong, and J. Zhang. “NNMF-based factorization techniques for high-accuracy privacy protection on non-negative-valued datasets.” Sixth IEEE International Conf. on Data Mining-Workshops (ICDMW'06), 2006.
- [52]C. Wang, S. S. M. Chow, Q. Wang, K. Ren and W. Lou, “Privacy-preserving public auditing for secure cloud storage”, IEEE Transactions on computers 62, no. 2, pp. 362-375,2013.
- [53] SL. Warner, “Randomized response: A survey technique for eliminating evasive answer bias”, Journal of the American Statistical Association 60.309, pp. 63-69, 1965.
- [54] P. Zhou, J. Lei, and W. Ye. “Large-scale data sets clustering based on MapReduce and Hadoop”, Journal of Computational Information Systems 7.16, pp. 5956-5963, 2011.

# Resource Allocation in Cloud Environment Using Approaches Based Particle Swarm Optimization

Vahid Asadzadeh Chalack  
Department of Computer  
Azad Islamic University  
Germi, Iran

Seyed Naser Razavi  
Computer Engineering  
Department, Faculty of  
Electrical and Computer  
Engineering, University of  
Tabriz, Iran

Sajjad Jahanbakhsh Gudakahriz  
Computer Engineering  
Department of Computer  
Azad Islamic University  
Germi, Iran

---

**Abstract:** It is obvious that in emerging computing paradigms such as cloud computing systems, scheduling is one of the main phases to take advantages of capabilities. The cloud computing environment is a dynamic environment which allows services to be shared among many users. Scheduling methods of traditional systems are ill-suited for the cloud computing systems, and this new environment requires new methods tailored to its specifications. In this paper, we developed multiple algorithms for task scheduling in cloud computing systems. These algorithms are based on the particle swarm optimization (PSO) algorithm, which is a technique inspired by collective and social behavior of animal swarms in nature, and wherein particles search the problem space to find an optimal or near-optimal solution. The algorithms were developed with the aim of minimizing Makespan, Flowtime and the task execution cost simultaneously. Simulation and test results show the better efficiency of the proposed methods than other similar algorithms.

**Keywords:** Resources Allocation, Cloud computing, Particle swarm optimization, Makespan, Flowtime.

---

## 1. INTRODUCTION

Cloud computing is an internet-based computing paradigm that provides a new framework for provision, consumption and delivery of IT services (including software, information and shared computing resources). Cloud computing allows the IT resources to be provided through a flexible and scalable internet-based method at the moment and the scale that is demanded by the user. A cloud computing provider offers online commercial applications through a web browser or other software. Applications and information are stored in the servers and users are given access on request. Details stay hidden from the users and they need no expertise or knowledge on cloud infrastructure technology to use it. In a distributed system, some nodes may become heavily loaded while others remain inactive or underutilized. The task of distribution scheduler is to schedule processes to nodes (processors) in an optimal way. While there is an implicit distinction between task scheduling and task allocation, these two terms revolve around one problem. In terms of resources, the problem is how to allocate processors to the processes. From a user perspective, the problem is how to schedule the processes in the processors. Studies have shown that nature-inspired heuristic optimization methods are most effective for such applications. Most of these methods try to minimize the total time of execution. Collective Intelligence is an artificial intelligence method based on collective behavior, which has been used by many researchers to develop heuristic collective intelligence optimization algorithms, such as ant colony optimization (ACO) algorithm, particle swarm optimization (PSO) algorithm and firefly algorithm. Among the abovementioned algorithms, PSO is widely regarded as the best because of features such as rapid convergence, insensitivity to initiation values, flexibility and high tolerance to error. The major drawback of PSO however is in its local search, which often converges to local optima and lead to globally suboptimal solutions. It has been proven that this algorithm can be improved through combination with other methods.

## 2. LITERATURE

There are various types of task scheduling algorithm. The main goal of a scheduling algorithm is to achieve high computing performance and best system throughput. Traditional scheduling algorithms cannot operate in cloud environment (because of overhead costs), thus providers have resorted to heuristic or hybrid algorithms to fill this gap [1]. Effectiveness of task scheduling has a direct effect on the quality of cloud, thus many algorithms have been developed to resolve this particular problem [2]. In some studies, algorithms have been developed to optimize the resource efficiency. In this section, we review some of these scheduling algorithms. Efficient virtual machine allocation and task scheduling is one of the key issues of the cloud systems. ACO-based load balancing algorithm has been proposed for solving the load balancing issue of virtual machines in the task scheduling process. This algorithm can adapt to dynamic cloud environment and decrease task scheduling execution time while providing a balanced load for the virtual machines of data centers [3]. Dynamic load balancing strategy has been developed based on genetic algorithm. This algorithm increases scheduling speed, reduces switching between processors, selects the task to be executed and even identifies the details of the tasks when used in advanced state. Adaptability threshold of this algorithm contributes to the load balancing of dynamic processors [4].

The time-cost balancing algorithm has been developed to provide a number of functions based on cloud computing capabilities. The example of these functions is the compression of workflow aimed to reduce execution time and cost based on the user input information [5]. This algorithm tries to efficiently map tasks to resources in the cloud. The main phase of the algorithm includes the use of improved bee colony algorithm to assign priorities to tasks and then the use of an algorithm for task grouping based on their priorities. This scheduling algorithm computes the cost for acquisition of resources and the performance of calculations to complete workflow tasks. In this algorithm, the ratio of acquisition of

the resources to the cost of efficient communication to perform workflow tasks had a significant improvement. [9]

A heuristic workflow scheduler algorithm based on the group optimization of tasks is proposed. The heuristic workflow scheduling algorithm was designed based on the group optimization of tasks, heuristic approach of applications according to the resources in the cloud with the aim of reducing computation times and reducing data transfer times. The algorithm includes two main components: the use of heuristic algorithm in order to discover right resources and the use of group optimization algorithm for correct mapping of these resources to tasks. The experimental results show that the use of this algorithm will save costs and properly distribute the workload on resources. [9]

### 3. Scheduling Problem

The task scheduling problem consists of N tasks and M machines. Each task must be processed by one of the M virtual machines such that in the end the overall scheduling duration would be minimized. The proposed algorithm is focused on the Quality of Service parameters makespan, flowtime and task execution cost. Each task can be executed on just one resource and cannot be paused before the end. The algorithm utilizes the ETC matrix model explained in [6]. Since the proposed scheduling algorithm is static, the expected time for executing task j on resource i is assumed to be predetermined and set within ETC matrix [i, j].

Completion\_Time [i,j] is equal time that job j be completed on resource i and is computed as follows.

$$(1) \text{Completion\_Time}[i,j] = \text{ETC}[i,j]$$

Makespan: maximum completion\_Time [i,j], that is computed as follows.

$$(2) \text{Makespan} = \text{Max} (\text{Completion\_Time}[i,j]) \quad 1 \leq j \leq N, \quad 1 \leq i \leq M$$

Flowtime: sum of the completion time of jobs [i, j] over all resources, that is computed by follow equation.

$$(3) \text{flowtime} = \sum_{i=1}^m \sum_{j=1}^n \text{completion\_Time}[i,j]$$

The goal of scheduling in the proposed algorithm is to submit each of the jobs to each of the resources to minimize makespan and flowtime of the jobs at last.[10]

### 4. The proposed scheduling algorithm

In the proposed scheduling algorithm, the population diversity factor is used to alter the inertia weight for solving the task scheduling problem in cloud computing. Before presenting the algorithm, the parameters required for solving the scheduling problem with PSO algorithm need to be assessed.

#### 4.1 Representation of Particle Swarm

One of the important issues of using PSO algorithm for solving task scheduling problem is how to turn a scheduling problem into a solution or in other words how to form a mapping between the solution and the particles of PSO algorithm. In the PSO scheduling algorithm, each particle is a possible solution for task allocation problem. Each particle vector has a length of N (N is the number of input tasks); and each element inside this vector is a random integer between 1 and M (M is the total number of resources).

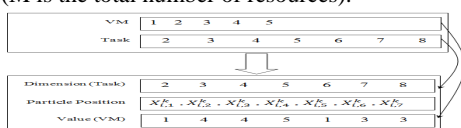


Figure 1. The mapping of tasks to resources

#### 4.2. Generating the Initial Population of Particles

In the proposed method, the initial population is generated at random. For this purpose, algorithm generates a random integer between 1 and M, representing the number of resource on which the task at hand will be executed. Randomness helps maintaining population diversity and gives all members of population an equal chance of being selected. Each particle vectors is as long as the number of tasks. The size of particle population determines the number of candidate solutions or the size of search in the problem space.

	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>
Particle <sub>1</sub>	VM <sub>2</sub>	VM <sub>5</sub>	VM <sub>3</sub>	VM <sub>2</sub>	VM <sub>1</sub>
Particle <sub>2</sub>	VM <sub>1</sub>	VM <sub>4</sub>	VM <sub>5</sub>	VM <sub>3</sub>	VM <sub>2</sub>
Particle <sub>j</sub>	VM <sub>1</sub>	VM <sub>5</sub>	VM <sub>5</sub>	VM <sub>4</sub>	VM <sub>2</sub>

Figure 2. Typical Particle

#### 4.3 Using the Population Diversity Factor for Changing the Inertia Weight

One of the problems of this algorithm is the premature convergence caused by rapid information exchange between particles which itself is because of lack of diversity in the particle population. Thus, one way of improving the overall performance of PSO algorithm is to maintain population diversity by improving indices such as inertia weight. In this method, we have used a regulator represented by kb [7, 8].

During the search, the regulator improves the static weight (w) by controlling the population diversity through negative feedback. The regulator index is shown with e and is expressed as follows:

$$e(t) = \frac{D_i(t) - D_0(t)}{D_i(t)} = \frac{D_0(t-1) - D_0(t)}{D_0(t-1)}$$

Where D<sub>0</sub>(t-1) is the extent of diversity at step t-1 and D<sub>0</sub> is the extent of diversity at step t.

Using the regulator index kb, static weight is defined as follows:

$$w(t) = kb \times e(t) = kb \left( 1 - \frac{D_0(t)}{D_0(t-1)} \right)$$

The suitable range for kb is between 0.9 and 1.6, and 1.4 is a good choice for most situations. In this version of method, the static weight will be adjusted according to population diversity (its decrease or increase). During the search, when population diversity D<sub>0</sub> decreases rapidly, inertia weight will increase to improve global search; when D<sub>0</sub> decreases very slowly, a small inertia weight w will be introduced to improve local search; and when D<sub>0</sub> fluctuates as a step progresses, a negative inertia weight will be used to reduce particle diversity. With this technique, the balance between global and local search will be maintained and the search will not converge to local optima. Diversity is measured using the following equation.

$$\text{diversity}(s(t)) = \frac{1}{n_s} (x + a)^n = \sum_{i=1}^{n_s} \sqrt{\sum_{j=1}^{n_s} (x_{ij}(t) - \hat{x}_j(t))^2}$$

Where (X<sub>j</sub>)(t) is the average of j-th dimension of all particles.

## 5. Simulation

In this section, the efficiency of the tasks scheduling algorithm presented in the previous section is evaluated. The objective of this algorithm is to schedule several independent tasks in a cloud environment. These tasks belong to an application that the user has requested to be executed on the cloud environment. Along with the task, the user specifies the QoS parameters, i.e. the time-cost optimization strategy that needs to be respected by the system. For example, the user can demand that cloud computing system should execute the application in the shortest time possible. Through simulation, the efficiency of time optimization algorithm can be compared with other algorithms.

### 5.1. Cloud computing environment model

A simulated environment for cloud computing includes one or more users and a number of resources. Users at any time may enter cloud computing and offer their application to run. This program includes a number of independent works, each of which can run on any desired resource. The user entity, after login, will be assigned its own broker entity, which performs its application scheduling and then returns results to the user after assigning tasks to resources and completing their implementation.

### 5.2. User model

Often in simulation experiments, the environment is modeled as single-user. The single-user version does not mean that only a user is logged in the system from simulation start to finish, but it means that the scheduling of any two users does not interfere with each other. For example, if a user is logged in at time 1500 and his application implementation finishes at time 4000 and another user is logged in at time 6000, then the system is called single-user. For comparison purposes, the system is presented as single-user.

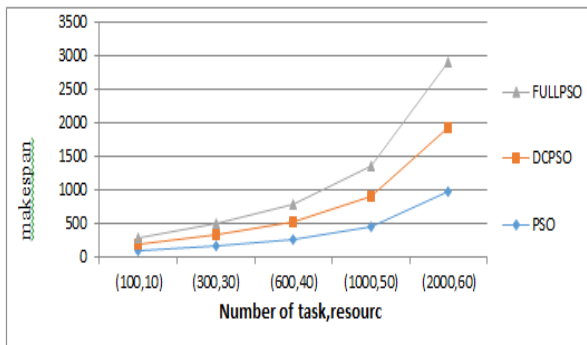


Figure 3. Diagram of makespan

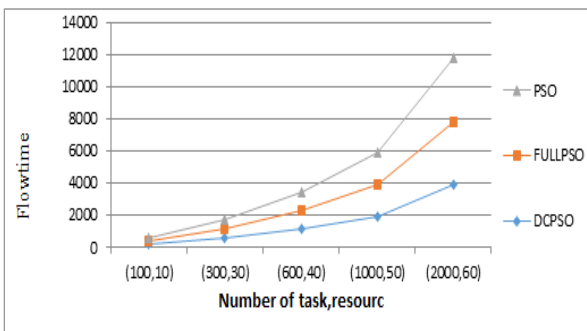


Figure 4. Diagram of flowtime

## 6. Conclusion

This paper provided a new method of resources allocation for cloud computing environment. The proposed algorithm utilized two high speed PSO-based methods which improved resource efficiency, selection of the best resource for task allocation, and minimized Makspan, while improving Flowtime parameter. The said parameters were simulated accurately. The simulation results showed the better performance of the proposed methods than other similar alternatives.

## 7. Suggestions

The future works are suggested to focus on development of a mechanism for partitioning applications or large problems into smaller sub-problems or tasks; as devising a mechanism for high speed and accurate portioning depending on task length (which is a determining load balancing factor in distributed systems especially cloud computing systems) will lead to improved speed and accuracy of scheduling and resource allocation, and thereby the increased efficiency of cloud computing system.

## 8. REFERENCES

- [1] Elzeki O. M. , M. Z. Rashad and M. A. Elsouid, "Overview of Scheduling Tasks in Distributed Computing Systems", International Journal of Soft Computing and Engineering (IJSCE), ISSN: 2231 - 2307, Vol. 2, Issue-3, pp. 470-475, 2012. Casavant, T.L., Kuhl, J.G. 1988 A taxonomy of scheduling in general-purpose distributed computing systems. IEEE Transactions on Software Engineering.
- [2] Pardeep Kumar and Amandeep Verma," Independent Task Scheduling in Cloud Computing by Improved Genetic Algorithm", International Journal of Advanced Research in Computer Science and Software Engineering, ISSN: 2277 128X, Vol 2, Issue 5, May 2012,pp:111-114.
- [3] Tinghuai Ma, Ya Chu, Licheng Zh., Otgonbayar A., "Resource Allocation and Scheduling in Cloud Computing: Policy and Algorithm", IETE Technical, Vol. 31, No. 1, 2014.
- [4] Bilgaiyan S., Sagnika S., Das M., "An Analysis of Task Scheduling in Cloud Computing using Evolutionary and Swarm-based Algorithms", International Journal of Computer Applications (0975 – 8887), Vol. 89 , No.2, pp. 11-18, 2014.
- [5] Selvarani S, Sadhasivam G., "Improved Cost-based Algorithm for Task Scheduling in Cloud Computing", Computational Intelligence and Computing Research (ICCIC), pp. 1-5, 2010.
- [6] Bittencourt LF, Madeira E., "A Cost Optimization Algorithm for Workflow Scheduling in Hybrid Clouds", In Journal of Internet Services and Applications, pp. 207-227, 2011.
- [7] J. Kennedy and R.Mendes, "Neighborhood topologies in fully-informed and best of neighborhood particle swarms,"vol. 12, pp. 50, 2003.
- [8] R.NEVES, "Watch Thy Neighbor Or How The Swarm Can Learn Form its Environment. " vol. 10, pp. 50, 202.
- [9] Dillon T,Wu Ch,Chang E.2010.Cloud Computing: Issues and Challenges.In Advanced Information Networking and Applications (AINA).24th IEEE International Conference on,pp:27-33.

- [10] Sajjad, Asadzadeh. Ch., Seyed Naser, R., Ali H.,(2014),  
“Job Scheduling on the Grid Environment using Max-  
min Firefly Algorithm”, International Journal of  
Computer Applications Technology and Research, Vol3.  
Issue 1, p.p 63-67

# Emotions in Humans and Intelligent systems

## (In terms of ‘Classical Indian Philosophy’ and Marvin Minsky’s ‘The Emotion Machine’)

G. Gnaneswari  
Asst. Professor, New Horizon College,  
Research Scholar (Jain University),  
Bangalore, India

M.V. Vijayakumar,  
Professor and PG Coordinator,  
Dr. Ambedkar Institute of Technology  
Bangalore, India

### Abstract:

Today ‘Emotional Intelligence’ has become a popular field of research. This paper compares and contracts the cross cultural philosophies with cognition in today’s digital world Artificial Intelligence. In this paper we discuss about what Indian philosophy has to say about human beings and their emotion which is completely different from what is believed to be the right thing in today’s psychology. That is because beliefs change with time. The belief of our ancestors is not what we believe today. And with these changing times if we have to build a machine with human intelligence that can understand emotions then it is a challenge. This makes the study quite interesting.

**Keywords** cognition, conversational agent, emotion, cognitive architecture, speech processing, behaviour, belief, philosophy.

### 1. INTRODUCTION:

There is a Common belief that ‘emotions are the one thing that will differentiate humans from machines’, so people who do not emote are criticized to be behaving like robots. The reason for human behaviour, can be due to the environment, the state of mind, the state of body and his/her belief/approach towards life. This will change from birth as he grows and experiences new things.

On the other hand, machines or robots are usually used only for automation. We hear

people say ‘But nothing made of mechanical stuff could ever have genuine feelings like love.’ They are not expected to have emotions for that kind of a job. But things are changing today; service industry wants to know all about the customer. In places like social networks he/she is been closely watched/tracked in order to predict his/her behaviour. So we may need robots with cognition as that of a human.

So generally it is believed that humans can be humans only when they have emotions, and machines will not have emotions. So these are the

two beliefs that will be discussed in this paper keeping in mind the eras they belong to. Thus in this paper we will discuss about What Indian philosophy has to say about human emotions? , What cognitive scientists have to say about emotions in machines?, And how similar/dissimilar are the answer to both the questions?



## 2. INDIAN PHILOSOPHY:

Indian Philosophy talks how human should get rid of emotions to reach the almighty. But in today’s world, human beings without any emotions are criticized as robots.

In Indian classical philosophy, term used for cognition is ‘gnana’ which is a mental phenomenon. The various emotions in Sanskrit are vedanā (feeling), bhāva (feeling), rāga (love, attraction), dveṣa (hatred), harṣa (joy), bhaya (fear) and śoka (sorrow). [1]

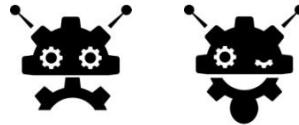
Many Indian philosophers regard ‘emotion’ is an obstacle to rational thinking and acquiring right knowledge.

For instance, ‘Nayaya-vaishesila’ philosophy involves a strict division into cognition and mental phenomena like feelings. It also treats emotions and dosas (defects) or upadha (impurities) which is the result of mithyajnana (ignorance). Thus according to Nayaya-vaishesila ‘ final

liberation is the ultimate aim and final liberation means the end of cognition.’

The ‘vedanta’ philosophy refers to mental qualities, such as manas (mind), buddhi (intellect), vijñāna (cognition) or citta (consciousness).

It also says that liberation from emotions and realisation of true self are the key to the vedantha philosophy.



Both these philosophies distinguishes emotion and feelings. But Sāṃkhya philosophy account to puruṣa (consciousness) and prakṛti (“reflection”, “activity”, “inertia” etc..). [1] It says that in order to experience a conscious intellect with cognition, emotion is an obstacle. The goal remains attaining pure contentless consciousness.

Interestingly, Patañjali Yoga philosophy claims that “the mind is ‘coloured’ by all of the objects it knows, including cognitions and emotions” which is very close to the western philosophy. [1] Western philosophy preaches the idea having cognition that free of emotions.

Thus most Indian philosophies wanted to get rid of ‘emotions’ which prevent liberation. Even ‘cognition’ is not accepted or eradicated completely; only the concept of vijñāna (knowledge part of cognition) was accepted.

In the so called modern philosophies such as the “ISHA” bases its programme known as the “Inner Engineering” on “Anaithukkam Assai Padu” meaning desire (emotion) everything.

Contrary to all the above beliefs, the psychologists of today talk on the need for emotion and cognition in humans. In psychology, emotions are attached to cognitive and biological changes taking place inside the human. They say emoting and expressing oneself will save you from loneliness and depression.

### **3. MARVIN MINSKY’S EMOTION MACHINE:**

Here is what Dr. Marvin Minsky says in his book ‘The Emotion Machine’ about why and how can machines have emotions. The main difference between human and machine is that humans evolve in every way intellectually, emotionally etc., Minsky argues how can cognitions such as emotion, commonsense, consciousness can be programmed to a robot. He says that ‘none of those popular psychology words refers to any single, definite process instead of those words attempts to describe the effects of large networks of processes inside our brain’. There cannot be anything that is only pure logic and rational because in humans cognition is affected by assumptions, values and purpose.

### **4. SPECTRUM OF EMOTIONS USED IN BUILDING A CONVERSATIONAL AGENT:**

A speech contains information such message, speaker details, language specific intricate details, emotion etc. In a conversation, Agent 1’s dialogue contains all these details; if Agent 2 has considered all these element except his emotions while replying, then reply will be robotic. Thus reading the emotions from a conversation becomes pertinent. This paper talks about the need for the study of emotions in developing a Cognitive Conversational Agent. One of the key tasks in building a Cognitive agent is to understand the emotions of the Agent1 in order to efficiently reply to Agent 2. Since emotions cannot be simply pinpointed a spectrum representation of emotion is essential for the study of emotions.

Study of emotions is becomes a must in the process of building an emotional conversational agent. Researchers have come up with various analytical methods to understand emotions from speech. Today Conversational agents are not about just conversing with human with proper syntax and semantics. Pragmatics has become a significant part of this particular study since the other parts are already been taken care.

Though many research work has been done on emotions in speech; this study uses the following methods: (i) 15 sentences pertaining to different emotions (ii) The speech signal is sampled at 16 kHz, and represented as 16 bit numbers (iii) LSTM / RNN / HMM is used for sampling (iv) Audio is represented by spectral features like Mel-frequency cepstral coefficients (MFCC) and linear prediction cepstral coefficients (LPCC) (v) Training and



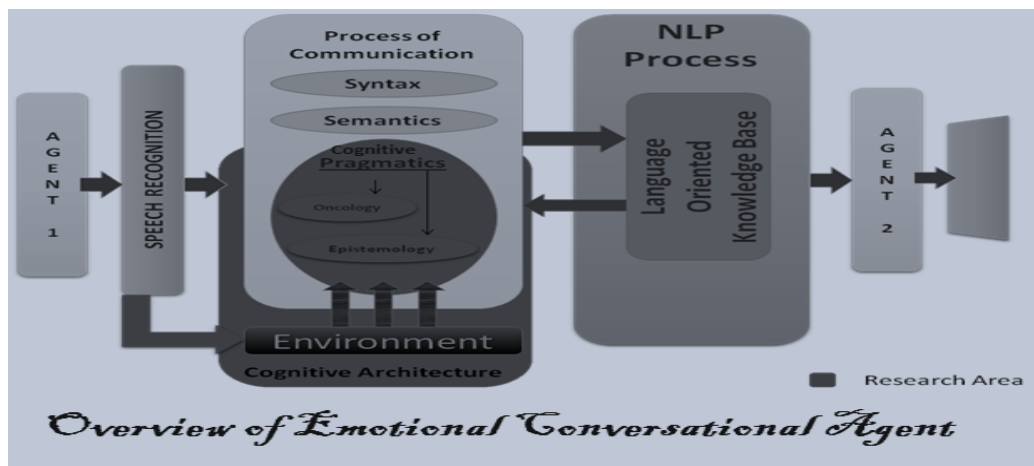
Testing code written using Matlab DSP.

For example in the following research, each experiment, the highest recognition accuracy achieved, its variance, the inputs features and clustering parameters used, is listed in Table.[2]

Thus there is possibility of measuring humans emotions.[3] These can be used by agents as input and make it capable of decision making.

All Speakers						
Experiment	Features	Distance Measure	Centroid	Iterations	Recognition Accuracy	Variance
despair-elation	MFCC	L1 norm	UDC	100	75.76%	1.74%
happy-sadness	MFCC	L1 norm	UDC	1	77.91%	14.34%
interest-boredom	Pitch	L1 norm	UDC	100	71.21%	2.48%
shame-pride	MFCC	L1 norm	UDC	1	73.15%	3.23%
hot anger-elation	MFCC	L1 norm	UDC	1	69.70%	10.75%
cold anger-sadness	MFCC	L1 norm	UDC	1	75.66%	3.35%

#### 4. EMOTIONAL CONVERSATIONAL AGENT ARCHITECTURE:



In the above paper a model of the Conversational Agent for communication between two perceptual agents keeping in mind all the various factors that are involved in developing a Conversational Agent is designed. This research

paper emphasizes on ‘Cognitive Pragmatics’ for developing a Conversational Agent based on Cognitive Architecture. [6]

There understanding of the conversation can be based on the content of the speech, facial

expressions of the speaker or by the features extracted from the emotional speech[5]. These three elements comprises of the environment.

#### **Future study:**

The study of human mind and its development helps in designing a cognitive architecture for robots. How can an architecture be designed by using Indian philosophy’s idea of a working mind and also the psychology behind the human mind. Such cognition studies paves way to building intelligent systems in terms of their thinking and their actions. A very interesting research area building intelligent systems in terms of the way they think and converse. Problems that work is happening now are common sense, consciousness, belief-intention-desire, decision making, goal driven, etc in solving simple and complex problems.

#### **7. REFERENCES:**

- [1] Tuske, Joerg, "The Concept of Emotion in Classical Indian Philosophy", The Stanford Encyclopedia of Philosophy (Spring 2011 Edition), Edward N. Zalta (ed.), URL = <<http://plato.stanford.edu/archives/spr2011/entries/concept-emotion-india/>>.
- [2] Shah Hewlett, "Emotion Detection from Speech", URL = <http://cs229.stanford.edu/proj2007/ShahHewlett-EmotionDetectionfromSpeech.pdf>
- [3] R.Banse, K.R.Scherer, "Acoustic profiles in vocal emotion expression", Journal of Personality and Social Psychology, Vol.70, 614- 636, 1996
- [4] V.A Petrushin, "Emotional Recognition in Speech Signal: Experimental Study, Development, and Application", ICSLP-2000, Vol.2, 222-225, 2000, L.R.Rabiner
- [5] S. G. Koolagudi and K. S. Rao, "Exploring speech features for classifying emotions along valence dimension," in The 3rd international Conference on Pattern Recognition and Machine Intelligence (PReMI-09), Springer LNCS (S. C. et al., ed.), (IIT Delhi), pp. 537–542, Springer-verlag, Heidelberg, Germany, December 2009

#### **6. CONCLUSION:**

Thus it was right when it was believed yesterday that suppressing all the human emotions is a way to reach almighty and those machines will and need not emote. But today beliefs have changed. Psychologists say that it is best to express your emotions and not suppress them. Cognitive scientists have also understood the need for machines with cognition and are working on developing cognitive machine. Only these researches can pave the way towards creating robots which behave just like humans not just physically and mentally but emotionally too. This paper also concludes that conversational agents should be built keeping in mind the emotions.

[6] G Gnaneswari, Dr.M.V.Vijayakumar, “Emotional Conversational Agents based on Cognitive Pragmatics”, International Journal of Pharma & Bio Sciences, June 2016 Issue, Special Edition on Computational Data Science, ISSN 0975-6299, PP: 15-19

# Machine Learning Algorithms for Recommender System - a comparative analysis

Satya Prakash Sahu  
University of Hyderabad  
Hyderabad, India  
spsahu@uohyd.ac.in

Anand Nautiyal  
University of Hyderabad  
Hyderabad, India  
anandnautiyal@uohyd.ac.in

Mahendra Prasad  
University of Hyderabad  
Hyderabad, India  
je.mahendra@uohyd.ac.in

---

**Abstract:** Recommendation system is one of the most popular applications of Artificial Intelligence which attracts many researchers all over the globe. The advent of the Internet era has brought wide implementation of recommendation system in our everyday lives. There are many machine learning techniques which can be used to realize the recommendation system. Among all these techniques we are dealing with Content Based Filtering, Collaborative Based Filtering, Hybrid Content-Collaborative Based Filtering,  $k$ -mean clustering and Naive Bayes classifier. We have exploited these algorithms to their extreme in order to achieve the best possible precision and have presented a comprehensive comparative analysis. The strength of all these algorithms can be clearly realized by the significant enhancement in the accuracy, depicted by the experimental analysis taking cold start problem into consideration.

**Keywords:** Recommender System, Classifier, Content Based, Collaborative Based, Cluster, Correlation.

---

## 1. INTRODUCTION

Recommendation system[10] is an application which is used for prediction in various domains throughout the internet. A large amount of data flows through the internet and it gives away a lot of information regarding the user searching activity. The information extracted from the pattern of previously searched data can be molded into the prediction of relevant data for the user[1]. The implementation of the system can be performed by various techniques. In this paper, we have discussed Content Based Filtering, Collaborative Filtering[10], Hybrid Content-Collaborative Based Filtering,  $k$ -mean clustering Based and Naive-Bayes Classifier based techniques.

The Content Based Filtering approach takes into account a user's profile which is constructed based on his previous ratings[2]. His ratings determine his inclination and interests, forming the basis for recommending a new item. A higher rating denotes a higher likelihood of the user to visit similar items. So, a new item is recommended according to the maximum number of ratings given by the user in a genre[3].

In the Collaborative Based Filtering, recommendation for a user is governed by other users' profiles. An item is recommended based on the ratings of other users who have similar interests as the user under consideration[2][4]. In another approach, the content and collaborative based filtering are combined to form the Hybrid Content-Collaborative Based Filtering. It includes the advantages of both the methods and outperforms both of them.

In the  $k$ -mean clustering, the similarity between the objects is calculated by the means of various distance measures such as Euclidean distance[5], Pearson Correlation, etc. The value of  $k$  determines the number of clusters to be formed[6]. The nearest  $k$  objects are the most similar to one another. These clusters of similar objects drive the recommendation of new arriving objects. Naive Bayes is another popular and efficient classifier based on Bayes theorem. It is a conditional probability based classifier. The prior knowledge of the classifier assists learning. The naive assumption is that the features are conditionally independent[7].

In this paper, we have used the MovieLens dataset[8]. All the above algorithms deal with this dataset in order to recommend the movies and calculate the precision along with tackling the cold-start problem[3]. Cold-start problem is one of the most commonly encountered challenges of the recommendation system. It is also known as the new user problem as it creates problem of generating recommendations for the new user. We have divided this analysis into various sections. Section II describes the different state-of-the-art techniques for the recommendation system. Section III gives the experimental results for all these techniques. Section IV concludes the study. Section V describes the future work that we propose.

## 2. ALGORITHMS

### 2.1 Content Based Filtering

The Content Based Filtering considers the items rated by a user to formulate the future recommendations while exploring the internet services. A user tends to rate an item which he likes or dislikes. His ratings reflect his response towards that item. If he likes an item, he rates it higher and lesser ratings denote that he is not much interested. These rated items serve as the 'content' in the Content Based Filtering[2][3]. Based on this content, the user is recommended future items which he might approve of. Here, the user is recommended movies which fall in a particular genre of his liking.

**Algorithm 1.** Content Based Filtering

---

**Input:** users  $X$ , movies  $m$ , rating  $r$ , movie genre  $m_g$ , Number of movies to be recommended( $\mu$ ).

**Output:** Recommended movies  $R$

---

1. for all users do
2. Select seen movies  $s$ , unseen movies  $s'$ , association of unseen movies  $as_i'$  w.r.t  $X$ , association of each genre  $ag_j$  w.r.t  $s'$ , where  $i$  is 1 to  $n$  and  $j$  is 1 to  $m$ .
3. Calculate  $score_j$ .
4. Select highest three  $score_j$
5. Select  $m' \subset s'$  according to highest three  $score_j$

6. Calculate score  $m_e'$  where  $e \in m'$
7. Return top  $\mu$  score recommendations.
8. end for

In this algorithm, the notations used have the following meaning : association of each movie  $as_i$  represents total number of users who rated movie  $i \in s'$ , association of each genre  $ag_j$  represents total number of movie belonging to genre  $j$ .

$$score_j = ag_j / m$$

$$score(m_e') = am_e / \text{total count of } m'$$

## 2.2 Collaborative Filtering

There can be many users who must be having the same pattern of rating an item as the user intended. This similar pattern of their ratings with the user guides the Collaborative Filtering[2][3][10]. The notion behind the Collaborative Filtering is the recommendation of an item based on the preferences of like-minded users.

### Algorithm 2. Collaborative Filtering

**Input:** users  $X$ , movies  $m$ , rating  $r$ , Number of movies to be recommended( $\mu$ ).

**Output:** Recommended movies  $R$ .

1. for all users do
2. Select seen movies  $s$ , unseen movies  $s'$
3. Find similarity ( $sim_i$ ) w.r.t  $s$ , where  $i = 1$  to  $n$ .
4. Select highest  $sim_i$  user
5. Select  $m' \in s$  of user obtained in step 4 and  $s'$  of  $i^{th}$  user.
6. Calculate weight  $W(m_e')$  where  $e \in m'$
7. Return top  $\mu$  weight recommendations.
8. end for

In this algorithm, the notations used have the following meaning :  $sim_i$  represents common movies between user  $i$  and other users.

$$weight(m_e') = \text{rating of particular movie}_e / \text{max rating.}$$

## 2.3 Hybrid Filtering

To cater better precision, a hybrid filtering method is used which can provide the advantages of both the content and the collaborative approaches[4] and can overcome their shortcomings. Suppose, the user appreciates mostly movies in  $g \subset G$  genres, and the collaborating users also give high ratings to the  $g \subset G$  genres, then  $g$  will be taken as the metric to recommend movies to the user.

### Algorithm 3 .Hybrid Filtering

**Input:** users  $X$ , movies  $m$ , rating  $r$ , movie genre  $m_g$ , Number of movies to be recommended( $\mu$ ).

**Output:** Recommended movies  $R$ .

1. for all users do
2. Select seen movies  $s$ , unseen movies  $s'$ , association of each genre  $ag_j$  w.r.t  $s'$ , where  $i$  is 1 to  $n$  and  $j$  is 1 to  $m$ .

3. Calculate  $score_j$ .
4. Select highest three  $score_j$
5. Select  $m'' \in s$  of the  $i^{th}$  user according to highest three  $score_j$
6. Find similarity ( $sim_j$ ) w.r.t  $m''$
7. Select highest  $sim_j$  user.
8. Select  $m'$  according to its highest three  $score_j \in s$  of user obtained in step 7 and  $s'$  of the  $i^{th}$  user under consideration.
9. Calculate weight  $W(m_e')$  where  $e \in m'$
10. Return top  $\mu$  weight recommendations.
11. end for

## 2.4 K-Mean Clustering

The  $k$ -mean is a non parametric classification technique. It distributes the items into  $k$  clusters according to their proximity to one another. In this paper, this proximity is being measured by using the Euclidean distance[11]. For calculating the Euclidean distance we have taken rated and unrated movies as binary. Each cluster possesses a centroid which is the mean of all the items in the cluster. All the objects in a cluster move towards the centroid and the centroid is updated in each iteration. The iteration continues until a saturation point arrives, when the centroid stops altering. By following this approach we are decreasing the search space which results in reduced computational complexity[6]. These computations are performed off-line which helps the classification to be efficient in terms of time complexity.

### Algorithm 4. $k$ -mean clustering

**Input:** users  $X$ , movies  $m$ , rating  $r$ , Number of movies to be recommended  $\mu$ , value of  $k$ .

**Output:** Recommended movie  $R$ .

1. begin
2. Randomly select  $k$  centroids.
3. Calculate euclidean distance ( $euclid$ ) for  $X$  from  $k$  centroids.
4. Allocate  $X$  to  $k^{th}$  cluster according to  $euclid$ .
5. Update centroid for each cluster with (summation( $k_i$ ) from 1 to  $p$ )/ $p$ , where  $p$  is the number of members in  $k_i$  cluster
6. Repeat step 3 to step 5 until centroid( $t$ )  $\neq$  centroid ( $t+1$ ).
7. for all users do
8. Select seen movies  $s$ , unseen movies  $s'$ .

9. Find similarity ( $sim_i$ ) w.r.t  $s$ , where  $i = 1$  to  $p$ .
10. Select highest  $sim_i$  user.
11. select  $m' \subset s$  of highest  $sim_i$  and  $s'$  of  $i^{th}$  user.
12. Calculate weight  $W(m_e')$  where  $e \in m'$
13. Return top  $\mu$  weight recommendations.
14. end for
15. end

## 2.5 Naive Bayes

The Naive Bayes is based on the Bayes theorem. The probabilistic approach followed by Naive Bayes Classifier determines the probability of the classification and helps in finding the uncertainty about the model[9]. It is an efficient learning algorithm which uses the prior knowledge of the observed data. The Naive assumption is that the features are conditionally independent[1].

### Algorithm 5. Naive Bayes

**Input:** users  $X$ , movies  $m$ , rating  $r$ , number of movies to be recommended( $\mu$ )

**Output:** Recommended movies  $R$ .

1. for all users do
2. Select seen movies  $s$ , unseen movies  $s'$ .
3. Find similarity ( $sim_i$ ) w.r.t  $s$ , where  $i = 1$  to  $n$ .
4. Select  $x' \subset X$  where  $sim_i > 10$ .
5. Calculate association of unseen movies  $as_i'$  w.r.t to  $x'$
6. Calculate score ( $s_e'$ ) where  $e \in s'$ .
7. Return top  $\mu$  score recommendations.
8. end for

## 3 EXPERIMENTAL RESULT

We now illustrate the analysis of the experiments performed and provide a comparison of all the state-of-the-art methods described above. To compare their accuracy we have used the MovieLens dataset of 10K, 50K and 100K. The dataset varies in sparsity. For example, the 100K MovieLens dataset has 100K ratings, 943 users and 1682 movies of 19 different genres. The analysis of these algorithms is demonstrated based on precision measure. For each test user, we convert 30% of the user's seen movies into unseen movies and apply the algorithms described above. Out of the total number of

recommendations (T), the ones which are also present in the converted movies are the correct recommendations(tc).

$$\text{Precision} = (\Sigma tc / \Sigma T) * 100$$

For all the experiments, we are taking value of  $\mu = 5$  and value of  $k = 10$ .

Algorithm/Size	10K	50K	100K
Content Based	18.45	18.66	19.10
Collaborative	17.97	18.69	19.95
Hybrid	20.31	21.03	22.20
K-Mean	21.05	21.93	22.67
Naive Bayes	24.62	25.19	25.73

Table 1. Precision of Different Algorithms.

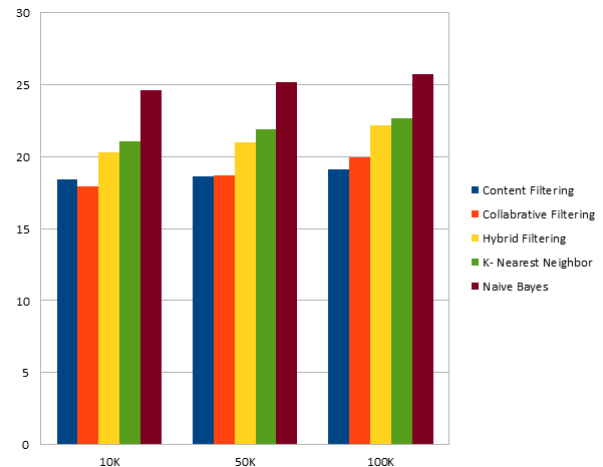


Fig. 1. Precision Comparison

## 4 CONCLUSION

All the algorithms described in this paper are compared with respect to their precision rates. This comprehensive analysis depicts the strength and the weakness of each one of them in different versions of the MovieLens dataset. The experiments performed are the witness of the sparsity handling by these algorithms. Our experiments have shown promising results and this paper conforms that out of all these approaches Naive Bayes gives the best precision.

## 5 FUTURE WORK

With this paper, we have achieved encouraging results from all these algorithms. In the real time sophisticated recommendation systems there is a need of high accuracy. Such systems still have space for improvement. There are several machine learning algorithms which can be applied to these real time systems. It is worthwhile to examine those other algorithms to improve the precision further.

## 6. REFERENCES

- [1] Katore L.S., and Umale J.S. Comparative Study of Recommendation Algorithms and Systems using WEKA, International Journal of Computer Applications, Volume 110 – No. 3, pp14-17, 2015.
- [2] Tewari A.S., Kumar A., and Barman A.G., Book Recommendation System Based on Combine Features of Content Based Filtering, Collaborative Filtering and Association Rule Mining, IEEE, 978-1-4799-2572-8, pp 500-503, 2014.
- [3] Wanaskar U.H., Vij S.R., and Mukhopadhyay D. A Hybrid Web Recommendation System Based on the Improved Association Rule Mining Algorithm, Journal of Software Engineering and Applications, 6, pp 396-404 2013.
- [4] Shinde U., and Shedge R. Comparative Analysis of Collaborative Filtering Technique, IOSR Journal of Computer Engineering, Volume 10, pp 77-82 2013.
- [5] CHENG X., WANG J., Danqian LU. Research of Question Analysis Based on HNC and K Nearest Neighbor, Journal of Computational Information Systems, 6:10, pp 3449-3455, 2010.
- [6] Campos P.G., Bellogín A., Díez F., and Chavarriaga J.E. Simple Time-Biased KNN-based recommendations, ACM, 978-1-4503-0258-6, 2010.
- [7] Puntheeranurak S., and Pitakpaisarnsin P. Time-aware Recommender System Using Naïve Bayes Classifier Weighting Technique, International Symposium on Computer, Communication, Control and Automation, 3CA, pp 266-269 ,2013.
- [8] Bellogín A. Castells P. and Cantador I. Precision-Oriented Evaluation of Recommender Systems: An Algorithmic Comparison, ACM, 978-1-4503-0683-6, 2011.
- [9] Puntheeranurak S. and Sanprasert S. Hybrid Naive Bayes Classifier Weighting and Singular Value Decomposition Technique for Recommender System, IEEE, 978-1-4244-9698-3, pp 473-476, 2011.
- [10] Laishram A., Sahu S.P., Padmanabhan V., and Udgata S. K. Collaborative Filtering, Matrix Factorization and Population Based Search: The Nexus Unveiled. ICONIP , Part III, LNCS 9949, pp. 352–361, 2016.
- [11] Prasad M. and Singh A. A Novel Hybrid Ant Colony Optimization Approach to Terminal Assignment Problem. ACM, AICTC '16, Article 29, August 12-13, 2016.

# Signal Jump Detection Process

Rohit Chandrashekhar Iyer  
Department of Electronics  
Engineering  
Priyadarshini College of  
Engineering Nagpur,  
Maharashtra, India

Rohit S. Garode  
Department of Electronics  
Engineering  
Priyadarshini College of  
Engineering Nagpur,  
Maharashtra, India

Aniket P. Bhojar  
Department of Electronics  
Engineering  
Priyadarshini College of  
Engineering Nagpur,  
Maharashtra, India

---

**Abstract:** Now a days we see that various types of accidents happens on the road. In India many accidents are caused due to human negligence. License Plate detection and recognition is a key technique in most of the traffic related applications such as signal jumping road traffic monitoring, airport gate monitoring, speed monitoring and Automatic parking access control .It is simply the ability to automatically extract and recognition of the vehicle license number plate's character from a captured image .In this paper, we try to give an enhance view of the signal jump detection and recognition of number plate.

**Keywords:** Signal jump detection, Number plate detection, Plate Recognition, Image processing, ANPR system, Vehicles detection, character detection

---

## 1. INTRODUCTION

With the growth of the urbanization, industrialization and population, there has been a tremendous growth in the traffic. There is occurrence of bundle of problems too, these problems include signal jump, traffic jams, accidents and traffic rule violation. In 1868, the traffic lights only installed in London and today these have installed in every cities around the world. Today red light violation is one of the most common and serious problem which results in the collision of millions of vehicles at the traffic light signals every year. A red light violation occurs when a vehicle try to cross the intersection at the red traffic light. So to give the punishment to the drivers of these vehicles, we must identify the vehicle that violates the traffic light signals

Number Plate Recognition (NPR) is an image technology used to identify plates for their vehicles. This technology is gaining popularity in security and traffic facilities. The purpose of NPR was to build a system capable of automatically recording of the license plate numbers of signal jump traveling down a roadway.

## 2. EASE OF USE

### 2.1 Use of Microcontroller ATmega16[5]

Standard for a microcontroller based low cost platform. It consists of an Atmel ATmega16 is a low power CMOS 8-bit microcontroller based on the AVR enhanced RISC architecture. By executing powerful instructions in a single clock cycle, the ATmega16 achieves throughputs approaching 1 MIPS per MHz allowing the system designed to optimize power consumption versus processing speed. The AVR core combines a rich instruction set with 32 general purpose working registers. All the 32 registers are directly connected to the Arithmetic Logic Unit (ALU), allowing two independent registers to be accessed in one single instruction executed in one clock cycle. The resulting architecture is more code efficient while achieving throughputs up to ten times faster than conventional CISC microcontrollers.

The ATmega16 provides the following features: 16K bytes of In-System Programmable Flash Program memory with Read While Write capabilities, 512 bytes EEPROM, 1K byte SRAM, 32 general purpose I/O lines, 32 general purpose working

registers, On-chip Debugging support and programming, three flexible Timer/Counters with compare modes, Internal and External Interrupts, a serial programmable USART, a byte oriented Two-wire Serial Interface, an 8-channel, 10-bit ADC with optional differential input stage with programmable gain, a programmable Watchdog Timer with Internal Oscillator, an SPI serial port, and six software selectable power saving modes. The Idle mode stops the CPU while allowing the USART, Two wire interface, A/D Converter, SRAM, Timer/Counters, SPI port, and interrupt system to continue functioning, The ADC Noise Reduction mode stops the CPU and all I/O modules except Asynchronous Timer and ADC, to minimize switching noise during ADC conversions.

1. Atmega 16 run at 16Mhz clock.8051 run at lower clock speed.
2. Atmega16 has inbuilt ADC.8051 doesn't has ADC.
3. RAM and ROM memory of Atmega16 is more than 8051.
4. Form Factor of Atmega16 is cheaper than the 8051.
5. Programmer for Atmega16 is cheaper than the 8051.
6. Atmega16 have 16kb flash memory and Atmega have 32 kb.
7. Atmega have 1kb SRAM and Atmega have 2kb SRAM.
8. Atmega have 512bytes EEPROM and Atmega have 1kb.

In this project we doesn't need more memory and RAM so we are not using Atmega32. And the cost of Atmega32 is also high.

### 2.2 Use of IC-MAX 232(Level shifter)

MAX 232 converts signals from a RS 232serial port to signals suitable for use in TTL-compatible digital logic circuits. The MAX232 is a dual transmitter / dual receiver that typically is used to convert the RX, TX, CTS, RTS signals. The drivers provide TIA-232 voltage level outputs about  $\pm 7.5$  TO 12 Volts from a single 5-volt supply by on-chip charge pump and external capacitors. When a MAX232 IC receives a TTL level to convert, it changes a TTL logic 0 to between +3 and +15 V,



and changes TTL logic 1 to between  $-3$  and  $-15$  V, and vice versa.

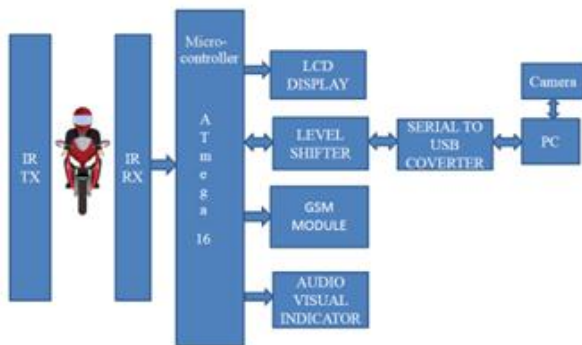
### 2.3 Use of IR Sensor

IR stands for infrared sensor. IR Sensors work by using a specific light sensor to detect a select light wavelength in the Infra-Red (IR) spectrum. By using an LED which produces light at the same wavelength as what the sensor is looking for, you can look at the intensity of the received light. When an object is close to the sensor, the light from the LED bounces off the object and into the light sensor. This results in a large jump in the intensity, which we already know can be detected using a threshold. A electronic remote device mainly consists of this IR transmitter and receiver. The IR signal is modulated during transmission. And demodulator during reception

### 2.4 Use of GSM module

The GSM module has the SIM900A microcontroller, a SIM slot, a 12V input, RS232 (Serial) interface and Pin outs for interfacing with controllers. It is suitable for SMS, Voice as well as DATA transfer application in M2M interface.

## 3. Block Diagram of Signal Jump Detection Process

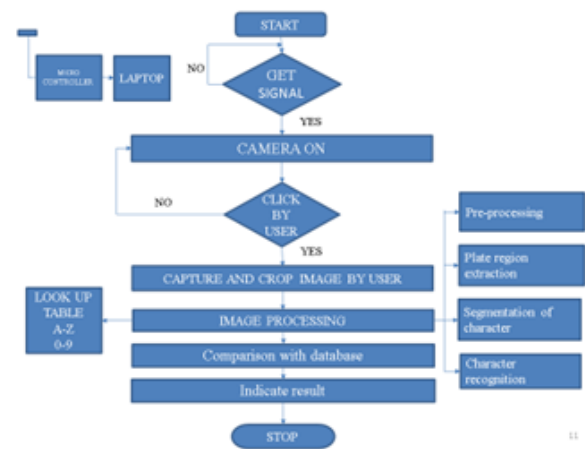


## 4. Working

IR Transmitter and Receiver are connected at line of sight on the two sides of the road. If the vehicle is between Transmitter and Receiver there is deflection in output of IR receiver and this signal is given to microcontroller. After the sensing is done, the signal will be received by the Atmega 16 microcontroller which is used in interfacing the sensors and the computer which is used to find the number of the vehicle. The audio and visual indicator is there to give help to the traffic inspector on the signal that someone has broken the signal. As the signal is broken by the vehicle the Atmega 16 will start working and the

LCD display connected is there to display that from where side the signal is broken. At the output of Atmega16 there is level shifter connected to it. The level shifter is mainly connected to administer the appropriate voltage to the computer (Laptop) as it runs on 12v and the output of the Atmega16 is approx. 5-6v. And if the level shifter is not connected there will be garbage value or unwarranted results will show up. The tool which we are using for NPR (Number Plate Recognition) is MATLAB12.0. The flowchart of processes in MATLAB is as follows.

## 5. PC, MATLAB AND CAMERA TYPESET TEXT



The main and the most important portion of this system is the software model. The software model uses series of image processing techniques which are implemented using MATLAB. The NPR algorithm is broadly divided into following parts.

1. Capture image.
2. Pre-processing.
3. Plate region extraction.
4. Segmentation of character.
5. Character recognition.
6. Comparison with database.
7. Indicate result.

### 5.1 Capture Image

The first step is the capturing of an image using camera. For this project, the test images of vehicles are taken with a camera which is there on the Laptop. The images will be stored as colour JPEG format in system. The next step is to use the Matlab function to convert the captured vehicle JPEG image into gray scale format.



Figure 1 Original Captured Image

## 5.2 Pre-processing

**Filtering:** When the image is saved, there is noise present in the image. To remove noise from the image median filters are used so that image becomes free from noise. Noise removal is necessary step in License plate recognition system because it greatly affects the recognition rate of the system.

**Gray Processing:** It involves conversion of color image into a gray image. The method is based on different color transform. According to the R, G, B value in the image, it calculates the value of gray value, and obtains the gray image at the same time.



Figure 2 Grey Scale image

## 5.3 Plate Region Extraction[2]

The most important step in the process of automatic number plate recognition is a detection of a number plate area. We can include algorithms that are able to detect a rectangular area of the number plate in an original image. Human beings define a number plate in a natural language as a “small plastic or metal plate attached to a vehicle for official identification purposes”, but machines do not understand this definition as well as they do not understand what “vehicle”, “road”, or whatever else is. Because of this, there is a need to find an alternative definition of a number plate based on descriptors that will be comprehensible for machines. Let us define the number plate as a “rectangular area with increased occurrence of

horizontal and vertical edges”. The high density of horizontal and vertical edges on a small area is in many cases caused by contrast characters of a number plate, but not in every case. This process can sometimes detect a wrong area that does not correspond to a number plate. Because of this, we often detect several candidates for the plate by this algorithm, and then we choose the best one by a further heuristic analysis. The edges for an image are always the important characteristics that offer an indication for a higher frequency. Detection of edges for an image may help for image segmentation, data compression, and also help for well matching, such as image reconstruction and so on.

There are many methods to make edge detection. The most common method for edge detection is to calculate the differentiation of an image. The first-order derivatives in an image are computed using the gradient, and the second order derivatives are obtained using the Laplacian. Another method for edge detection uses Hilbert Transform.

A number plate can be extracted by using image segmentation method. There are numerous image segmentation methods available in various literatures. In most of the methods image binarization is used. To find the region of image we will calculate the centroid and boundary and we have some condition we will apply the further procedure as shown to extract the number plate in MATLAB. Following figure shows the extracted number plate.

Equations:-

$$b = a(R/3, 1, C);$$

a = Original Image, R= Row, C= column.

Find the area number using following equations

$$B = \text{STATS. BoundingBox};$$

$$X_{\min} = B(2);$$

$$X_{\max} = B(2) + B(4);$$

$$Y_{\min} = B(1);$$

$$Y_{\max} = B(1) + B(3);$$

$$LP = b(X_{\min} + 25: X_{\max} - 20, Y_{\min} + 10: Y_{\max} - 10);$$



Figure 3 Cropped Image Number Plate



Figure 4 Filtered Extracted Number Plate

### 5.4 Segmentation of Character[1]

The next step after the detection of the number plate area is a segmentation of the plate. The segmentation is one of the most important processes in the automatic number plate recognition, because all further steps rely on it. If the segmentation fails, a character can be improperly divided into two pieces, or two characters can be improperly merged together. We can use a horizontal projection of a number plate for the segmentation, or one of the more sophisticated methods, such as segmentation using the neural networks. If we assume only one-row plates, the segmentation is a process of finding horizontal boundaries between characters. The second phase of the segmentation is an enhancement of segments. The segment of a plate contains besides the character also undesirable elements such as dots and stretches as well as redundant space on the sides of character. There is a need to eliminate these elements and extract only the character.

Segmentation of plate using a horizontal projection. Since the segmented plate is deskewed, we can segment it by detecting spaces in its horizontal projection. We often apply the adaptive thresholding filter to enhance an area of the plate before segmentation. The adaptive thresholding is used to separate dark foreground from light background with non-uniform illumination.



Figure 5 Character Segmentation

### 5.5 Look up table[3]

0	1	2	3	4	5
6	7	8	9	A	B
C	D	E	F	G	H
I	J	K	L	M	N
O	P	Q	R	S	T
U	V	W	X	Y	Z

### 5.6 Character recognition and Comparison with look up table[3]

The segmented character is now used to compare with individual character against the complete alphanumeric look up table. It match individual character and finally the number is identified and stored in string format in a variable. The string is then compared with the stored database for the vehicle authorization then recognized number plate string is compare with authenticated database file, , if the both value is same means it will display the authorized otherwise it will display the unauthorized.

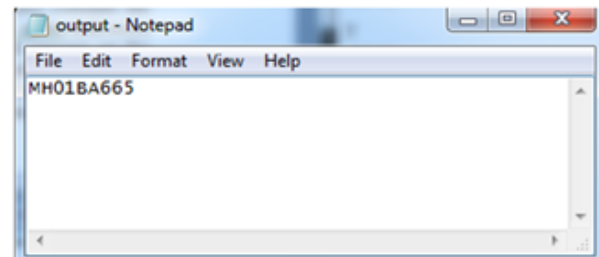


Figure 6 Recognized Number Plate

## 6. Conclusion

In this system, an application software is designed for the detection of number plate of vehicles using their number plate. At first plate location is extracted using edge detection operation then separated the plate characters individually by segmentation. Finally template matching is applied with the use of correlation for recognition of plate characters. In this project we are trying to reduce traffic violation in India which is the biggest problem facing in India. And take strict action on them. Some of possible difficulties:

1. Broken number plate.
2. Blurry images.
3. Number plate not within the legal specification.
4. Low resolution of the characters.
5. Poor maintenance of the vehicle plate.

## 7. REFERENCES

- [1] Chetan Sharma and Amandeep Kaur “Indian Vehicle License Plate Extraction and Segmentation” International Journal of Computer Science and Communication Vol. No. 2, pp. 593-599, July-December 2011
- [2] Kumar Parasuraman, Member, IEEE and P.Vasantha Kumar, "An Efficient Method for Indian Vehicle License Plate Extraction and Character Segmentation", 2010 IEEE International Conference on Computational Intelligence and Computing Research.
- [3] P.Subbuthai, Azha Perisamy and S.Muruganand, “Identifying the Character by Applying PCAMethod using Matlab” International Journal of Computer Applications (0975-8887), Volume 60-No.1, December 2012.
- [4] Chitode. J. S, Rupali Kate, “Number Plate Recognition Using Segmentation”, International Journal of Engineering Research & Tehnology (IJERT), Vol. 1 Issue 9, and ISSN: 2278-0181, 2012.
- [5] Microcontroller Datasheet At Mega 16, [www.atmel.com/devices/atmega16](http://www.atmel.com/devices/atmega16)
- [6] Chirag N. Paunwala, Suprava Patnaik, “A Novel Multiple License Plate Extraction Technique for Complex Background in Indian Traffic Conditions”, In Proceedings of International Journal of Image Processing, vol.4, issue2, 2007

# Towards Ontology Lifecycle: Building, Matching and Evolution to Semantically Integrate Application Ontologies

Razika Driouche  
National High College  
of Biotechnology, Taoufik  
Khaznadar, Algeria

---

**Abstract:** Semantic interoperability among applications, systems, and services are mostly based on ontology. Its increase usage in Information Systems and knowledge sharing systems raises the importance of ontology development and maintenance. It is essential for sharing information among independent organizations, exchange of information among heterogeneous systems. To make this possible, we need to carefully model the domain knowledge while preserving its semantics. Ontologies are complex in nature and often structured. Their development and maintenance incorporate research areas like: building, evolution, versioning, matching and integration where these are fundamentally different. We uncover the gap in the current research area of ontology building, matching and evolution. We propose a research direction based on ontology construction using knowledge extraction, matching evolution between versions. This paper presents system architecture to manage the lifecycle of the application ontology incorporating building, matching and evolution processes. This solution is integrated in the source ontology since its creation in order to make it possible to evolve and to be versioned.

**Keywords:** ontology lifecycle; ontology building; ontology matching; ontology evolution; application ontology.

---

## 1. INTRODUCTION

With growing business globalization and worldwide collaboration of manufacturing companies, a seamless exchange of products, services and information, within and across enterprises is urgently required. Both vendors and users are making serious efforts to improve enterprise interoperation [1].

Modern organizations are increasingly operating upon distributed and heterogeneous information systems, as they continuously build new autonomous systems, powered by the rapid advancement of information technology. They are facing challenges to integrate heterogeneous applications. The need to integrate heterogeneous applications, both within and across organizations, is indeed becoming pervasive.

Every day, organizations all over the world generate reports, articles, books, emails, and all kind of textual data concerning several topics. The increase of the storage capacity of computers and servers enable these organizations to keep all files. They produce without the need of deleting anything. One mainly problem they face is to know what kind of information they have, and how it is related.

The fundamental aspect of information exchange among applications, systems, and services is the development of a consistent and comprehensive model for representing the domain knowledge [2]. It is essential for sharing information among independent organizations, and exchange information among heterogeneous applications of Information Systems. To make this possible, we need to model the domain knowledge while preserving its semantics [3]. The development of ontologies is becoming a crucial part of semantic web and knowledge management in the organizations.

Interoperability among different ontologies becomes essential to gain from the power of the Semantic Web. Thus, matching of ontologies becomes a core question.

Ontology matching is a key interoperability enabler for the semantic web, as well as a useful tactic in integration tasks dealing with the semantic heterogeneity problem. It takes the ontologies as input and determines as output a set of correspondences between the semantically related entities of those ontologies.

Ontology matching is seen as a solution provider in today's landscape of ontology research. As the number of ontologies that are made publicly available and accessible on the Web increases steadily, so does the need for applications to use them. A single ontology is no longer enough to support the tasks envisaged by a distributed environment like the Semantic Web. Multiple ontologies need to be accessed from several applications. Matching could provide a common layer from which several ontologies could be accessed and hence could exchange information in semantically sound manners [4].

Thus the use of ontology is increasing in Information Systems, which in response increases the significance of ontology maintenance. Ontologies need to be kept updated for the dependent systems to remain usable. With the increase of changes occurring in the represented domains, ontology evolution becomes a necessary process.

Ontologies are often large and complex structures, whose development and maintenance give rise to certain interesting research problems. For many practical applications, ontologies change over time according to some factors, such as domain changes, adaptations to different applications, and changes to our conceptualisation or understanding of a domain. Support for change management is vital to support distributed ontologies. Preserving consistency, while accommodating new changes, is a crucial task that needs special attention [3]. Also, matching between ontologies are easily affected by changes in the ontologies because a change in one ontology could effects the others.

The paper mainly addresses the problem of cooperating enterprises trying to solve the interoperability problem by introducing ontology based reconciliation solutions. Our purpose focuses on ontology lifecycle for building, matching and evolution ontologies in the enterprise Information Systems.

The goal of this research is to present a system architecture to describe the lifecycle of the application ontology incorporating building, matching and evolution processes. The paper discusses the main features of these processes and their contributions to address the problem of interoperability. The building process is based on knowledge extraction from corpuses and databases to generate the domain ontology. For this purpose, we have developed a set of ontologies intended to capture the semantics for applications integration. The matching process tries to find semantic relationships between entities of ontologies. It takes the ontologies as input and determines as output a set of correspondences to build the matching ontology. Typically, similarity measurement strategies become necessary. In evolution process, the main focus is on keeping ontology and its dependents consistent when changes occur. It includes two sub-processes. The first one is related to the application ontology evolution to guarantee its consistency. The second one concentrates on the matching evolution to highlight consequent effects of ontology evolution on dependent ontologies.

The remainder of this paper is organized as follows. Section 2 presents the building, matching and evolution system architecture. Section 3 sketches out the proposed ontology building process. Then, we present an iterative matching process in section 4. Next, we describe the five (5) major steps of the ontology evolution process and the matching evolution process to highlight consequent effects of ontology evolution on dependent ontology. Just after that, many ideas concerning mapping evolution are mentioned where the

matching evolution process is showed. Section 6 discusses some related work on ontology building as well as ontology evolution. Finally, Section7 provides concluding remarks and sketches some future work.

## 2. BUILDING, MATCHING AND EVOLUTION SYSTEM ARCHITECTURE

EIS (Enterprise Information Systems) is defined as an enterprise application system or an enterprise data source that provides the information infrastructure for an enterprise. An EIS can have many different types including batch applications, traditional applications, client/server applications, web applications, relational databases, and so on. These systems are often materialized in enterprise reality in the form of relational databases, ERP (Enterprise Resource Planning), CRM (Customer Relationship Management), SCM (Supply Chain Management), and legacy systems [5].

The proposed system architecture aims at offering a support for integrating heterogeneous and distributed applications, and accessing multiple ontologies (Figure. 1). It includes building matching and evolution management of application ontologies ensured by three (3) levels respectively.

**Building level:** A company model is a computational representation of the structure including, activities, processes, information, resources, people, behavior, goals and business constraints. The goal is to capture the sets of the enterprise applications, the activities that they perform, the required resources, the manipulated data and the invoked messages. Then, we identify the information flow, their structure and the technical infrastructure to support them for building the application ontologies.

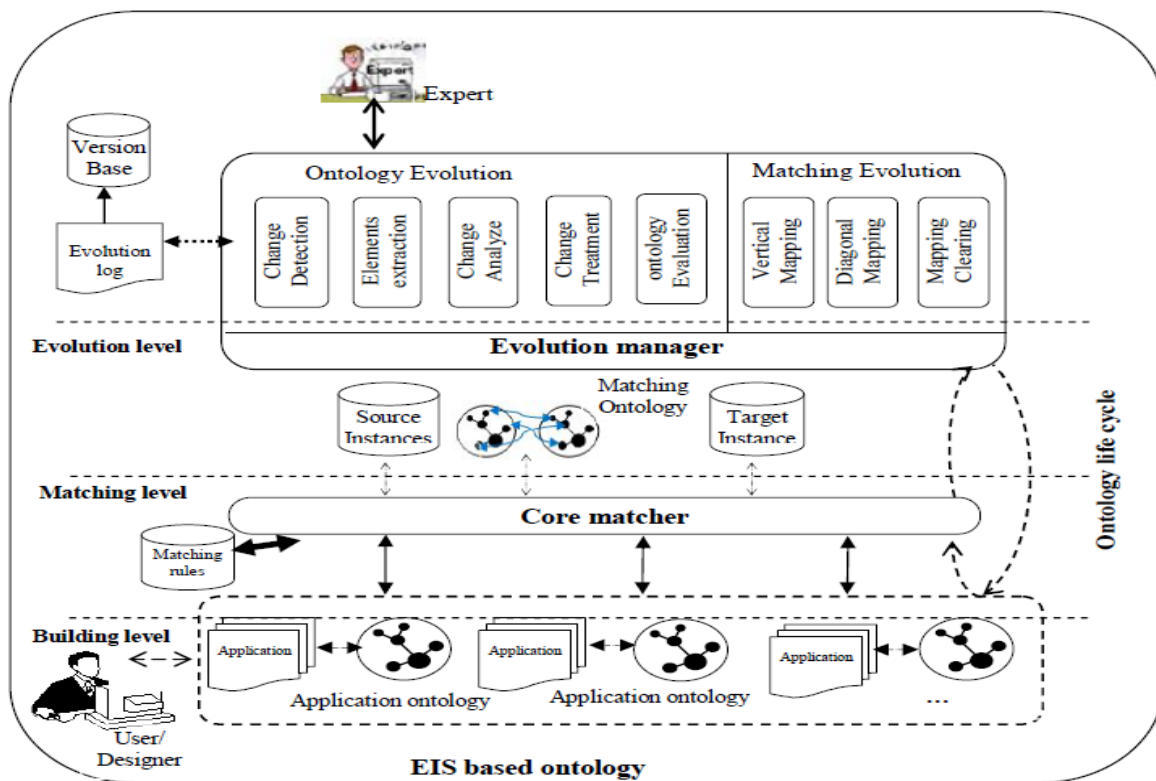


Figure 1. Building, matching and evolution system architecture.

**Matching level:** It concerns the application ontologies integration. The EIS based ontology consists of heterogeneous, autonomous and distributed application. Each application has its own ontology. The application ontologies are related to each other with a matching ontology. We aim at overcoming the gap between application ontologies, according to the semantic relations. A special component, named matcher, is used to perform the tasks of building the matching ontology, and transforming instances of the source ontology into instances of the target ontology.

The Matching Ontology (MO) is formally defined by a 4 tuple:

**MO= ( E, O, M, RT)**

**E:** set of entities such as concepts, relations and attributes.

**O:** set of applications ontologies in the system.

**M:**  $O_s \rightarrow O_t$

Mapping relation between source ontology ( $O_s$ ) and target ontology ( $O_t$ )

**RT:** rules transformation of source instances to target instances.

Additionally, an overview about the matcher component is given in the following. The main task of the matcher is to find semantic relations between concepts of application ontologies. It involves the following tasks:

- Tries to find related concepts or attributes of ontologies and the relations between them. This can be done automatically, semi-automatically or manually with the help of domain experts.
- Represents the identified relations between ontologies based on semantic relations. It combines many algorithms to measure the similarity. Then, it adopts a multi-strategy approach to compute the concepts similarity at various levels, such as lexical, properties (roles and attributes), hierarchical and instances similarities.
- Transforms instances from the source application ontology into instances of the target application ontology by evaluating the equivalence relations defined earlier by the adaptor. Two problems that may arise are that the mappings are incomplete or the that the mapped entities differ in the context. The missing mappings can be gained through inference mechanism.

**Evolution level:** It concerns the evolution management. The evolution manager is composed of two parts, ontology evolution and matching evolution. The first part encompasses the set of activities which ensures that the ontology continues to meet organizational objectives and users' needs in an efficient and effective way. It includes five (5) steps; detection, elements extraction, analysis, treatment of needed change and evaluation. The second part focuses on matching evolution because dynamic environment and applications changes often have consequent effects on dependent ontologies. The role of matching evolution is to detect the new mapping between the old and new versions of the updated ontology.

### 3. BUILDIND PROCESS

Every day, organizations over the entire world generate reports, articles, books, emails, and all kind of textual data concerning several topics. The increase of the storage capacity of computers and servers enable these organizations to keep all files they produce without the need of deleting anything. Although this is an obvious advantage for everybody, it also

implies some problems. One mainly problem they face is to know what kind of information they have, and how it is related. One way to organize information in computer science is in ontology form. This is similar to a conceptual map in which the main topics or concepts are related to each other by some kind of relations [6].

We propose an extraction and building process which includes four main phases [5]: the linguistic study and knowledge extraction, the specification, the ontology conceptualization and formalization and finally, the ontology implementation and validation (cf. figure 2).

The proposed process begins with the linguistic study and the knowledge extraction. It introduces the following steps: corpus pre-processing, extraction of terms, cleaning and filtering, and finally classification. The second phase is the ontology specification phase. It identifies the knowledge domain and the purpose of the ontology, including the operational goal, the intended users and the scope of the ontology which contain the set of terms to be represented and their characteristics. Then, the third phase consists in the conceptualization and formalization. The last phase concerns the ontology implementation and validation test. The arrival of a new corpus of text expresses the need of evolution to maintain our domain ontology. The process of evolution is invoked to guarantee the consistency and the coherence of the ontology and ensure the evolution of the data and/or the domain.

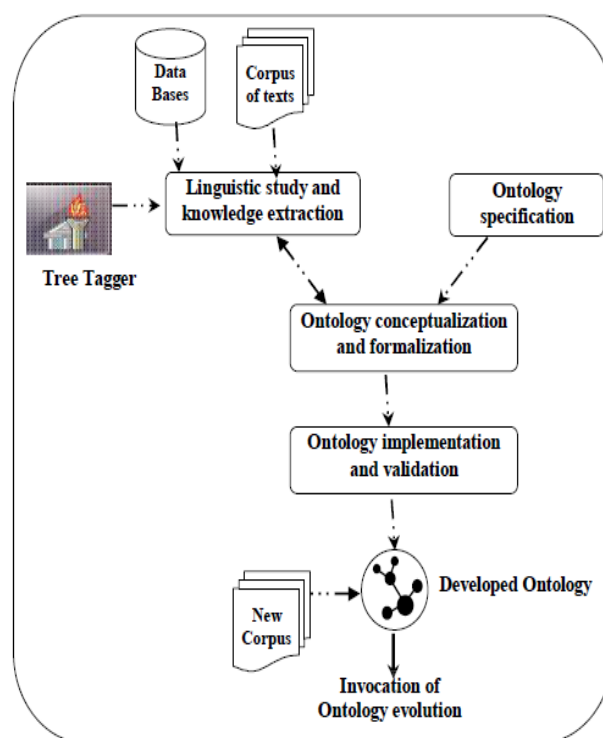


Figure 2. Application ontology building process.

#### 3.1 Linguistic study and knowledge extraction

This phase relies on corpus work and involves the following tasks:

- Corpus pre-processing. It aims to define a strategy to treat the missing data. It consists in normalizing the text to obtain coherent results and also, as possible, to correct human errors by the assistant of linguistic experts. This task serves to normalize the diverse manners of writing the same word, to

correct the obvious spelling mistakes or the typographic incoherence and to clarify certain lexical information expressed implicitly in texts. The textual or linguistic analysis of the corpus means systematizing and making more effective the search of terms in the texts. We also used the spellchecker to avoid errors in the corpus. Then, the text is divided into a set of sentences to allow the use of the morphosyntactic analyzer Tree Tagger [7].

b. Extraction of terms and cleaning. It aims at listing all the terms contained in a corpus. To achieve this goal, we use Tree Tagger, version 3.2 [7]. It is a tool of morphosyntactic labelling and lemmatization. It serves to assign to each term in the corpus its morph syntactic category (name, verb, adjective, article, proper noun, pronoun, abbreviation, etc.) and give for each term its lemmatization. As input, corpus of texts must be organized into a set of sentences, and stored them in a file of .txt extension. Tree Tagger is used to classify extracted terms (concepts/relations) using the annotation and lemmatization information [7]. After the mining of text, we perform the cleaning operations, such as remove the stop words, change the upper case characters to lower case and remove the irrelevant and abbreviation terms.

Several measures are usually used to select the candidate terms, we can quote the number of appearances of a term within a corpus, as well as more complex measures such as the mutual information, tf-idf, is still used in statistical distributions methods [8]. The method is based on the syntactic analysis, and uses the grammatical techniques. They put the hypothesis that the grammatical dependences reflect semantic dependences [9].

c. Classification of terms. The terms extracted from the previous step, were then classified into two categories of terms, following this idea, we try to classify the semantic elements extracted according into two categories: the concepts and the relations. Basing on the information provided by the TreeTagger tool, we classify NAME (proper nouns) as concepts and the terms of type (verb) as relations.

### 3.2 Ontology specification

This phase aims at supplying a clear description of the studied problem and at establishing a document of requirements specification. We need to determine why the ontology is being built, and what is its intended uses and final users.

### 3.3 Ontology conceptualization and formalization

The conceptualization step comprises the following tasks:

a. Glossary of terms. It contains the definition of all terms extracted in the previous phase (concepts, instances, attributes, relations). It contains all the terms and linguistic description.

b. Concept taxonomies. The hierarchy of concepts classification shows the organization of the ontology concepts in a hierarchical order which expresses the relations sub-class and super-class.

c. Definition of binary relations diagram. It specifies which concepts are linked by each relation.

d. Concept dictionary. It contains some of the domain concepts, instances of such concepts, class and instance attributes of the concepts, relations whose source is the concept and, optionally, concept synonyms and acronyms

e. Definition of binary relations tables. The binary relations are represented in the form of properties which attach a

concept to another. For each relation, we define: its name, the name of source concept, the name of target concept, the cardinality and the name of the inverse relation if it exists.

f. Definition of the attributes tables. The attributes are properties which take it values in the predefined types (String, Integer, Boolean ...). For each attribute appearing in the concepts dictionary, we specify its name, the type and the domain.

g. Definition of the logic axioms table. We define for each axiom, its description in natural language, the name of the concept to which the axiom refers, attributes used in the axiom and the logic expression.

h. Definition of the instances table. For each instance identified in the concepts dictionary, we specify the instance name, the concept name to belong to it, the attributes and their values.

The formalization step consists of two parts: terminological language TBOX in which concepts and relations are defined; and an assertion language ABOX in which we introduce the instances.

a. TBOX construction: We define here concepts and relations relating to our domain, by using the constructors provided by description logic to give structured descriptions at concepts and relations [10].

b. ABOX construction: The assertion language is dedicated to the description of facts, by specifying the instances (with their classes) and the relations between them.

### 3.4 Ontology implementation and validation

The implementation step involves the representation of the captured concepts and its relationship in a formal language. Protégé OWL [11] is a development environment with functionality for editing classes, slots (properties) and instances. Protégé is highly extensible and customizable. To evaluate correctness and completeness of domain ontology, we use query and visualization provided by PROTÉGÉ OWL. We use the built in query engine for simple query searches and query plug-in to create more sophisticated searches. We also use visualization plug-ins to browse the application ontology and ensure its consistency. The problems of coherence, correctness and completeness are then verified using the RACER inference engine [11].

OOPS is a web application based on Java [12], used by ontology developers during the ontology validation activity [13]. OOPS! scans ontologies looking for potential pitfalls that could lead to modelling errors. We enter the URL pointing the OWL document describing the ontology to be tested. Once the ontology is parsed using the Jena API the model is scanned looking for pitfalls, from those available in the pitfall catalogue. Therefore, the ontology elements involved in potential errors are detected as well as warnings regarding OWL syntax and some modelling suggestions are generated as well as explanations describing the pitfalls [13].

Once the constructed ontology is validated, it is ready to be invoked by users' requests using SWRL language. The Protégé SWRL Editor is an extension to Protégé OWL that permits interactive editing of SWRL rules [14]. It is tightly integrated with Protégé OWL and is primarily accessible through it. When editing rules, we can directly refer to OWL classes, properties, and individuals within an OWL knowledge base.



#### 4. MATCHING PROCESS

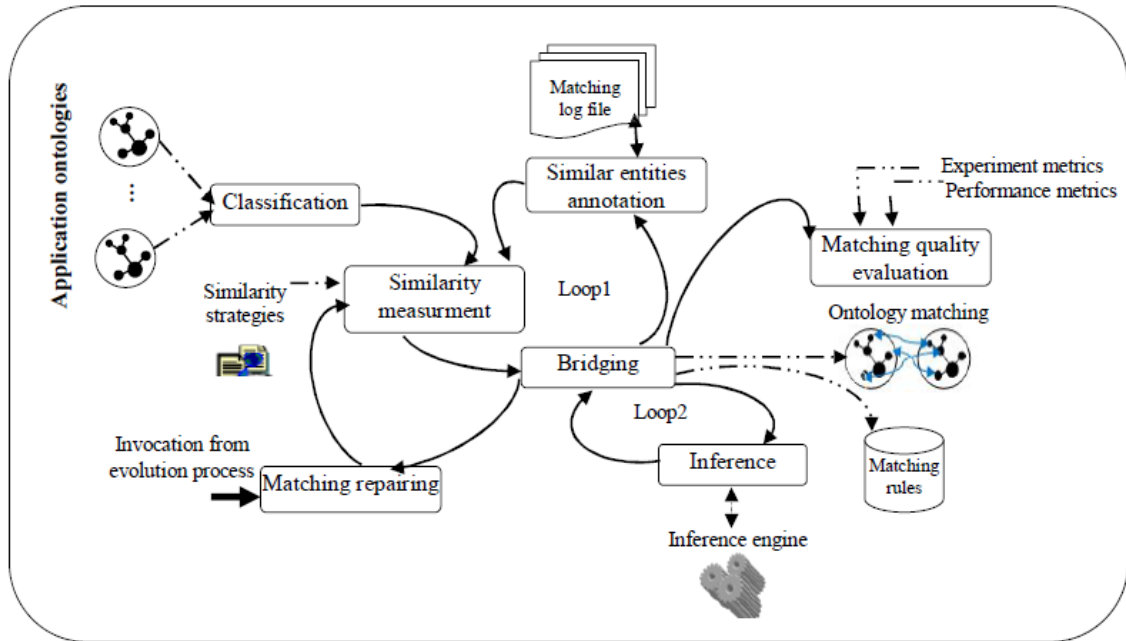


Figure 3. Iterative matching process.

Ontology matching is the process whereby semantic relations are defined between two or more ontologies to align them. It is the set of activities required to transform instances of source ontology into instances of target ontology. In this paper, we propose a matching process which includes seven main steps: classification, similarity measurement, similar entities annotation, bridging, inference, matching quality evaluation and matching repairing. In the proposed process, the classification step tries to filter the ontologies entities in order to obtain candidates entities. It is an iterative process, as described in figure 3, with a primary loop and a secondary loop. At every iteration  $i$ , a semantic bridge is created between entity  $e_i$  of the source ontology and entity  $e_j$  of the target ontology. In the main loop, at every iteration  $i$ , the process executes three steps: it first computes similarity  $\text{sim}(e_i, e_j)$  using similarity strategies, then annotates similar entities, and finally collects similar entities, selecting the most similar entity and defines a bridge. The loop ends when it becomes impossible to create a bridge between entities. The second loop concerns two steps, bridging and inference. It tries to detect new bridges basing on matching rules and human experts. These matching are then used to translate instances of source ontology into instances of target ontology. Finally, the last step focuses on experimental study to deduce some criteria to evaluate matching quality (For more detailed description of this process, see [15]).

#### 5 EVOLUTION PROCESS

Ontology evolution is defined by Haase and Stojanovic, [16][17] as the “timely adaptation of an ontology to the arisen changes and the consistent management of these changes”. Ontology evolution is a process that supports the enrichment of the ontology by adding new entities (concepts, properties, and instances) or by modifying existing entities when new knowledge is acquired.

The usage of ontology is wide spread in Information Systems especially when building a lingua franca for resolving the terminological and conceptual incompatibilities between the

enterprise applications. Ontology evolution takes place when the perspective under which the domain is viewed has changed. More specifically, ontology evolution means modifying or upgrading the ontology when there is a certain need for change as communities of practice concerned with a field of knowledge develop a deeper understanding of the domain. Ontology change management deals with the problem of deciding the modifications to perform in ontology, implementation of these modifications, and the management of their effects in dependent data structures, ontologies, services and applications [18].

One of the crucial tasks faced by practitioners and researchers in the area of knowledge representation is to efficiently encode the human knowledge in ontologies. Maintenance of usually large and dynamic ontologies and in particular adaptation of these ontologies to new knowledge is one of the most challenging problems in the Semantic Web research. This has led to the emergence of several different, but closely related, research areas such as ontology integration, merging, and versioning [3].

In our study, we focus on the ontology evolution and mappings evolution between related ontologies because they stand for the basis of all types of relations between ontologies, such as merging, integration and alignment. For this purpose, we describe two sub-processes. The first one is related to the application ontologies evolution. The second one concentrates on the matching evolution.

##### 5.1 Ontology evolution process

Ontologies are not static entities but evolve over time. We aim in our work to propose an evolution management system to allow evolving, versioning and exploiting application ontologies in dynamic environments. This system helps the designer, user, and expert to supervise the required changes and provides interfaces to participate to the ontology evolution process (cf. Figure. 4).

### 5.1.1 Change detection

An evolution process requires some modifications to occur. It is, thus, necessary to identify the needs of evolution and the compatible changes to apply to the existing ontology. These modifications are expressed informally by different ontology actors (User, expert and ontology designer). The actors can express ambiguous, vague or redundant modifications. These needs will be expressed semi-formally according to one or several types of changes to be applied to create the new version of ontology.

The interview is commonly used to capture the possible changes. It enables the ontology designer to ask periodically questions and allows to the ontology user and experts some freedom to express their answers. The interview includes specific and general questions. The first one concerns the experts to capture specific information. This kind of questions is structured and has the dichotomous (Yes/No) answer. The second one concerns the ontology user to explore an issue or a specific need. This kind of questions is unstructured and its answer is a corpus of text.

Example 1: Expert question

- a- Does change affect the ontology properties?  
 -Yes  
 -No
- b- Does change concern ontology instances?  
 -Yes  
 -No

Example 2: User question

- a- Does change need to adapt functional requirements?
- b- How can the needed change improve the ontology use?

The output of this step is a set of corpus of texts. They enable the ontology designer to capture the needed change(s).

### 5.1.2 Elements extraction

We refer to linguistic study and knowledge extraction phase in the ontology building process to discover the pertinent terms and the type of change(s).

### 5.1.3 Change analysis

To resolve changes, we must identify and represent them in a suitable format. Changes must be formally expressed through types of changes. The composed changes which express a sequence of several elementary changes forming only one logical entity together [17].

### 5.1.4 Change treatment

During this step, it is necessary to determine the direct and indirect types of changes to be applied. In case of ambiguity or in presence of several possibilities, the ontology actors (user, expert and designer) decide on the action to occur.

All changes, and derived ones, confirmed by the designer are applied to the ontology. Consequently, the changes are physically applied to the ontology. The implemented changes need to be propagated to all interested parts in the ontology.

### 5.1.5 Ontology evaluation

It is essential to verify the consistency of the ontology in relation to the semantics of the ontology changes. At the end of the evolution process, a new ontology version is created. At this level, we decide whether to preserve the old version of ontology in the version base or not. The last task in this step is to keep track of the performed changes in the evolution log. The latter records the history of applied ontology changes as an order sequence of information.

A change in one application ontology in the system could have extensive effects on other related ontologies. This is especially important when ontologies are used as basis for semantic integration of enterprise applications. To handle this problem, we have proposed a matching evolution process that defines ontologies versions mappings and new mappings.

In order to avoid performing undesired changes, before applying a change to ontology, a temporary version of the ontology is created to support the change activities.

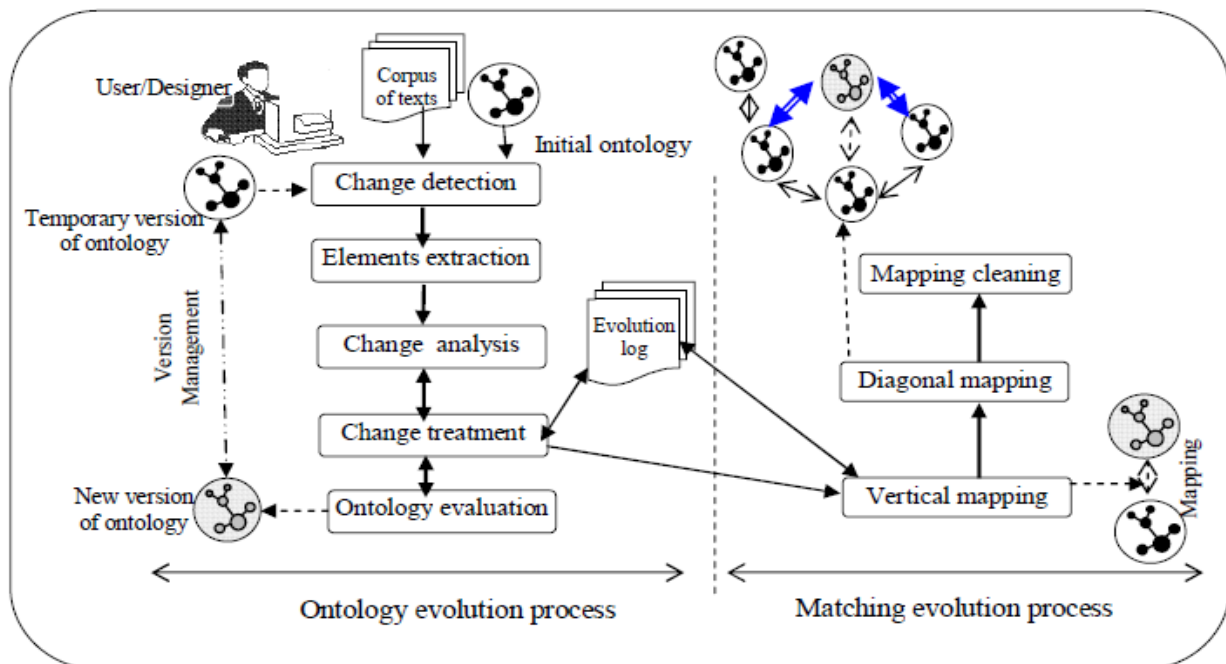


Figure 4. Ontology and matching evolution process.

It enables the ontology engineer to accept or reject the suggested changes and eliminate the changes that can cause ontology inconsistency. Moreover, ontology designer should check the results of a change request on the temporary version to ensure consistency. At the end, the designer can perform successively all the changes on the concerned ontology.

## 5.2 Matching evolution process

Multi-ontology means the existence of multiple ontologies related to each other in many ways [19]; reuse fusion, alignment and integration, to adapt to the various tasks of the EIS. These ontologies must be accessible by different applications and must even exchange semantic information. That is achieved through the mapping of ontologies which is necessary for the management of multiple and heterogeneous ontologies. It is due to its capacity to provide a common layer allowing the access to ontologies and the semantic exchange of information. The problem is how to manage the ontology versions in the system when the ontology evolves?

Additionally, our work is articulated around the ontological mappings evolution after an ontology evolution. Therefore, we define the three following types of mappings:

-The horizontal mappings are the set of existing mappings between the old version of evolved ontology and related ontologies.

-The vertical mappings are mappings between the old version of the evolved ontology and its new version.

-The diagonal mappings are those mappings that exist between the new version of the evolved ontology and all related ontologies. These diagonal mappings are new mappings that are generated when ontology evolves. Therefore, the detection of these diagonal mappings in an automatic way constitutes the principal objective of our work. The diagonal mapping is the composition of horizontal mapping and the vertical mapping.

We have proposed a matching evolution process composed of three (3) steps. The first one is the detection of the vertical mappings between the evolved ontology versions (old, new). For that reason, we studied the effects and the correspondences derived from the application of the change operations. Then, the diagonal mappings are obtained by composition of vertical mappings with the horizontal ones existing between evolved ontology and related ontologies. Finally, we eliminate the invalid and useless correspondences of the obtained mappings.

## 5. DISCUSS AND RELATED WORK

A range of methods and techniques have been reported in the literature regarding ontology building methodologies. We have selected some methodologies whose proposals meet the design criteria mentioned above. Given that ontologies are mainly used in ontological engineering, many of the existing methodologies are geared to the organization and exchange of information in computer systems, as well as in the Semantic Web. Nevertheless, we consider that it is possible to adapt those methodologies to the aims of data and text-mining. Some of them are, for instance, Uschold and King's [20], METHONTOLOGY [21], On-To-Knowledge [22] and Noy and McGuinness' [23]. Other methodologies arose from the work by terminology researchers interested in taking advantage of the features of ontologies for extracting

knowledge from local resources. The most relevant is TERMINAE [24]

Based on these results, METHONTOLOGY meets the most criteria, with the exception of corpus based knowledge extraction. TERMINAE also complies with all the requirements. Therefore, we propose to create a methodology that combines the best characteristics of METHONTOLOGY, on one side, and of TERMINAE on the other. The proposed process is completely suitable to the domain of ontological engineering and knowledge extraction from corpora and databases.

According to Stojanovic [17], "Ontology Evolution is the timely adaptation of ontology to the arisen changes and the consistent propagation of these changes to dependent artefacts (i.e. Dependent ontologies, ontology instances, applications using ontology)". Ontology evolution is a complex process, due to the variety of sources and consequences of changes. Ontology evolution requires taking into account the effects of each change on the ontology to ensure uniformity in the basic ontology and all dependent objects.

Research on ontology evolution is being carried out by different researcher's groups, and their approaches overlap with each other. The current state of the art can be found in [3], [16]. While some of these tools are ontology editors, others provide more specialized features to the user, like the support for evolution strategies, collaborative edits, change propagation, transactional properties, intuitive graphical interfaces, undo/redo operations etc.

Despite these features, our work focuses on the ontology change and matching evolution between related ontologies. Furthermore, we propose two processes. The first one is related to the application ontologies evolution. The second one concentrates on the matching evolution. We have also address the problem of undo/redo operations by using temporary version of the concerned ontology.

## 6. CONCLUSION

Semantic interoperability among applications, systems, and services are mostly based on ontology. A solution is to use an ontology based approach associated to enterprise applications. It provides a semantic layer to encapsulate the applications' heterogeneity. In this paper, we have outlined architecture for application ontologies lifecycle for building, matching and evolution management.

The goal of this research study is to extract knowledge by mining corpus of text to build application ontology. This article deals with knowledge extraction using a text mining approach. More precisely, we concentrate on the extraction and construction process which includes four main phases: the specification, the linguistic study and knowledge extraction, the ontology conceptualization and finally, the implementation of the developed ontology. We use also tools of terminological extraction such as Tree Tagger for the morpho-syntactic labelling and Protégé OWL for the implementation of the ontology.

We have also developed an evolution management system to allow evolving, versioning and exploiting application ontologies in dynamic environments. This system allows the designer, the user and the expert to supervise the required changes, and provides interfaces to participate to the ontology evolution process.

For the future, we identified a number of open issues, we to address in future work. We will improve tool support in the

building process by investigating ways of automatic ontology extraction from data base schema. Particularly interesting is the question of how to combine a top-down modeling approach (the way humans think) with a bottom-up approach (which results from automatic ontology extraction). Furthermore, we intend to integrate our matching tool with (semi-) automatically generated data dictionaries, in order to help domain and/or modeling experts faster understand foreign domains, during the matching process.

## 7. REFERENCES

- [1] Jochem, R. 2010. Enterprise interoperability assessment. In Proceedings of the 8th International Conference of Modeling and Simulation, Hammamet Tunisia. 10-12,
- [2] T. R. Gruber, “Towards principles for the design of ontologies used for knowledge sharing”, *International Journal of Human-Computer studies*, Vol 43, 907-928, 1995.
- [3] A. M. Khattak, K. Latif, and S. Y. Lee, "Change management in evolving web ontologies", *Knowledge based Systems*, (SCI IF: 1.574), ISSN: 0950-707051, 2012.
- [4] Hong-Hai. Do. Melnik. S. Erhard. R. 2002. Comparison of Schema Matching Evaluations. In Proceedings of International Workshop on Web Databases, German, Informatics Society.
- [5] Driouche, R. Bensassi, H. Kemcha, N. 2015. Domain ontology building process based on text Mining from medical structured corpus. In Proceedings of the International Conference on Digital Information Processing, Data Mining and Wireless Communications, Dubai.
- [6] J. I. Toledo-Alvarado, A. Guzmán-Arenas, G. L. Martínez-Luna, “Automatic building of an ontology from a corpus of text documents using data mining tools “ *Journal of Applied Research and Technology*, Vol. 10(3), 398-404, 2012.
- [7]<http://www.cis.unimuenchen.de/~schmid/tools/TreeTagger/>
- [8] Winkler. W. E. 1990. String comparator metrics and enhanced decision rules in the Fellegi-Sunter model of record linkage. In Proceedings of the Section on Survey Research Methods (American Statistical Association). 354–359.
- [9] H. Luong, S. Gauch, Q. Wang, and A. Maglia, “An ontology learning framework using focused crawler and text mining”. *International Journal on Advances in Life Sciences*, Vol 1(23), 99-109, 2009.
- [10] F. Baader, D. Calvanese, D. L. McGuinness, D. Nardi, P. F. Patel-Schneider. 2003. *The description logic handbook: Theory, implementation and applications*. Cambridge University Press, Cambridge, UK.
- [11] <http://protege.stanford.edu/>
- [12]<http://www.oracle.com/technetwork/java/javae/overview/index.html>
- [13] M. Poveda. M.C. Suárez-Figueroa. A. A. Gómez-Pérez. 2010. Double classification of common pitfalls in ontologies. In Proceedings of the Workshop on Ontology Quality at the 17th International Conference on Knowledge Engineering and Knowledge Management. 1-12. Lisbon, Portugal.
- [14] Golbreich, C. Imai, A. 2004. Combining SWRL rules and OWL ontologies with Protégé OWL Plugin, Jess, and Racer. In Proceedings of the 7th International Protégé Conference, Bethesda, MD.
- [15] Driouche. R. 2012. A Mapping process for semantic integration of enterprise applications. In Proceedings of the 3rd international Arab conference on e-technology, Zarqa University, Jordan, 100-107.
- [16] Haase, P. and Sure. Y. 2004. State of the art on ontology evolution. SEKT Deliverable.
- [17] Stojanovic, L 2004 Methods and tools for ontology evolution. Doctoral thesis, University of Karlsruhe.
- [18] G. Flouris, D. Manakanatas, H. Kondylakis, D. Plexousakis, G. Antoniou. “Ontology change: classification and survey”. *Knowledge Engineering Review*, Vol 23(2), 117-152, 2008.
- [19] N. Choi, I. Song, H. Han, “A survey on ontology mapping”, *ACM SIGMOD Record*, Vol 35 (3), 2006, 34-41.
- [20] Uschold, M. and King, M. 1995. Towards a methodology for building ontologies. In Proceedings of the Workshop on Basic Ontological Issues in Knowledge Sharing, D. Skuce (Ed.), Montreal, Canada.
- [21] Gómez-Pérez, A., Fernández-López, M., and Corcho, Ó. 2004. *Ontological engineering: with examples from the areas of knowledge management*. London: Springer Verlag. Greenwood, Edition. *Metodología de la investigación social*. Buenos Aires, Argentina: Paidós.
- [22] S. Staab, H. P. Schnurr, R. Studer, and Y. “Sure, knowledge processes and ontologies”. *IEEE Intelligent Systems*, Vol. 16 (1), 26-34, 2001.
- [23] Noy, N. F. and McGuinness, D. L. 2001. *Ontology development 101: A guide to creating your first ontology* [online]. Technical Report Stanford Knowledge Systems Laboratory.  
<http://www.ksl.stanford.edu/people/dlm/papers/ontology101/ontology101-noymcguinness.html>.
- [24] Aussenac-Gilles, N. Després, S. and Szulman, S. 2008. The TERMINAE method and platform for ontology engineering from texts. In Proceedings of the Conference on Ontology Learning and Population: Bridging the Gap between Text and Knowledge. Amsterdam: IOS Press.

# Towards Ontology Lifecycle: Building, Matching and Evolution to Semantically Integrate Application Ontologies

Razika Driouche  
National High College  
of Biotechnology, Taoufik  
Khaznadar, Algeria

---

**Abstract:** Semantic interoperability among applications, systems, and services are mostly based on ontology. Its increase usage in Information Systems and knowledge sharing systems raises the importance of ontology development and maintenance. It is essential for sharing information among independent organizations, exchange of information among heterogeneous systems. To make this possible, we need to carefully model the domain knowledge while preserving its semantics. Ontologies are complex in nature and often structured. Their development and maintenance incorporate research areas like: building, evolution, versioning, matching and integration where these are fundamentally different. We uncover the gap in the current research area of ontology building, matching and evolution. We propose a research direction based on ontology construction using knowledge extraction, matching evolution between versions. This paper presents system architecture to manage the lifecycle of the application ontology incorporating building, matching and evolution processes. This solution is integrated in the source ontology since its creation in order to make it possible to evolve and to be versioned.

**Keywords:** ontology lifecycle; ontology building; ontology matching; ontology evolution; application ontology.

---

## 1. INTRODUCTION

With growing business globalization and worldwide collaboration of manufacturing companies, a seamless exchange of products, services and information, within and across enterprises is urgently required. Both vendors and users are making serious efforts to improve enterprise interoperability [1].

Modern organizations are increasingly operating upon distributed and heterogeneous information systems, as they continuously build new autonomous systems, powered by the rapid advancement of information technology. They are facing challenges to integrate heterogeneous applications. The need to integrate heterogeneous applications, both within and across organizations, is indeed becoming pervasive.

Every day, organizations all over the world generate reports, articles, books, emails, and all kind of textual data concerning several topics. The increase of the storage capacity of computers and servers enable these organizations to keep all files. They produce without the need of deleting anything. One main problem they face is to know what kind of information they have, and how it is related.

The fundamental aspect of information exchange among applications, systems, and services is the development of a consistent and comprehensive model for representing the domain knowledge [2]. It is essential for sharing information among independent organizations, and exchange information among heterogeneous applications of Information Systems. To make this possible, we need to model the domain knowledge while preserving its semantics [3]. The development of ontologies is becoming a crucial part of semantic web and knowledge management in the organizations.

Interoperability among different ontologies becomes essential to gain from the power of the Semantic Web. Thus, matching of ontologies becomes a core question.

Ontology matching is a key interoperability enabler for the semantic web, as well as a useful tactic in integration tasks dealing with the semantic heterogeneity problem. It takes the ontologies as input and determines as output a set of correspondences between the semantically related entities of those ontologies.

Ontology matching is seen as a solution provider in today's landscape of ontology research. As the number of ontologies that are made publicly available and accessible on the Web increases steadily, so does the need for applications to use them. A single ontology is no longer enough to support the tasks envisaged by a distributed environment like the Semantic Web. Multiple ontologies need to be accessed from several applications. Matching could provide a common layer from which several ontologies could be accessed and hence could exchange information in semantically sound manners [4].

Thus the use of ontology is increasing in Information Systems, which in response increases the significance of ontology maintenance. Ontologies need to be kept updated for the dependent systems to remain usable. With the increase of changes occurring in the represented domains, ontology evolution becomes a necessary process.

Ontologies are often large and complex structures, whose development and maintenance give rise to certain interesting research problems. For many practical applications, ontologies change over time according to some factors, such as domain changes, adaptations to different applications, and changes to our conceptualisation or understanding of a domain. Support for change management is vital to support distributed ontologies. Preserving consistency, while accommodating new changes, is a crucial task that needs special attention [3]. Also, matching between ontologies are easily affected by changes in the ontologies because a change in one ontology could effects the others.

The paper mainly addresses the problem of cooperating enterprises trying to solve the interoperability problem by introducing ontology based reconciliation solutions. Our purpose focuses on ontology lifecycle for building, matching and evolution ontologies in the enterprise Information Systems.

The goal of this research is to present a system architecture to describe the lifecycle of the application ontology incorporating building, matching and evolution processes. The paper discusses the main features of these processes and their contributions to address the problem of interoperability. The building process is based on knowledge extraction from corpuses and databases to generate the domain ontology. For this purpose, we have developed a set of ontologies intended to capture the semantics for applications integration. The matching process tries to find semantic relationships between entities of ontologies. It takes the ontologies as input and determines as output a set of correspondences to build the matching ontology. Typically, similarity measurement strategies become necessary. In evolution process, the main focus is on keeping ontology and its dependents consistent when changes occur. It includes two sub-processes. The first one is related to the application ontology evolution to guarantee its consistency. The second one concentrates on the matching evolution to highlight consequent effects of ontology evolution on dependent ontologies.

The remainder of this paper is organized as follows. Section 2 presents the building, matching and evolution system architecture. Section 3 sketches out the proposed ontology building process. Then, we present an iterative matching process in section 4. Next, we describe the five (5) major steps of the ontology evolution process and the matching evolution process to highlight consequent effects of ontology evolution on dependent ontology. Just after that, many ideas concerning mapping evolution are mentioned where the

matching evolution process is showed. Section 6 discusses some related work on ontology building as well as ontology evolution. Finally, Section7 provides concluding remarks and sketches some future work.

## 2. BUILDING, MATCHING AND EVOLUTION SYSTEM ARCHITECTURE

EIS (Enterprise Information Systems) is defined as an enterprise application system or an enterprise data source that provides the information infrastructure for an enterprise. An EIS can have many different types including batch applications, traditional applications, client/server applications, web applications, relational databases, and so on. These systems are often materialized in enterprise reality in the form of relational databases, ERP (Enterprise Resource Planning), CRM (Customer Relationship Management), SCM (Supply Chain Management), and legacy systems [5].

The proposed system architecture aims at offering a support for integrating heterogeneous and distributed applications, and accessing multiple ontologies (Figure. 1). It includes building matching and evolution management of application ontologies ensured by three (3) levels respectively.

**Building level:** A company model is a computational representation of the structure including, activities, processes, information, resources, people, behavior, goals and business constraints. The goal is to capture the sets of the enterprise applications, the activities that they perform, the required resources, the manipulated data and the invoked messages. Then, we identify the information flow, their structure and the technical infrastructure to support them for building the application ontologies.

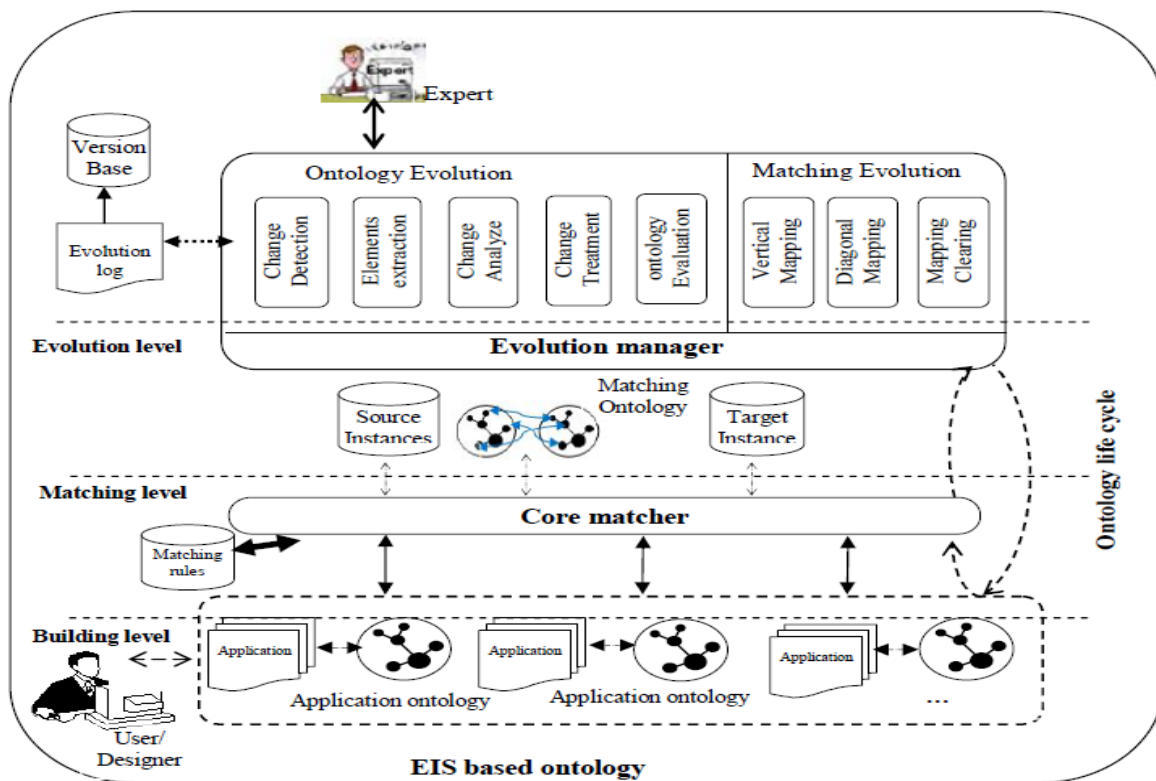


Figure 1. Building, matching and evolution system architecture.

**Matching level:** It concerns the application ontologies integration. The EIS based ontology consists of heterogeneous, autonomous and distributed application. Each application has its own ontology. The application ontologies are related to each other with a matching ontology. We aim at overcoming the gap between application ontologies, according to the semantic relations. A special component, named matcher, is used to perform the tasks of building the matching ontology, and transforming instances of the source ontology into instances of the target ontology.

The Matching Ontology (MO) is formally defined by a 4 tuple:

**MO= ( E, O, M, RT)**

**E:** set of entities such as concepts, relations and attributes.

**O:** set of applications ontologies in the system.

**M:**  $O_s \rightarrow O_t$

Mapping relation between source ontology ( $O_s$ ) and target ontology ( $O_t$ )

**RT:** rules transformation of source instances to target instances.

Additionally, an overview about the matcher component is given in the following. The main task of the matcher is to find semantic relations between concepts of application ontologies. It involves the following tasks:

- Tries to find related concepts or attributes of ontologies and the relations between them. This can be done automatically, semi-automatically or manually with the help of domain experts.
- Represents the identified relations between ontologies based on semantic relations. It combines many algorithms to measure the similarity. Then, it adopts a multi-strategy approach to compute the concepts similarity at various levels, such as lexical, properties (roles and attributes), hierarchical and instances similarities.
- Transforms instances from the source application ontology into instances of the target application ontology by evaluating the equivalence relations defined earlier by the adaptor. Two problems that may arise are that the mappings are incomplete or the that the mapped entities differ in the context. The missing mappings can be gained through inference mechanism.

**Evolution level:** It concerns the evolution management. The evolution manager is composed of two parts, ontology evolution and matching evolution. The first part encompasses the set of activities which ensures that the ontology continues to meet organizational objectives and users' needs in an efficient and effective way. It includes five (5) steps; detection, elements extraction, analysis, treatment of needed change and evaluation. The second part focuses on matching evolution because dynamic environment and applications changes often have consequent effects on dependent ontologies. The role of matching evolution is to detect the new mapping between the old and new versions of the updated ontology.

### 3. BUILDIND PROCESS

Every day, organizations over the entire world generate reports, articles, books, emails, and all kind of textual data concerning several topics. The increase of the storage capacity of computers and servers enable these organizations to keep all files they produce without the need of deleting anything. Although this is an obvious advantage for everybody, it also

implies some problems. One mainly problem they face is to know what kind of information they have, and how it is related. One way to organize information in computer science is in ontology form. This is similar to a conceptual map in which the main topics or concepts are related to each other by some kind of relations [6].

We propose an extraction and building process which includes four main phases [5]: the linguistic study and knowledge extraction, the specification, the ontology conceptualization and formalization and finally, the ontology implementation and validation (cf. figure 2).

The proposed process begins with the linguistic study and the knowledge extraction. It introduces the following steps: corpus pre-processing, extraction of terms, cleaning and filtering, and finally classification. The second phase is the ontology specification phase. It identifies the knowledge domain and the purpose of the ontology, including the operational goal, the intended users and the scope of the ontology which contain the set of terms to be represented and their characteristics. Then, the third phase consists in the conceptualization and formalization. The last phase concerns the ontology implementation and validation test. The arrival of a new corpus of text expresses the need of evolution to maintain our domain ontology. The process of evolution is invoked to guarantee the consistency and the coherence of the ontology and ensure the evolution of the data and/or the domain.

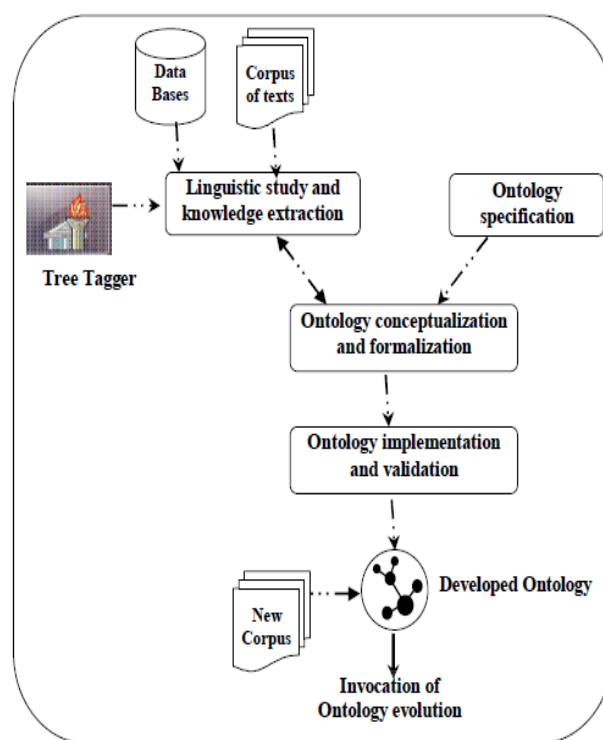


Figure 2. Application ontology building process.

#### 3.1 Linguistic study and knowledge extraction

This phase relies on corpus work and involves the following tasks:

- Corpus pre-processing. It aims to define a strategy to treat the missing data. It consists in normalizing the text to obtain coherent results and also, as possible, to correct human errors by the assistant of linguistic experts. This task serves to normalize the diverse manners of writing the same word, to

correct the obvious spelling mistakes or the typographic incoherence and to clarify certain lexical information expressed implicitly in texts. The textual or linguistic analysis of the corpus means systematizing and making more effective the search of terms in the texts. We also used the spellchecker to avoid errors in the corpus. Then, the text is divided into a set of sentences to allow the use of the morphosyntactic analyzer Tree Tagger [7].

b. Extraction of terms and cleaning. It aims at listing all the terms contained in a corpus. To achieve this goal, we use Tree Tagger, version 3.2 [7]. It is a tool of morphosyntactic labelling and lemmatization. It serves to assign to each term in the corpus its morph syntactic category (name, verb, adjective, article, proper noun, pronoun, abbreviation, etc.) and give for each term its lemmatization. As input, corpus of texts must be organized into a set of sentences, and stored them in a file of .txt extension. Tree Tagger is used to classify extracted terms (concepts/relations) using the annotation and lemmatization information [7]. After the mining of text, we perform the cleaning operations, such as remove the stop words, change the upper case characters to lower case and remove the irrelevant and abbreviation terms.

Several measures are usually used to select the candidate terms, we can quote the number of appearances of a term within a corpus, as well as more complex measures such as the mutual information, tf-idf, is still used in statistical distributions methods [8]. The method is based on the syntactic analysis, and uses the grammatical techniques. They put the hypothesis that the grammatical dependences reflect semantic dependences [9].

c. Classification of terms. The terms extracted from the previous step, were then classified into two categories of terms, following this idea, we try to classify the semantic elements extracted according into two categories: the concepts and the relations. Basing on the information provided by the TreeTagger tool, we classify NAME (proper nouns) as concepts and the terms of type (verb) as relations.

### 3.2 Ontology specification

This phase aims at supplying a clear description of the studied problem and at establishing a document of requirements specification. We need to determine why the ontology is being built, and what is its intended uses and final users.

### 3.3 Ontology conceptualization and formalization

The conceptualization step comprises the following tasks:

a. Glossary of terms. It contains the definition of all terms extracted in the previous phase (concepts, instances, attributes, relations). It contains all the terms and linguistic description.

b. Concept taxonomies. The hierarchy of concepts classification shows the organization of the ontology concepts in a hierarchical order which expresses the relations sub-class and super-class.

c. Definition of binary relations diagram. It specifies which concepts are linked by each relation.

d. Concept dictionary. It contains some of the domain concepts, instances of such concepts, class and instance attributes of the concepts, relations whose source is the concept and, optionally, concept synonyms and acronyms

e. Definition of binary relations tables. The binary relations are represented in the form of properties which attach a

concept to another. For each relation, we define: its name, the name of source concept, the name of target concept, the cardinality and the name of the inverse relation if it exists.

f. Definition of the attributes tables. The attributes are properties which take it values in the predefined types (String, Integer, Boolean ...). For each attribute appearing in the concepts dictionary, we specify its name, the type and the domain.

g. Definition of the logic axioms table. We define for each axiom, its description in natural language, the name of the concept to which the axiom refers, attributes used in the axiom and the logic expression.

h. Definition of the instances table. For each instance identified in the concepts dictionary, we specify the instance name, the concept name to belong to it, the attributes and their values.

The formalization step consists of two parts: terminological language TBOX in which concepts and relations are defined; and an assertion language ABOX in which we introduce the instances.

a. TBOX construction: We define here concepts and relations relating to our domain, by using the constructors provided by description logic to give structured descriptions at concepts and relations [10].

b. ABOX construction: The assertion language is dedicated to the description of facts, by specifying the instances (with their classes) and the relations between them.

### 3.4 Ontology implementation and validation

The implementation step involves the representation of the captured concepts and its relationship in a formal language. Protégé OWL [11] is a development environment with functionality for editing classes, slots (properties) and instances. Protégé is highly extensible and customizable. To evaluate correctness and completeness of domain ontology, we use query and visualization provided by PROTÉGÉ OWL. We use the built in query engine for simple query searches and query plug-in to create more sophisticated searches. We also use visualization plug-ins to browse the application ontology and ensure its consistency. The problems of coherence, correctness and completeness are then verified using the RACER inference engine [11].

OOPS is a web application based on Java [12], used by ontology developers during the ontology validation activity [13]. OOPS! scans ontologies looking for potential pitfalls that could lead to modelling errors. We enter the URL pointing the OWL document describing the ontology to be tested. Once the ontology is parsed using the Jena API the model is scanned looking for pitfalls, from those available in the pitfall catalogue. Therefore, the ontology elements involved in potential errors are detected as well as warnings regarding OWL syntax and some modelling suggestions are generated as well as explanations describing the pitfalls [13].

Once the constructed ontology is validated, it is ready to be invoked by users' requests using SWRL language. The Protégé SWRL Editor is an extension to Protégé OWL that permits interactive editing of SWRL rules [14]. It is tightly integrated with Protégé OWL and is primarily accessible through it. When editing rules, we can directly refer to OWL classes, properties, and individuals within an OWL knowledge base.



#### 4. MATCHING PROCESS

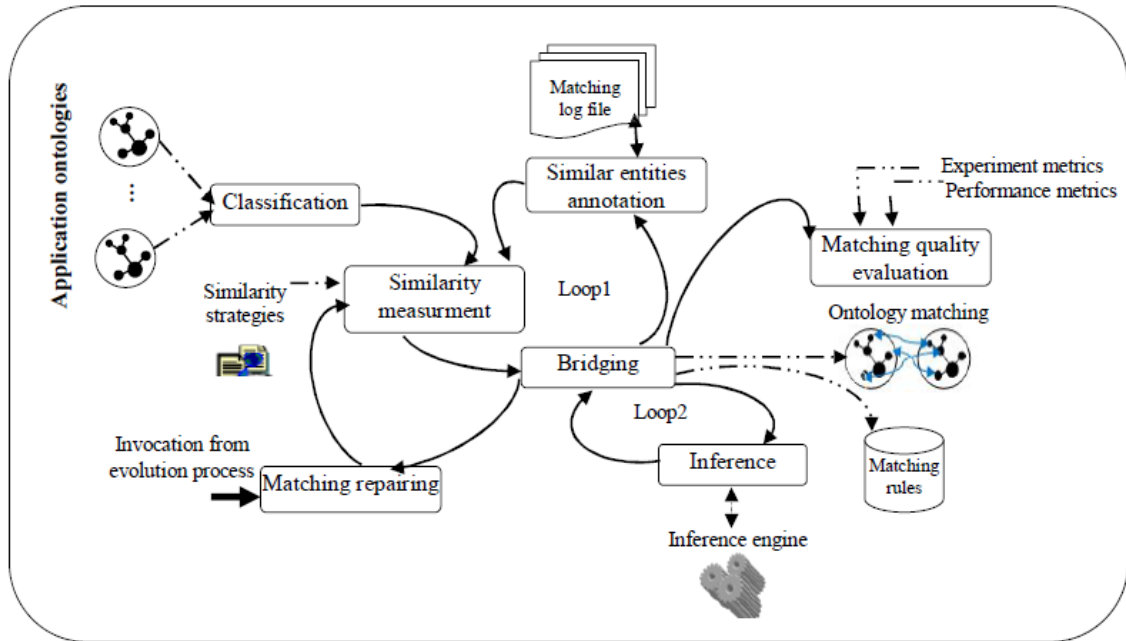


Figure 3. Iterative matching process.

Ontology matching is the process whereby semantic relations are defined between two or more ontologies to align them. It is the set of activities required to transform instances of source ontology into instances of target ontology. In this paper, we propose a matching process which includes seven main steps: classification, similarity measurement, similar entities annotation, bridging, inference, matching quality evaluation and matching repairing. In the proposed process, the classification step tries to filter the ontologies entities in order to obtain candidates entities. It is an iterative process, as described in figure 3, with a primary loop and a secondary loop. At every iteration  $i$ , a semantic bridge is created between entity  $e_i$  of the source ontology and entity  $e_j$  of the target ontology. In the main loop, at every iteration  $i$ , the process executes three steps: it first computes similarity  $sim(e_i, e_j)$  using similarity strategies, then annotates similar entities, and finally collects similar entities, selecting the most similar entity and defines a bridge. The loop ends when it becomes impossible to create a bridge between entities. The second loop concerns two steps, bridging and inference. It tries to detect new bridges basing on matching rules and human experts. These matching are then used to translate instances of source ontology into instances of target ontology. Finally, the last step focuses on experimental study to deduce some criteria to evaluate matching quality (For more detailed description of this process, see [15]).

#### 5 EVOLUTION PROCESS

Ontology evolution is defined by Haase and Stojanovic, [16][17] as the “timely adaptation of an ontology to the arisen changes and the consistent management of these changes”. Ontology evolution is a process that supports the enrichment of the ontology by adding new entities (concepts, properties, and instances) or by modifying existing entities when new knowledge is acquired.

The usage of ontology is wide spread in Information Systems especially when building a lingua franca for resolving the terminological and conceptual incompatibilities between the

enterprise applications. Ontology evolution takes place when the perspective under which the domain is viewed has changed. More specifically, ontology evolution means modifying or upgrading the ontology when there is a certain need for change as communities of practice concerned with a field of knowledge develop a deeper understanding of the domain. Ontology change management deals with the problem of deciding the modifications to perform in ontology, implementation of these modifications, and the management of their effects in dependent data structures, ontologies, services and applications [18].

One of the crucial tasks faced by practitioners and researchers in the area of knowledge representation is to efficiently encode the human knowledge in ontologies. Maintenance of usually large and dynamic ontologies and in particular adaptation of these ontologies to new knowledge is one of the most challenging problems in the Semantic Web research. This has led to the emergence of several different, but closely related, research areas such as ontology integration, merging, and versioning [3].

In our study, we focus on the ontology evolution and mappings evolution between related ontologies because they stand for the basis of all types of relations between ontologies, such as merging, integration and alignment. For this purpose, we describe two sub-processes. The first one is related to the application ontologies evolution. The second one concentrates on the matching evolution.

##### 5.1 Ontology evolution process

Ontologies are not static entities but evolve over time. We aim in our work to propose an evolution management system to allow evolving, versioning and exploiting application ontologies in dynamic environments. This system helps the designer, user, and expert to supervise the required changes and provides interfaces to participate to the ontology evolution process (cf. Figure. 4).

### 5.1.1 Change detection

An evolution process requires some modifications to occur. It is, thus, necessary to identify the needs of evolution and the compatible changes to apply to the existing ontology. These modifications are expressed informally by different ontology actors (User, expert and ontology designer). The actors can express ambiguous, vague or redundant modifications. These needs will be expressed semi-formally according to one or several types of changes to be applied to create the new version of ontology.

The interview is commonly used to capture the possible changes. It enables the ontology designer to ask periodically questions and allows to the ontology user and experts some freedom to express their answers. The interview includes specific and general questions. The first one concerns the experts to capture specific information. This kind of questions is structured and has the dichotomous (Yes/No) answer. The second one concerns the ontology user to explore an issue or a specific need. This kind of questions is unstructured and its answer is a corpus of text.

Example 1: Expert question

- a- Does change affect the ontology properties?  
 -Yes  
 -No
- b- Does change concern ontology instances?  
 -Yes  
 -No

Example 2: User question

- a- Does change need to adapt functional requirements?
- b- How can the needed change improve the ontology use?

The output of this step is a set of corpus of texts. They enable the ontology designer to capture the needed change(s).

### 5.1.2 Elements extraction

We refer to linguistic study and knowledge extraction phase in the ontology building process to discover the pertinent terms and the type of change(s).

### 5.1.3 Change analysis

To resolve changes, we must identify and represent them in a suitable format. Changes must be formally expressed through types of changes. The composed changes which express a sequence of several elementary changes forming only one logical entity together [17].

### 5.1.4 Change treatment

During this step, it is necessary to determine the direct and indirect types of changes to be applied. In case of ambiguity or in presence of several possibilities, the ontology actors (user, expert and designer) decide on the action to occur.

All changes, and derived ones, confirmed by the designer are applied to the ontology. Consequently, the changes are physically applied to the ontology. The implemented changes need to be propagated to all interested parts in the ontology.

### 5.1.5 Ontology evaluation

It is essential to verify the consistency of the ontology in relation to the semantics of the ontology changes. At the end of the evolution process, a new ontology version is created. At this level, we decide whether to preserve the old version of ontology in the version base or not. The last task in this step is to keep track of the performed changes in the evolution log. The latter records the history of applied ontology changes as an order sequence of information.

A change in one application ontology in the system could have extensive effects on other related ontologies. This is especially important when ontologies are used as basis for semantic integration of enterprise applications. To handle this problem, we have proposed a matching evolution process that defines ontologies versions mappings and new mappings.

In order to avoid performing undesired changes, before applying a change to ontology, a temporary version of the ontology is created to support the change activities.

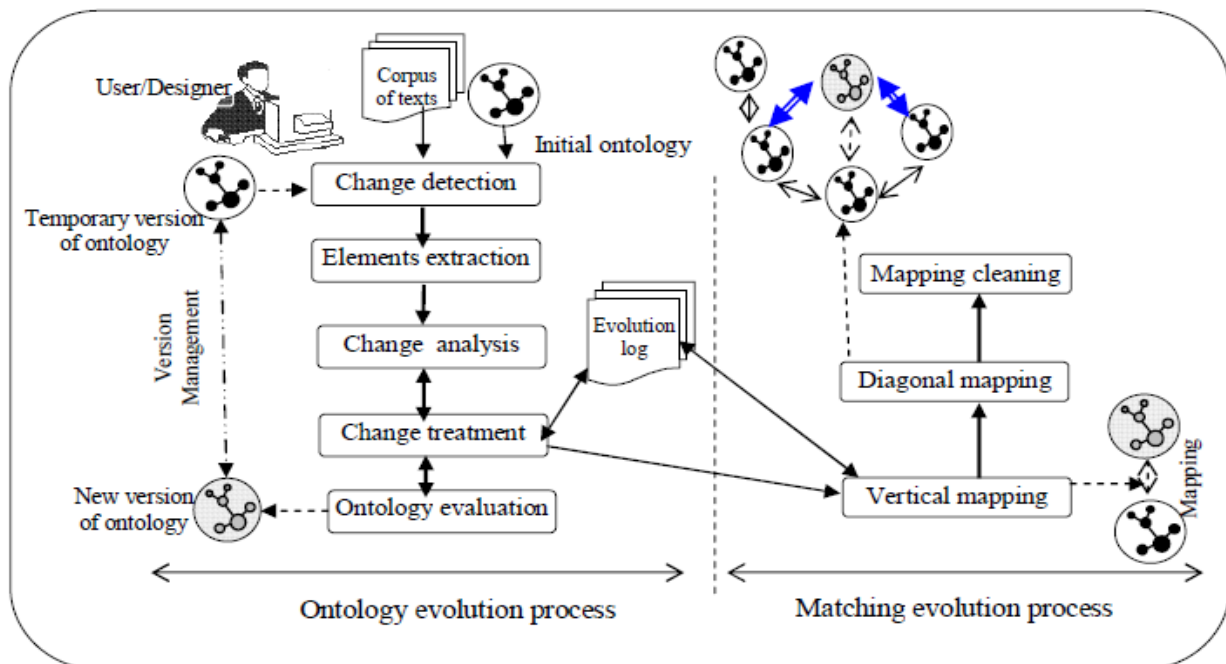


Figure 4. Ontology and matching evolution process.

It enables the ontology engineer to accept or reject the suggested changes and eliminate the changes that can cause ontology inconsistency. Moreover, ontology designer should check the results of a change request on the temporary version to ensure consistency. At the end, the designer can perform successively all the changes on the concerned ontology.

## 5.2 Matching evolution process

Multi-ontology means the existence of multiple ontologies related to each other in many ways [19]; reuse fusion, alignment and integration, to adapt to the various tasks of the EIS. These ontologies must be accessible by different applications and must even exchange semantic information. That is achieved through the mapping of ontologies which is necessary for the management of multiple and heterogeneous ontologies. It is due to its capacity to provide a common layer allowing the access to ontologies and the semantic exchange of information. The problem is how to manage the ontology versions in the system when the ontology evolves?

Additionally, our work is articulated around the ontological mappings evolution after an ontology evolution. Therefore, we define the three following types of mappings:

-The horizontal mappings are the set of existing mappings between the old version of evolved ontology and related ontologies.

-The vertical mappings are mappings between the old version of the evolved ontology and its new version.

-The diagonal mappings are those mappings that exist between the new version of the evolved ontology and all related ontologies. These diagonal mappings are new mappings that are generated when ontology evolves. Therefore, the detection of these diagonal mappings in an automatic way constitutes the principal objective of our work. The diagonal mapping is the composition of horizontal mapping and the vertical mapping.

We have proposed a matching evolution process composed of three (3) steps. The first one is the detection of the vertical mappings between the evolved ontology versions (old, new). For that reason, we studied the effects and the correspondences derived from the application of the change operations. Then, the diagonal mappings are obtained by composition of vertical mappings with the horizontal ones existing between evolved ontology and related ontologies. Finally, we eliminate the invalid and useless correspondences of the obtained mappings.

## 5. DISCUSS AND RELATED WORK

A range of methods and techniques have been reported in the literature regarding ontology building methodologies. We have selected some methodologies whose proposals meet the design criteria mentioned above. Given that ontologies are mainly used in ontological engineering, many of the existing methodologies are geared to the organization and exchange of information in computer systems, as well as in the Semantic Web. Nevertheless, we consider that it is possible to adapt those methodologies to the aims of data and text-mining. Some of them are, for instance, Uschold and King's [20], METHONTOLOGY [21], On-To-Knowledge [22] and Noy and McGuinness' [23]. Other methodologies arose from the work by terminology researchers interested in taking advantage of the features of ontologies for extracting

knowledge from local resources. The most relevant is TERMINAE [24]

Based on these results, METHONTOLOGY meets the most criteria, with the exception of corpus based knowledge extraction. TERMINAE also complies with all the requirements. Therefore, we propose to create a methodology that combines the best characteristics of METHONTOLOGY, on one side, and of TERMINAE on the other. The proposed process is completely suitable to the domain of ontological engineering and knowledge extraction from corpora and databases.

According to Stojanovic [17], "Ontology Evolution is the timely adaptation of ontology to the arisen changes and the consistent propagation of these changes to dependent artefacts (i.e. Dependent ontologies, ontology instances, applications using ontology)". Ontology evolution is a complex process, due to the variety of sources and consequences of changes. Ontology evolution requires taking into account the effects of each change on the ontology to ensure uniformity in the basic ontology and all dependent objects.

Research on ontology evolution is being carried out by different researcher's groups, and their approaches overlap with each other. The current state of the art can be found in [3], [16]. While some of these tools are ontology editors, others provide more specialized features to the user, like the support for evolution strategies, collaborative edits, change propagation, transactional properties, intuitive graphical interfaces, undo/redo operations etc.

Despite these features, our work focuses on the ontology change and matching evolution between related ontologies. Furthermore, we propose two processes. The first one is related to the application ontologies evolution. The second one concentrates on the matching evolution. We have also address the problem of undo/redo operations by using temporary version of the concerned ontology.

## 6. CONCLUSION

Semantic interoperability among applications, systems, and services are mostly based on ontology. A solution is to use an ontology based approach associated to enterprise applications. It provides a semantic layer to encapsulate the applications' heterogeneity. In this paper, we have outlined architecture for application ontologies lifecycle for building, matching and evolution management.

The goal of this research study is to extract knowledge by mining corpus of text to build application ontology. This article deals with knowledge extraction using a text mining approach. More precisely, we concentrate on the extraction and construction process which includes four main phases: the specification, the linguistic study and knowledge extraction, the ontology conceptualization and finally, the implementation of the developed ontology. We use also tools of terminological extraction such as Tree Tagger for the morpho-syntactic labelling and Protégé OWL for the implementation of the ontology.

We have also developed an evolution management system to allow evolving, versioning and exploiting application ontologies in dynamic environments. This system allows the designer, the user and the expert to supervise the required changes, and provides interfaces to participate to the ontology evolution process.

For the future, we identified a number of open issues, we to address in future work. We will improve tool support in the

building process by investigating ways of automatic ontology extraction from data base schema. Particularly interesting is the question of how to combine a top-down modeling approach (the way humans think) with a bottom-up approach (which results from automatic ontology extraction). Furthermore, we intend to integrate our matching tool with (semi-) automatically generated data dictionaries, in order to help domain and/or modeling experts faster understand foreign domains, during the matching process.

## 7. REFERENCES

- [1] Jochem, R. 2010. Enterprise interoperability assessment. In Proceedings of the 8th International Conference of Modeling and Simulation, Hammamet Tunisia. 10-12,
- [2] T. R. Gruber, “Towards principles for the design of ontologies used for knowledge sharing”, International Journal of Human-Computer studies, Vol 43, 907-928, 1995.
- [3] A. M. Khattak, K. Latif, and S. Y. Lee, "Change management in evolving web ontologies", Knowledge based Systems, (SCI IF: 1.574), ISSN: 0950-707051, 2012.
- [4] Hong-Hai. Do. Melnik. S. Erhard. R. 2002. Comparison of Schema Matching Evaluations. In Proceedings of International Workshop on Web Databases, German, Informatics Society.
- [5] Driouche, R. Bensassi, H. Kemcha, N. 2015. Domain ontology building process based on text Mining from medical structured corpus. In Proceedings of the International Conference on Digital Information Processing, Data Mining and Wireless Communications, Dubai.
- [6] J. I. Toledo-Alvarado, A. Guzmán-Arenas, G. L. Martínez-Luna, “Automatic building of an ontology from a corpus of text documents using data mining tools “ Journal of Applied Research and Technology, Vol. 10(3), 398-404, 2012.
- [7]<http://www.cis.unimuenchen.de/~schmid/tools/TreeTagger/>
- [8] Winkler. W. E. 1990. String comparator metrics and enhanced decision rules in the Fellegi-Sunter model of record linkage. In Proceedings of the Section on Survey Research Methods (American Statistical Association). 354–359.
- [9] H. Luong, S. Gauch, Q. Wang, and A. Maglia, “An ontology learning framework using focused crawler and text mining”. International Journal on Advances in Life Sciences, Vol 1(23), 99-109, 2009.
- [10] F. Baader, D. Calvanese, D. L. McGuinness, D. Nardi, P. F. Patel-Schneider. 2003. The description logic handbook: Theory, implementation and applications. Cambridge University Press, Cambridge, UK.
- [11] <http://protege.stanford.edu/>
- [12]<http://www.oracle.com/technetwork/java/javae/overview/index.html>
- [13] M. Poveda. M.C. Suárez-Figueroa. A. A. Gómez-Pérez. 2010. Double classification of common pitfalls in ontologies. In Proceedings of the Workshop on Ontology Quality at the 17th International Conference on Knowledge Engineering and Knowledge Management. 1-12. Lisbon, Portugal.
- [14] Golbreich, C. Imai, A. 2004. Combining SWRL rules and OWL ontologies with Protégé OWL Plugin, Jess, and Racer. In Proceedings of the 7th International Protégé Conference, Bethesda, MD.
- [15] Driouche. R. 2012. A Mapping process for semantic integration of enterprise applications. In Proceedings of the 3rd international Arab conference on e-technology, Zarqa University, Jordan, 100-107.
- [16] Haase, P. and Sure. Y. 2004. State of the art on ontology evolution. SEKT Deliverable.
- [17] Stojanovic, L 2004 Methods and tools for ontology evolution. Doctoral thesis, University of Karlsruhe.
- [18] G. Flouris, D. Manakanatas, H. Kondylakis, D. Plexousakis, G. Antoniou. “Ontology change: classification and survey”. Knowledge Engineering Review, Vol 23(2), 117-152, 2008.
- [19] N. Choi, I. Song, H. Han, “A survey on ontology mapping”, ACM SIGMOD Record, Vol 35 (3), 2006, 34-41.
- [20] Uschold, M. and King, M. 1995. Towards a methodology for building ontologies. In Proceedings of the Workshop on Basic Ontological Issues in Knowledge Sharing, D. Skuce (Ed.), Montreal, Canada.
- [21] Gómez-Pérez, A., Fernández-López, M., and Corcho, Ó. 2004. Ontological engineering: with examples from the areas of knowledge management. London: Springer Verlag. Greenwood, Edition. Metodología de la investigación social. Buenos Aires, Argentina: Paidós.
- [22] S. Staab, H. P. Schnurr, R. Studer, and Y. “Sure, knowledge processes and ontologies”. IEEE Intelligent Systems, Vol. 16 (1), 26-34, 2001.
- [23] Noy, N. F. and McGuinness, D. L. 2001. Ontology development 101: A guide to creating your first ontology [online]. Technical Report Stanford Knowledge Systems Laboratory.  
<http://www.ksl.stanford.edu/people/dlm/papers/ontology101/ontology101-noymcguinness.html>.
- [24] Aussenac-Gilles, N. Després, S. and Szulman, S. 2008. The TERMINAE method and platform for ontology engineering from texts. In Proceedings of the Conference on Ontology Learning and Population: Bridging the Gap between Text and Knowledge. Amsterdam: IOS Press.

# ASIC Design of Reversible Multiplier Using Adiabatic Technique

Minal Gholpe

Department Electronics and Telecommunication  
G.H.Raisoni Institute of Engineering and  
Technology, Nagpur

Prasad Sangare

Department Electronics and Telecommunication  
G.H.Raisoni Institute of Engineering and  
Technology, Nagpur

---

**Abstract**-From past few decades, VLSI technology has been growing to the large extent. All credit for this goes to the increasing usage of integrated circuits for every embedded system, mobile technologies, computing systems, etc. Growth and use of technology has increased the thirst for low energy or power consumption. An Adiabatic approach is perfect solution for the designing of power and energy efficient designs. The word 'Adiabatic' is the change of state that occurs without the loss or gain of heat. Reversible computing performed on Toffoli gate having adiabatic design techniques promises more reduced in power consumption as compared to traditional adiabatic CMOS circuits. Tanner EDA tool is used for designing the schematic and analysis. S-EDIT is used to design the schematic and T-SPICE is used to Simulate and check the results of Power Dissipation. W-EDIT is used to display the simulation results in the form of waveform.

---

## 1. INTRODUCTION

The electronics industry has achieved a phenomenal growth over the last few decades, mainly due to the rapid advances in integration technologies and large systems design. Use of integrated circuits in high-performance computing, telecommunications, and consumer electronics has been growing at a very fast pace. Advances in device manufacturing technology allow steady reduction of minimum feature size (such as minimum channel length of a transistor or an interconnect width realizable on chip). In 1980, at the beginning of the VLSI era, the typical minimum feature size was 2  $\mu\text{m}$ , and a feature size of 0.3  $\mu\text{m}$  was expected around the year 2000. A minimum feature size of 0.25  $\mu\text{m}$  was achieved by 1995. When we compare integration density of integrated circuits, a clear distinction which is made between the memory chips and logic chips. The number of transistors per chip has continued to increase at an exponential rate over last three decades, effectively confirming "Gordon Moore's" prediction on the growth rate of chip complexity, which was made in the early 1960s (Moore's Law). It has been observed that in terms of transistor count, logic chips which contain significantly fewer transistors in any given year mainly due to large consumption of chip area for complex interconnects. Today we are going through an advanced IC technology. In this we have VLSI technology.

CMOS is referred to as Complementary Metal Oxide Semiconductor, CMOS technology is becoming the mainstream fabrication technology for memories and microcomputers is only because of its high density and low power features. CMOS is a technology for constructing integrated circuits. CMOS technology is used in microprocessors, microcontrollers, static RAM, and various digital logic circuits. CMOS technology is also used for several analog circuits such as image sensors (CMOS sensor), data converters, and highly integrated transreceivers for many types of communication. But some applications such as computer and communication systems require better speed performance than that

obtained by CMOS technology so bipolar LSI's have been used in such fields.

## 2. LITERATURE REVIEW

Since last few decades the main challenges were Area, cost, and performance. But these days power is an important factor instead of cost, performance and area. The device which consumes very less power irrespective of speed such as heart pacemaker, RFID etc. works on the principle of adiabatic logic. The aim of reduction in power consumption is application specific. The authors have tried to decrease the power by combining the adiabatic and reversible technique[1]. The power consumed in traditional CMOS design can be given as,

$$P = CL \cdot VDD^2 \cdot f \quad [6]$$

Here the power (P) is proportional to switching frequency (f), capacitance (CL), and square of supply voltage (VDD). Power consumption can be reduced by minimizing power supply, capacitance and switching frequency of operation. But as soon as these parameters reduces, it may deteriorate the performance of the circuit.

Design using adiabatic principle helps in reducing power consumption at the cost of reduced performance. A method based on adiabatic technique uses an ac power supply rather than dc for energy recovery. Theoretically adiabatic circuits

consume zero power, it shows energy loss due to nonzero resistance in the switches. There are so many papers which describe different types of adiabatic technique such as ECRL, 2PASCL, PFAL etc. by which we can reduce power consumption of the circuit [2]. These technique consume less power as compare to other CMOS circuits

### 3.CONCEPTS

#### I.Adiabatic Circuits

The term “adiabatic” refers to the thermodynamic process that exchanges no energy with environment, and therefore there is no occurrence of power or energy dissipation . During the switching process, adiabatic technology reduces the power or energy dissipation and reuses some part of the energy by recycling it from the load capacitance.

Adiabatic circuits are basically low power circuits which use to conserve the energy by returning back its output energy to input, so that the same energy can be used for next operation.

Fig. 1 and Fig. 2 shows the Charging and Discharging in conventional CMOS circuit and Adiabatic System.

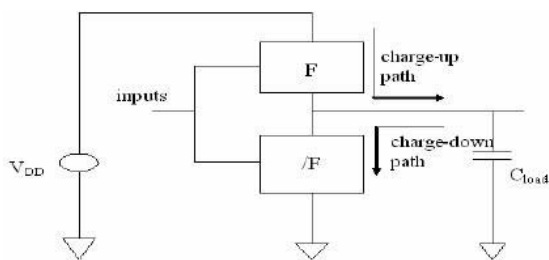


Figure 1. Charging and Discharging in Conventional System

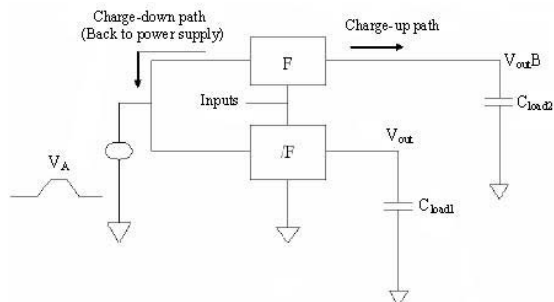


Figure 2. Charging and Discharging in Adiabatic system

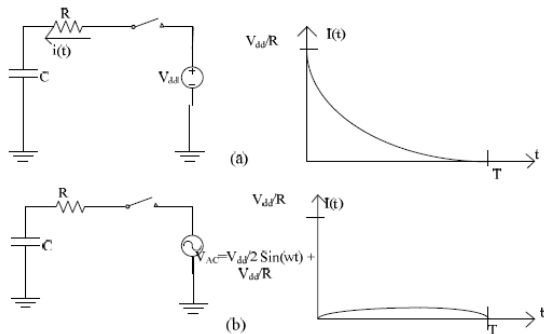


Fig. 3. (a) Switching of CMOS (b) Switching of Adiabatic Logic.

Adiabatic circuits aims to conserve the charges by following essential rules,

- 1) Avoiding turning on of transistor whenever there is a potential difference across the drain and source ( $V_{DS} > 0$ ).
- 2) Avoiding turning off of Transistor whenever there is a flow of current through drain and source. ( $I_{DS} \sim 0$ ).
- 3) The current should not pass through diode.

#### Adiabatic Logic Types

During literature survey, we found different types of adiabatic circuits . They can be grouped into two fundamental kinds:

- Fully Adiabatic Circuit
- Partially energy recovery Adiabatic Circuit (Quasi)

Partially Adiabatic families include the following

- Efficient Charge Recovery Logic
- 2N-2N2P Adiabatic Logic
- Positive Feedback Adiabatic Logic
- NMOS Energy Recovery Logic
- Clocked Adiabatic Logic
- True Single-Phase Adiabatic Logic

### II.Reversible Gates

#### Goals of Reversible Logic:

**A. Quantum Cost:** Quantum cost of a circuit is measure of implementation cost of quantum circuits. More precisely, quantum cost is defined as number of elementary quantum operations needed to realize a gate.

**B. Speed of Computation:** The time delay of circuits should be as low as possible as there are numerous computations which have to be done in a system involving a quantum processor; hence speed of computation is very important parameter while examining such systems.

**C. Garbage Outputs:** Garbage output are those output signals which do not take contribution in driving further blocks in that design. These outputs become redundant as they are not required for computation at later stage. The garbage outputs make the system slower; hence for better efficiency it is very necessary to minimize number of garbage outputs.

**D. Feedback:** Looping is strictly prohibited when we are designing reversible circuits.

**E. Fan-out:** The output of a certain block in the design can only drive at most one block in design. Hence it can be said that the Fan-out is restricted to 1.

Feynman gate (FG) and Toffoli gate (TG), are universal reversible gate. FG shown in Fig.2.1 has QC equals 1 and hardware complexity is  $1\alpha$ . TG shown in Fig. 2.2 has QC equals 5 and hardware complexity is  $1\alpha + 1\beta$ .

#### 2.1 Feynman Gate

It is a  $2 \times 2$  Feynman gate . The input vector is  $I(A, B)$  and output vector  $O(P, Q)$ . The outputs are defined by  $P=A,$

$Q=A\oplus B$ . Quantum cost of a Feynman gate is 1. Figure 2.1 shows a 2\*2 Feynman gate..

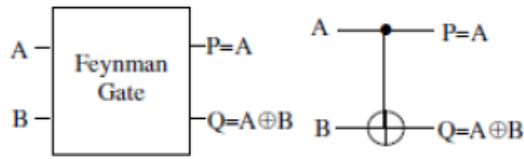


Figure 2.1: Feynman Gate

### 2.2 Toffoli Gate

It is a 3\*3 Toffoli gate . The input vector I(A, B, C) and the output vector O(P,Q,R). The outputs are defined by  $P=A$ ,  $Q=B$ ,  $R=AB\bar{C}$ . Quantum cost of a Toffoli gate is 5. Figure 2.2 shows a 3\*3 Toffoli gate.

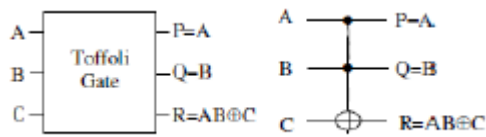


Figure 2.2: Toffoli gate

### Reversible Multiplier

Multiplier circuits are divided into two types : unsigned and signed. Several approaches have been presented to multiply signed numbers, such as 2's complement, Baugh-Wooley, and modified Baugh-Wooley methods . In modified Baugh-Wooley method, number quantity is considered as 2's complement and shows how the multiplication operation takes place, which need AND gates and NAND gates to produce a signed multiplier.

### III.Wallance Tree Multiplication Algorithm

The well-known Wallace high-speed multiplier use carry save adders to reduce an N-row bit product matrix to an equivalent two row matrix that is then summed with carry propagating adder to give product .

The common multiplication method is “add and shift” algorithm. Multiplication algorithm for an N bit multiplicand by N bit multiplier is shown below:

$$\begin{array}{r}
 A = A_3 A_2 A_1 A_0 \text{ Multiplicand} \\
 B = B_3 B_2 B_1 B_0 \\
 \text{Multiplier} \\
 \hline
 Y = Y_7 Y_6 Y_5 Y_4 Y_3 Y_2 Y_1 Y_0 \\
 \text{Multiplication of A and B}
 \end{array}$$

	A <sub>3</sub>	A <sub>2</sub>	A <sub>1</sub>	A <sub>0</sub>	Inputs
	x	B <sub>3</sub>	B <sub>2</sub>	B <sub>1</sub>	B <sub>0</sub>
	C	B <sub>0</sub> x A <sub>3</sub>	B <sub>0</sub> x A <sub>2</sub>	B <sub>0</sub> x A <sub>1</sub>	B <sub>0</sub> x A <sub>0</sub>
+	B <sub>1</sub> x A <sub>3</sub>	B <sub>1</sub> x A <sub>2</sub>	B <sub>1</sub> x A <sub>1</sub>	B <sub>1</sub> x A <sub>0</sub>	
	C	sum	sum	sum	sum
+	B <sub>2</sub> x A <sub>3</sub>	B <sub>2</sub> x A <sub>2</sub>	B <sub>2</sub> x A <sub>1</sub>	B <sub>2</sub> x A <sub>0</sub>	
	C	sum	sum	sum	sum
+	B <sub>3</sub> x A <sub>3</sub>	B <sub>3</sub> x A <sub>2</sub>	B <sub>3</sub> x A <sub>1</sub>	B <sub>3</sub> x A <sub>0</sub>	
	C	sum	sum	sum	sum
	Y <sub>7</sub>	Y <sub>6</sub>	Y <sub>5</sub>	Y <sub>4</sub>	Y <sub>3</sub>
	Y <sub>2</sub>	Y <sub>1</sub>	Y <sub>0</sub>		Outputs

For the Wallace reduction method, once partial product array is formed, adjacent rows are collected into nonoverlapping groups of three. Each group of three rows can be reduced by 1) applying a full adder to each column that contains three bits, 2) applying half adder to each column that contains two bits, and (3) passing any single bit columns to next stage without processing.

This reduction method is applied to each successive stage until two rows remain. The final two rows are summed with a carry propagating adder.

### 4.CONCLUSION

It can be seen that the performance of digital circuits can be enhanced using reversible gates. Adiabatic circuits are low power circuits which use to conserve the energy . Reversible multiplier designed using reversible gates and adiabatic logic families reduces power dissipation and leakage current. Wallace approach will minimize the number of required half adder and full adder which will reduce area of circuit.

### REFERENCES

- [1] Gaurav Kumar , Trailokya Nath Sasamal “Design and Analysis of Toffoli gate using Adiabatic Technique” International Conference on Computing, Communication and Automation (ICCCA2015) ISBN:978-1-4799-8890-7/15/\$31.00 ©2015 IEEE.
- [2] Sakshi Goyal , Gurvinder Singh, Pushpinder Sharma “Variation of Power Dissipation for Adiabatic CMOS and Conventional CMOS Digital Circuits” 2ND INTERNATIONAL CONFERENCE ON ELECTRONICS AND COMMUNICATION SYSTEM (ICECS 2015) 978-1- 4788-7225-8/15/\$31.00 ©2015 IEEE.

- [3] Arpan Chaudhuri, Mamia Saha, Moumita Bhowmik “Implementation of circuit in Different Adiabatic Logic” 2ND INTERNATIONAL CONFERENCE ON ELECTRONIC AND COMMUNICATION SYSTEM(ICECS 2015) 978-1-4788-7225-8/15/\$31.00 ©2015.
- [4] P. Rajashekhar Reddy, S. Raghavendra Swami, S. Ravi Kumar “Implementation of High Speed Low Power Combinational and Sequential Circuits using Reversible logic” International Journal of Electrical, Electronics and Computer Systems (IJEECS) ISSN (Online): 2347- 2820, Volume -3, Issue-12 2015.
- [5] A.A. Hatkar, A.P.Hatkar, N.P. Narkhede “ASIC Design of Reversible Multiplier Circuit International Conference on Electronic Systems, Signal Processing and Computing Technologies 978-1- 4799-2102-7/14 \$31.00 © 2014 IEEE.
- [6] B.Dilli Kumar and M. Bharati “Design of energy efficient Arithmetic circuits using charge recovery Adiabatic logic” International journal of engineering Trends and Technology, Volume 4, Issue 1, pp.31-40, April 2013.
- [7] Kartikeya Bhardwaj, Pravin S. Mane, Jorg Henke “Power- and Area-Efficient Approximate Wallace Tree Multiplier for Error-Resilient Systems” 15th International Symposium on Quality Electronic Design 978-1-4799-3946-6/14/\$31.00 ©2014 IEEE.
- [8] Yogesh M. Motey, Tejaswini G. Panse “Hardware Implementation of Truncate Multiplier Based on Multiplexer Using FPGA” International Conference on communication and Signal processing (ICCSP) ISBN 978-1-4673-4866-9 © 2013.
- [9] Ron S. Waters, Earl E. Swartzlander “A Reduced Complexity Wallace Multiplier Reduction” IEEE TRANSACTIONS ON COMPUTERS, VOL. 59, NO. 8, AUGUST 2010.
- [10] Sung-mo Kang, Yusuf Leblebici “Cmos Digital Integrated Circuit’s Analysis And Design” Tata McGraw-hill Third Edition- 2008.
- [11] Douglas A. Pucknell “Basic Of VLSI Design” Prentice-hall, Inc., 2007.
- [12] Neil H.E. Waste, David Harris., Ayan Banerjee, “Principles Of CMOS VLSI Design- A Circuit And Systems Perspective”, Third Edition, Addison-Wesley, MA, May 2004.