# New Algorithm for Computing Transitive Closure of Data Records Application

Affariza binti Musa #1, Mohamed Faidz Mohamed Said #2

<sup>#</sup> Universiti Teknologi MARA
70300 Seremban, Negeri Sembilan, MALAYSIA
<sup>1</sup> affariza.musal0@gmail.com
<sup>2</sup> faidzms@ieee.org

Abstract— Information quality change has gotten to be definitive issue for some associations and organizations since poor information quality debases authoritative execution though enhanced information quality results in consumer loyalty and cost sparing. The performance such as recognize and removing "duplicate" database records from a single database, and correlating records from different databases that identify the same real world "entity" are used routinely to develop data quality. These calculations dodge the repetitive calculations and high stockpiling cost found in various comparable calculations. Utilizing reproduction, this paper looks at the execution of the new calculations with those found in writing and shows plainly the prevalence of the new calculations. Due to huge databases having several hundred million to several billion records, and continuously growing, efficient techniques and algorithms are needed. This paper is present the new algorithm for computing the transitive closure of large database relations and present the simulation results that show that these direct algorithms perform uniformly better than the best of the iterative algorithms. The alleged "transitive conclusion issue" is a formal detailing of what should be done in step one. This paper presents a record gathering issue called transitive conclusion and proposes calculations to take care of the transitive conclusion issue. The paper additionally writes about the exact investigation of the proposed calculations and remarks on their executions.

Keywords: transitive closure, algorithm, large database

# I. INTRODUCTION

The importance of data quality has been realized by most companies. [15] State that over 80% agreed improving customer data quality was their top priority. Although, the complexity of implementing such task will add the high cost of cleansing data. Depending on company size and the amount of data you need to clean, the reported cost ranges from \$100,000 and \$500, 000 [1].

The large database relations are a similar research effort is required investigation of algorithms for computing transitive closure. There are two categories know in the transitive closure algorithms which are iterative algorithms and direct algorithm [1]

The key for applications in communication systems is the performance of expert database systems. Evaluation of

recursive queries is the main interest in expert database systems. Recursive queries are used efficient evaluation of transitive closure as a paradigm [7].

It is an important step towards the development of the future intelligent database systems by developing efficient algorithms for processing the transitive closure queries [13].

In this paper, a present new algorithm for implementing transitive closure in database corner [1] and demonstrates the analysis tools for achieving data quality in terms of, consistency, data accuracy, redundancy, currency and completeness [15]. In this paper, the excessive rounds of communication are avoided and fast access paths are preserved to present the data of a binary relation by designed a special scheme [7].

The organization of this paper will be as follows. In Section 2, the well-known direct algorithms for computing Transitive Closure is described of Boolean matrix and the calculations could be utilized as a part of the setting of social databases is appeared. Segment 3 is the heart of this paper where we introduce new direct calculations for registering the transitive conclusion of extensive database relations. In Section 4, we propose a technique for assessing the execution of recursive questions and assess the execution of the calculations created in Section 3 against the best iterative calculation. At last, in Section 5, we outline the principle finishes of this study [1].

# II. NAÏVE DIRECT ALGORITHMS

In this section, this paper is briefing a review of some wellknown calculations that were initially foreseen to process the transitive conclusion of Boolean lattics. The new calculations that we propose in Section 3 have been propelled by these calculations. This paper likewise demonstrates for the reasons for comparability how one could acquire in basic adaption of these calculations to process the transitive conclusion of the database relations [1].

# 2.1 The Warshall Algorithm

### Given an initial v x v Boolean matrix of elements a<sub>ij</sub> over a v node graph, with u<sub>ij</sub> being 1 if there is an arc from node i to node j and 0 otherwise, its transitive closure can be obtained as: 1. CLOSURE ← MAT 2. for k = 1 to n 3. for i = 1 to n 4. for j = 1 to n 5 CLOSURE(i; j) ← CLOSURE(i, j) ∨ (CLOSURE(i; k) ∧ CLOSURE(k; j))

This algorithm involves three for loops, with two of them nested. Each loop iterates from 1 to n. This gives us a time complexity of  $O(n^3)$ . If we were to find the transitive closure using the matrix multiplication method we would get a time complexity of  $O(n^4)$ . Each time a matrix multiplication is performed the time complexity is  $O(n^3)$  as there are three loops running (two nested) from 1 to n. The matrix multiplications are carried out a total of n - 1 2 3 n

times to find matrices  $(M_R) \stackrel{\sim}{\Theta}$ ,  $(M_R) \stackrel{\sim}{\Theta}$ ...  $(M_R) \stackrel{\sim}{\Theta}$ , since  $M_{R\infty} = M_R \vee$ 

 $(M_R) \Theta \lor ... \lor (M_R) \Theta$ . So the number of steps involved are  $n^3(n-1)$ , giving us a time complexity of  $O(n^4)$  [5].

## 2.2 The Warren Algorithm

For i-1 to v For k-1 to i-1 For j-1 to v

 $a_{ij} = a_{ij} \lor (a_{ij} \land a_{kj})$ 

The main change is that the i and k circles have been traded. In any case, this trade could bring about a few ways being passed up a major opportunity thus the calculation now requires two "goes" before it finishes. The adjustment in the scope of the second circle record, k, is a streamlining that decreases the expense of two passes [1].

### 2.3 Other Direct Algorithm

This algorithm avoids the redundant computations found in Semi-naive and uses linear data structures to store its intermediate results. To process a query, Jiang's algorithm adopts a mixture of breadth-first search and some form of depth-first search strategies. To illustrate the basic idea behind Jiang's algorithm, consider the computation of the PTCs relevant to the source nodes a and h of the graph in Fig. 1.

The storage of the Partially instantiated Transitive Closure (PTC)s for those nodes in the base graph with more than one incoming edge ensures that a given node in the base graph (and the portion of the graph reachable from that node) will be processed at most once (such a processing is carried out during the computation of the PTC relevant to that node; any future reference to the same node will make use of the corresponding stored PTC), thus eliminating the need to traverse any edge in the base graph more than once [14].



Figure 1. A base graph [14]

## **III. EFFICIENT DIRECT ALGORITHMS**

In this segment, this paper is available a few alterations of the calculations examined in the past segment and propose some non-evident executions that hope to perform much superior to the guileless usage talked about above. Every one of the calculations beneath accept that the underlying connection has been sorted on the fields taking an interest in the transitive conclusion so that every one of the "successors" of a given hub can be found on an adjoining set of tuples in the connection.

This examination expects that the connection whose transitive conclusion is to be processed is huge contrasted with the memory accessible and must be divided into pieces each of which can fit in memory. A parcel will comprise of the successor arrangements of a few hubs. These rundowns develop as the conclusion is registered, and the underlying parcels may no more fit in the memory. To handle this circumstance, the greater part of the calculations beneath depends on element parcels. As the calculation continues, if the memory begins topping off, some successor records are erased or composed pull out of memory to be incorporated into the following segment that is perused in.

## Algorithm 1: Block Warshall

For each column partition

(columns j<sub>b</sub> to j<sub>c</sub> inclusive) /\* Processing of diagonal block \*/ For j - j<sub>b</sub> to j<sub>c</sub>

For i - j<sub>b</sub> to j<sub>c</sub> For i - j<sub>b</sub> to j<sub>c</sub>, If tuple <i, j > exists Add succ. list j to succ. list i

/\* Processing of off-diagonal rows \*/ For i - 1 to v ∧ i ∉ j<sub>b</sub> to j<sub>c</sub>, For j - j<sub>b</sub> to j<sub>c</sub>, If tuple <i, j > exists Add succ. list j to succ. list i

Figure 2 shows a 7 x 7 matrix and the order in which the elements of this matrix will be processed using the Blocked Warshall algorithm.

1	3	24	25	26	40 42	41
2	4	27	28	29		43
579	6	15	18	21	44	45
	8	16	19	22	46	47
	10	17	20	23	48	49
11	12	30	31	32	36	38
13	14	33	34	35	37	39

Figure 2. The order of computation in the Blocked Warshall algorithm [1]

Note that the equation is centered using a center tab stop. Be sure that the symbols in your equation have been defined before or immediately following the equation. Use "Eq. 1" or "Equation 1", not "(1)", especially at the beginning of a sentence: "Equation 1 is  $\ldots$ "

## Algorithm 2: Block Warshall

/\* First Pass \*/

For each row partition (rows ib to ic, inclusive)

For j - 1 to  $i_c$ , For  $i - i_b$  to  $i_c$ ,

> If tuple <i, j > exists Add succ. list j to succ. list i

/\* Second Pass \*/

For each row partition (rows  $i_b$  to  $i_c$  inclusive)

For  $j - i_b$  to n

For  $i - i_b$  to  $i_c$ ,

If tuple Xi, j > exists

Add succ. list j to succ. list i

In the second pass, one could recollect the whole last component in every column analyzed in the main pass (recall

that we assessed past the corner to corner component by and large) and inspect the rest. Moreover, the line dividing in the second pass need not be same as in the primary pass.

Figure 3 shows the order in which the elements of a 7 x 7 matrix will be processed using the Blocked Warren algorithm. The horizontal lines bracket the row partitions and the thick stair-way lines separate the two passes. Notice that the order of computation is significantly different from the straight Warren, straight Warshall, or Blocked Warshall.

1	3	34	35	36	37	38
2	4	39	40	41	42	44
5	8	11	14	17	43	45
9	9 10	12 13	15 16	18 19	46 47	48 49
20 21	22 23	24 25	26 27	28 29	30 31	32 33

Figure 3. The order of computation in the Blocked Warrren algorithm singular [1]

	Table 1. Synthetic Database[1]							
Туре	Name	Number of Nodes	Number of Arcs	Nominal Out Degree	Nominal Arc Length	Avg. Arc Length	Arcs in Result	
Uniformly	u.1	2700	2685	1	Large	564.2	49873	
Random	u.10	300	2165	10	Large	47.4	50601	
Random	h.1	2612	2596	1	1	1.6	46127	
with High Locality	h.10	390	1603	10	1	2.0	50092	
Random	m.1	7200	7152	1	10	10.5	50036	
with Medium Locality	m.10	230	2031	10	10	11.4	50601	
Tree	t.1	1720	1731	1 <sup>a</sup>			50282	
	t.10	1200	12196	10 <sup>a</sup>			50176	
Inverted	It.1	1720	1731	1 <sup>a</sup>				50282
Tree	It.10	1200	12196	10 <sup>a</sup>				50176

a: Average Out-degree of non-leaf nodes. b: Average In-degree of non-leaf nodes

IV. COMPARATIVE PERFORMANCE OF THE ALGORITHM



Figure 4. Comparative performance of the three algorithms (Total I/O) [1]

All in all, Block Warren performs superior to the next two calculations, most importantly as the level of the chart increments and region is missing. The explanation behind this activity is that every successor rundown is prone to be upgraded all the more frequently, the higher the level of the diagram. Without area, because of communication with a wide range of hubs which are all in various allotments these redesign. Without territory, these redesigns will all occur because of cooperation with a wide range of hubs which are all in various segments. On account of the two Warshall calculations, a successor rundown is composed back subsequent to being upgraded once for every allotment handled. In Warren, the rundown is composed back just once, when the parcel to which this rundown has a place is handled. The lower composing expenses of Warren, in this manner, are clear [1].

Figure 5 demonstrates 2 transitive terminations for these 12 records. The first incorporates R1, R2, R3, R5, R6, R9, R10, R11, and R12. The second one incorporates R4, R7 and R8. Figure 3 likewise demonstrates w key or keys make a couple of records identified with one and another. For instance, R4 and R7 are identified with each other by Key 2, which implies their second keys are distinguishable. Thus, R9 and R10 are identified with



Figure 5. Transitive Closure problem for Direct Algorithm [15]

The proposed calculations accept that they should figure both the connection amongst records and the transitive terminations these relations characterize from the information record document. As said before, these are two free yet related exercises. A marginally distinctive method for planning the transitive conclusion issue is that the relations are given and the comparing transitive terminations must be processed. One method for determining the relations is by giving an arrangement of record matches or key sets. A calculation was likewise produced for this rendition of the transitive conclusion issue. It utilizes the disjoint set find and union idea and productive information structures. It takes the model calculation under 21 minutes to complete the transitive conclusion calculation for an information document comprising of 122,915,040 records of number sets. The record size is around 8GB. The system kept running on a solitary 2.8 GHz Dell PC running Windows XP [14].

Journal of Advanced Computing Research Vol. 1, Issue 1 (2016) 1-6

CIKIV23	CIKIV23	CIKIV4	CIKIV4	CIKIV4	CIKIVS	CIKIVS	CIKIV7	CIKIV7	C2KIUIS	C2K1V16	C2K1V9
R11	R12	RI	R2	R3	R5	R.6	R9	R10	R4	RS	R7
11	12	1	2	3	5	6	9	10	- 4	8	7
I			I		1	(a)		l ,			a
CIKIV23	CIKIV23	CIKIV4	CIKIV4	CIKIV4	CIKIV5	CIKIV5	CIKIV7	CIKIV7	C2K1U18	C2K1V16	C2K1V
RII	R12	RI	R2	R3	R5	Rő	R9	R10	R4	RS	<b>R</b> 7
12	12	3	3	3	6	6	10	10	4	8	7
						<b>(</b> b)					
C1K2U13	C1K2U2	CIK2VI	C1K2V10	CIK2V4	C1K2V42	CIK2V:	5 CIK2VO	CIK2V7	C2K2V3	C2K2V3	C2K2V.
R9	R2	R12	RI	R6	R10	R3	R11	R5	R4	R7	RS
10	3	12	3	6	10	3	12	6	4	7	8
	134	107 - 194 194			10 II.	(c)	10	÷.			
CIK2UI3	C1K2U2	CIK2VI	CIK2V10	CIK2V4	CIK2V42	CIK2V	S CIK2Ve	CIK2V7	C2K2V3	C2K2V3	C2K2V
R9	R2	R12	Rl	R6	R10	R3	R11	R5	R4	<b>R</b> .7	RS
10	3	12	3	6	10	3	12	6	8	8	8
						(đ)			•	•	
CIK3UI4	CIK3U14	CIK3U14	CIKGUI4	CIK3U2	CIKSUS	CIK30	CIK3U	30 CIK3U	8 C2K3U	C2K3U6	C2K3V
R9	R10	RII	R12	R5	R2	R.6	RI	R3	R7	RS	R4
10	10	12	12	6	3	6	3	3	8	8	8
	Ţ				(e)		1				
CIK3UI4	CIK3U14	C1K3U14	C1K3U14	C1K3U2	1 CIK3U3	CIK3U	CIK3U	30 CIK3U	8 C2K3U	C2K3U6	C2K3V
R9	R10	R11	R12	R5	R2	Rő	Rl	R3	<b>R</b> 7	RS	R4
12	12	12	12	6	6	6	6	6	8	8	8
						(f)					
C1K4U15	C1K4U31	CIK4V2	CIK4V34	CIK4V4	CIK4V5	CIK4VS	CIK4V6	6 C1K4V7	C2K4V3	C2K4V4	C2K4V
R9	Rl	R2	R10	Rő	R3	R12	R11	R5	R4	R7	RS
12	6	6	12	6	6	12	12	6	8	8	8
	L	1		1	· 1	(g)		· I			
CIK4U15	C1K4U31	CIK4V2	C1K4V34	C1K4V4	CIK4V5	CIK4V	CIK4Ve	CIK4V7	C2K4V3	C2K4V4	C2K4V
R9	RI	R2	R10	Rő	R3	R12	R11	R5	R4	<b>R</b> 7	RS
12	12	12	12	12	12	12	12	12	8	8	8
			5 3		3 ÷	(h)	10	-	25	s 2	-

Figure 6. Disjoint-Set Find and Union (identical keys are indicated by same color) [15]

## V. CONCLUSION

In this paper, the significance of information quality is talked about. The paper studies and looks at the execution of a few calculations transitive conclusion issues is reasonable to enhance the execution of the investigation apparatuses, transitive conclusion issue is presented and planned. The consequences of the study recommend that the relative execution of the calculations is a solid capacity of the parameters which portray the prepared inquiry and the base connection referenced by that question. The velocity change accomplished by the super-TC calculation is credited to the way that the super-TC deals with its relegated segment of fundamental memory more shrewdly than the S-wavefront calculation does.

The advancement of proficient calculations to handle the diverse types of the transitive conclusion inquiries inside the setting huge database frameworks has as of late pulled in an extensive volume of examination exertion. The execution of the investigation apparatuses, transitive conclusion issue is presented and figured.

#### **ACKNOWLEDGEMENTS**

Alhamdulillah, first of all, I would like to thank God as finally I able to finish my conference paper that has been given by my lecturer, Dr Mohammad Faidz Bin Mohammad Said. Even though, there are a lot of problems during the process making this paper. Luckily, all the problems that faced by me have a solution with help from my friends. Besides that, I would like to thank you to all my friends because they help me a lot when I want to finish this conference paper. I also want to thank you to Universiti Tekologi MARA (UiTM) because give me an opportunity to learn more about this course, Parallel Programming (CSC 580).

#### REFERENCES

[1] Agrawal, R., & Jagadish, H. V., *Direct Algorithms for Computing the Transitive Closure of Database Relations*, Proceedings of the 13th VLDB Conference, Brighton, 255-266, 1987.

[2] Chakradhar, S. T., Agrawal, V. D., & Rothweiler, S. G., *A Transitive Closure Algorithm for Test Generation*, IEEE Transaction on Computer-Aided Design of Integrated Circuits and Systems, 12(7), 1015 – 1027, 1993.

[3] Chen, Y., A New Algorithm for Transitive Closures and Computation of Recursion in Relational Databases, Natural Sciences and Engineering Council of Canada, 2003.

[4] Dar, S., & Agrawal, R., *Extending SQL with Generalized Transitive Closure*, IEEE Transactions on Knowledge and Data Engineering, 5(5), 799-811, 1993.

[5] Eqbal, R., Lecture 23 Composition of Relations, *Transitive Closure and Warshall's Algorithm*, Department of Computer Science and Engineering IIT Kharagpur, 2008.

[6] Garmendia, L., Campo, R. G. d., López, V., & Recasens, J., An Algorithms to Compute the Transitive Closure, a Transitive Opening of a Fuzzy Proximity, Mathware and Soft Computing, 16, 175-191, 2009.

[7] Guh, K.-C., *Evaluation of Transitive Closure in Distributed Database System*, IEEE Journal On Selected Areas in Communications, 7(3), 399-407, 1989.

[8] Hirvisalo, V., Nuutila, E., & Soisalon-Soininen, E., *Transitive Closure Algorithm MEMTC* and Its Performance Analysis, Discrete Applied Mathematics, 110, 77-84, 2001.

[9] Ioannidis, Y. E., & Rantakrishnan, R., *Efficient Transitive Closure Algorithms*, Proceedings of the 14th VLDB Conference Los Angeles, California, 382-394, 1988.

[10] Naessens, H., Meyer, H. D., & Baets, B. D., Algorithms for the Computation of T-Transitive Closures, IEEE Transactions on Fuzzy System, 10(4), 541 – 551, 2002.

[11] Nuutila, E., & Soinien, E. S., A Single-Pass Algorithm for Transitive Closure, 1993.

[12] Purdom, P. W., A Transitive Closure Algorithm. BIT Numerical Mathematics, 10(1), 76-94. doi: 10.1007/BF01940892, 1968.

[13] Qadah, G. Z., Henschen, L. J., & Kim, J. J., *Efficient Algorithms for the Instantiated Transitive Closure Queries*, IEEE Transactions on Software Engineering, 17(3), 296-309, 1991.

[14] Toroslu, I. H., & Qadah, G. Z., *The Strong Partial Transitive* -*Closure Problem Algorithm and Performance Evaluation*, IEEE Transaction on Knowledge and Data Engineering, 8(4), 617-629, 1996.

[15] Zhang, J., Bheemavaram, R., & Li, W. N., *Transitive Closure of Data Records Application and Computation*, ALAR Conference on Applied Research in Information Technology, 71-81, 2006.