

Optimal Task Ordering in Chain Data Flows: Exploring the Practicality of Non-scalable Solutions

Georgia Kougka and Anastasios Gounaris^(✉)

Department of Informatics, Aristotle University of Thessaloniki,
Thessaloniki, Greece
{georkoug,gounaria}@csd.auth.gr

Abstract. Modern data flows generalize traditional Extract-Transform-Load and data integration workflows in order to enable end-to-end data processing and analytics. The more complex they become, the more pressing the need for automated optimization solutions. Optimizing data flows comes in several forms, among which, optimal task ordering is one of the most challenging ones. We take a practical approach; motivated by real-world examples, such as those captured by the TPC-DI benchmark, we argue that exhaustive non-scalable solutions are indeed a valid choice for chain flows. Our contribution is that we thoroughly discuss the three main directions for exhaustive enumeration of task ordering alternatives, namely backtracking, dynamic programming and topological sorting, and we provide concrete evidence up to which size and level of flexibility of chain flows they can be applied.

1 Introduction

Data analysis in a highly dynamic environment becomes more and more critical in order to extract high-quality information from raw data and derive actionable information in a timely manner. To this end, we typically employ fully automated *data-centric flows* (or *data flows*) both for business intelligence [4, 10] and scientific purposes [13], which typically execute under demanding performance requirements, e.g., to complete in a few seconds. These flows generalize traditional Extract-Transform-Load (ETL) and data integration flows through the incorporation of data analytics [8, 17]. Meeting the demanding performance requirements, combined with the volatile nature of the environment and the data, gives rise to the need for efficient data flow optimization techniques.

Data flow optimization techniques cover a wide spectrum from deciding on the order of the constituent tasks to detailed low-level configuration of the underlying execution engine [12]. In this work, we focus on the former aspect, namely the specification of the execution order of the constituent tasks. In practice, this is usually the result of a manual procedure, which, in many cases results in non-optimal flow execution plans. Furthermore, even if a data flow is optimal for a specific input data set, it may prove significantly suboptimal for another data set

with different characteristics [7]. We tackle this problem through the proposal of optimization algorithms that can provide the optimal execution order of the tasks in a *chain* (or *linear*) data flow in an efficient manner and relieve the flow designers from the burden of selecting the task ordering on their own. We consider a single optimization objective, namely the minimization of the sum of the task execution costs; we assume that the execution cost of each task depends on the volume of data to be processed, which in turn depends on the relative position of the task in the execution flow.

The main challenges in flow optimization that need to be addressed and differentiate the problem from that of traditional query optimization, discussed in [3, 9], are that the tasks need not belong to a set of well-defined algebra, such as the relational one, there exist arbitrary precedence constraints among operators, and flows can consist of dozens of tasks, whereas, typically, operators in query plans are fewer. The main implication is that query optimization techniques, which operate on plans with up to a few tens of operators that belong to the relational algebra (according to which operator reordering is typically permitted), are not applicable. Nevertheless, they are successful in their domain and this is the reason the data flow solutions proposed in this work are partially inspired by query optimization, as we explain later. Overall, to date, there are very few proposals that deal with (or are applicable to) task reordering in data flows [8, 16, 20]. A common characteristic of these proposals is that they are too slow to find an exact solution in small flows [8], or they can find significantly suboptimal (approximate) solutions for bigger flows [16, 20].

In this work, we go beyond the state-of-the-art with regards to exact solutions in chain data flows, i.e., data flows where the tasks form a sequence. Chain data flows are a main building block in generic data flows. Optimization of chain flows is a big step towards optimization of more generic ones. Exact solutions cannot scale in general, but, partially inspired by real worlds-like data flows, such as those in the TPC-DI benchmark [14], we show that they can be applied in several cases in practice. The main contribution of this work is the proposal of two additional exact solutions that significantly improve upon the technique in [8] in terms of time overhead and the size of data flows they can handle.¹

The remainder of this paper is structured as follows. In Sect. 2, we present the notation, the problem statement and the motivation from TPC-DI. Our exact solutions are explained and evaluated in Sects. 3 and 4, respectively. We mention related work in Sect. 5, and we conclude in the next section, which discusses the issues in optimizing more generic flows than chain ones.

2 Preliminaries

In this paper, we deal with the problem of re-ordering the tasks of a chain data flow without violating existing precedence constraints between tasks, while the performance of the flow is maximized. The data flow is represented as a

¹ An abstract of these ideas, without considering TPC-DI, have appeared in [11] in less than a page.

directed acyclic graph (DAG), denoted as $G = (T, E)$, where each task $t_i \in T$, $i = 1, \dots, n$ corresponds to a node in the graph and the edges between nodes represent intermediate data shipping among tasks; i.e., in data flows, the exchange of data between tasks is explicitly represented through edges. The tasks that have no incoming edges are termed as *sources*, and those without outgoing edges as *sinks*.

The precedence constraints are captured by another directed acyclic graph $PC = (T, D)$, such that each ordered pair of vertices (t_i, t_j) that either belongs to D or the transitive closure of PC corresponds to the requirement that, in any valid G , there must exist a directed path from t_i to t_j . In other words, the PC graph corresponds to a higher-level, non-executable flow representation, where the exact task ordering is not defined; only a partial ordering is defined instead.

We focus on linear or chain data flows, where the flow contains only one sink and one source and each task in between has exactly one incoming and outgoing edge. Chain data flows can be regarded as sub-flows within generic data flows. For example, each path from a source to a sink in a generic flow forms a chain.

Further, we assume that each task t_i has a processing cost per input record (or tuple) c_i , and selectivity sel_i . The selectivity denotes the average number of returned tuples per input tuple. For filtering tasks, $sel_i < 1$; for data sources and operators that just manipulate the input $sel = 1$, whereas, for operators that may produce more output records for each input record, $sel_i > 1$. Figure 1(left) shows an example of a chain data flow with 5 tasks, their precedence constraints, and one possible valid G that respects such constraints. In the middle, example task metadata are presented. As will be shown later (e.g., Fig. 5), there are multiple other task orderings that respect the constraints as well, e.g., placing t_3 after t_4 .

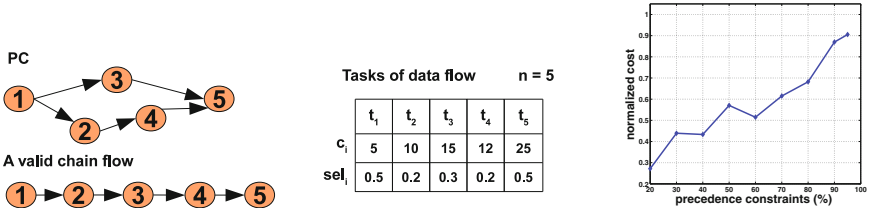


Fig. 1. An example of PC and G (left), task metadata (middle) and indicative performance improvements (right).

Problem Statement: Given a set of tasks T with known cost and selectivity values, and a corresponding precedence constraint graph PC , we aim to find a valid task ordering to form a chain G that minimizes the *sum cost metric* (SCM) per source tuple. SCM is defined as follows:

$$SCM(G) = \sum_{i=1}^n \left(\prod_{j \in Pred(t_i)} sel_j \right) c_i,$$

where $Pred(t_i)$ is the set of tasks that precede t_i in G ; if the set is null, the product of selectivities is set to 1. In the example of Fig. 1, $Pred(t_4) = \{t_1, t_2, t_3\}$. The optimal plan is denoted as P .

Figure 1(right) shows indicative performance improvements. We examine 100 randomly generated data flows consisting of $n = 15$ tasks with $c_i \in [1, 100]$, $sel_i \in (0, 2]$ and 20%–95% precedence constraints. A chain flow has 100% precedence constraints, when the transitive closure of its PC has $\frac{n(n-1)}{2}$ edges. The case of 100% precedence constraints is when there is a single valid ordering with no ordering alternatives. We use the percentage of precedence constraints as an efficient way to quantify the flexibility in reordering operators in a flow. The initial random ordering has normalized cost 1. As can be observed from Fig. 1(right), the improvements can be of a factor of 3 and more.

Note that the input set of tuples are processed by all the tasks of the chain data flow until they are filtered out, but typically, some of the input tuple attributes may not be required by every flow activity. According to [21], the unnecessary tuple attributes just run through the flow, resembling an assembly-line model. The execution of a flow activity is not affected by the unnecessary attributes. This implies that the tasks of a flow have the ability to be reordered as long as the precedence constraints between the tasks are preserved.

2.1 Problem Complexity

In [2], Burge et al. proved that finding the optimal ordering of tasks is an NP -hard problem when (i) each flow task is characterized by its cost per input record and selectivity; (ii) the cost of each task is a linear function of the number of records processed and that number of records depends on the product of the selectivities of all preceding tasks (assuming independence of selectivities for simplicity); and (iii) the optimization criterion is the minimization of the sum of the costs of all tasks. All the above conditions hold for our case, so our problem is intractable. Moreover, in [2] it is discussed that “*it is unlikely that any polynomial time algorithm can approximate the optimal plan to within a factor of $O(n^\theta)$* ”, where θ is some positive constant. Note that if we modify the optimization criterion, e.g., to optimize the bottleneck cost metric or the critical path renders the problem tractable [1, 18].

2.2 Chains in TPC-DI

TPC-DI [14] is the TPC standard for data integration flows. Although it focuses on extensions to ETLs rather than on both advanced ETLs and analytics, it is the closest standard to our scenarios. TPC-DI aims to model the data integration processes of a retail brokerage firm. More importantly, it has been implemented for an open-source data flow engine, namely Pentaho Kettle.² Figure 2 shows an example of part of the implementation, which is responsible for building one of the seven dimension tables in the underlying data warehouse defined by the standard.

² <http://www.essi.upc.edu/dtim/blog/post/tpc-di-etls-using-pdi-aka-kettle>.

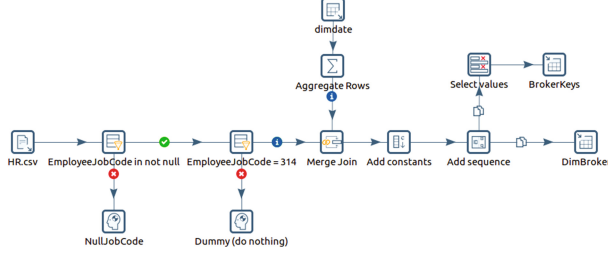


Fig. 2. A TPC-DI flow implemented in Kettle.

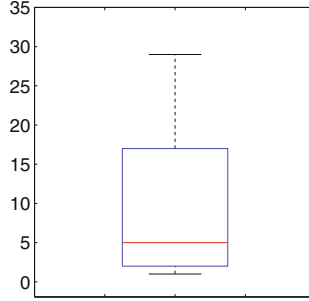


Fig. 3. Boxplot diagram for the length of chains in TPC-DI flows.

Next, we proceed to a simple analysis of the flows referring to historical loading, as these are implemented in Kettle. First, we extract chain sub-flows according to the following procedure: for each sink, we take the longest path from any of the connected sources. The boxplot diagram in Fig. 3 reveals significant information regarding the size of such chains. More specifically, most of the chains considered contain less than 20 tasks. Also, only 2 chains have size larger than 30 tasks; these chains are removed from the boxplot as outliers.

Analyzing the percentage of precedence constraints, it is found that there is a small degree of flexibility in re-ordering tasks. For example, in the largest chain in the boxplot, there are 97% precedence constraints. This pattern appears in all flows with size more than the median. For smaller chains, such as those contained in the flow in Fig. 2, the constraints drop to 87%, which is still high. However, according to the evidence in Fig. 1(right), significant improvements of 10–15% can still be achieved.

3 Accurate Algorithms for Linear Execution Plans

In this section, we present three accurate algorithms for reordering chain data flows in order to generate an optimal execution plan. The algorithms are based on backtracking, dynamic programming and generation of all topological sortings,

respectively. Our main novelty here is that we examine a topological sorting-based algorithm, despite its worst-case complexity. Counter-intuitively, as we show in the evaluation, the algorithm is practical not only for the type of flow chains appearing in TPC-DI, but also for much larger n , when there are many precedence constraints and, in general, can scale better than the two other options. However, still, it cannot be applied to arbitrary flows of very large size.

3.1 Backtracking

The *Backtracking* algorithm finds all the possible execution plans generated after reordering the tasks of a given data flow preserving the precedence constraints. The algorithm enumerates all the valid sub-flow plans after applying a set of recursive calls on these sub-flows until generating all the possible data flow plans. It backtracks when a placement of a task in a specific position violates the precedence constraints. The algorithm is proposed for flow optimization in [8].

Complexity: The worst case time complexity of *Backtracking* is factorial (i.e., $O(n!)$), since, if there are no dependencies, all orderings will be examined in a brute force manner.

3.2 Dynamic Programming

This algorithm is extensively used as part of the System R-type of query optimization to produce (linear) join orderings [15]. The rationale of the dynamic programming algorithm (termed as *DP* henceforth) for data flows remains the same, that is to calculate the cost of task subsets of size n based on subsets of size $n - 1$. For each of these subsets, we keep only the optimal solutions, which are valid with regards to the precedence constraints. Specifically, the *DP* algorithm considers each flow of size n as a flow of $(n - 1)$ tasks followed by the n th task; the key point is that the former part is the optimal subset of size $n - 1$, which has been found from previous step; then the algorithm exhaustively examines which of the n flow tasks is the one that, when added at the end, yields an optimal subplan of size n . For example, the algorithm starts by calculating subsets that consist of only one task $\{t_1\}$, then $\{t_2\}$, $\{t_3\}$ and so on. In a similar way, in the second step, it examines subsets containing two tasks, i.e., $\{t_1, t_2\}$, $\{t_1, t_3\}$ and so on, until it examines the complete flow $\{t_1, t_2, \dots, t_n\}$. The number of the optimal (non-empty) subsets of a flow is equal to $2^n - 1$.

Complexity: The time complexity is $O(n^2 2^n)$. This is because we examine all subsets of n tasks, which are $O(2^n)$. For each subset, which is up to size $O(n)$, we examine whether each element can be placed at the end of the subplan. Each such check involves testing whether any of the rest $n - 1$ tasks violate a precedence constraint, when placed before the n -th task. Overall, for each element, we make $O(n)$ comparisons. So, the overall time complexity is $O(2^n)O(n)O(n) = O(n^2 2^n)$. The space complexity is derived by the size of the auxiliary data structures employed. We use three vectors of size $2^n - 1$ as

Algorithm 1. Dynamic Programming

Require: 1. A set of n tasks, $T = \{t_1, \dots, t_n\}$. 2. A directed acyclic graph PC with precedence constraints.

Ensure: A directed acyclic graph P representing the optimal plan

{Initialize *PartialPlan*, *Costs* and *Sel* of size $2^n - 1$ }

```

1: for all  $i \in \{2, \dots, n\}$  do
2:    $\text{PartialPlan}(2^{i-1}) = t_i$ 
3:    $\text{Costs}(2^{i-1}) = c_i$ 
4:    $\text{Sel}(2^{i-1}) = \text{sel}_i$ ;
5: end for
6: for all  $s \in \{2, \dots, n\}$  do
7:    $R \leftarrow \text{Subsets}(T, s)$  { $X$  is a set with all subsets of  $T$  of size  $s$ }
   { $r$  is a specific subset of size  $s$ }
8:    $\text{tempBest} \leftarrow \infty$ 
9:   for each  $r \in R$  do
10:    for all  $i \in \{1, \dots, r.\text{length}()\}$  do
11:       $\text{tempSet} \leftarrow r - r(i)$ 
12:       $\text{pos1} \leftarrow \text{findIndex}(\text{tempSet})$ 
13:       $\text{pos2} \leftarrow \text{findIndex}(r(i))$ 
14:      if  $r(i)$  has all predecessors in  $\text{tempSet}$  then
15:         $\text{TempPlan} \leftarrow \text{tempSet}, r(i)$ 
16:         $\text{costTempPlan} \leftarrow \text{Costs}(\text{pos1}) + \text{Sel}(\text{pos1})\text{Costs}(\text{pos2})$ 
17:        if  $\text{costTempPlan} < \text{tempBest}$  then
18:           $\text{tempBest} \leftarrow \text{costTempPlan}$ 
19:           $k \leftarrow \text{pos1} + \text{pos2}$ 
20:           $\text{update}(\text{PartialPlan}(k), \text{Costs}(k), \text{Sel}(k))$ 
21:        end if
22:      end if
23:    end for
24:  end for
25: end for
26:  $P \leftarrow \text{PartialPlan}(2^n - 1)$ 

```

explained in the implementation details, one of which stores elements of size $O(n)$. So the space complexity is $O(n2^n)$.

Implementation Issues: In order to implement the algorithm, we use three vectors of size $2^n - 1$, namely *PartialPlan*, *Costs* and *Sel*. According to the algorithm implementation, the i -th cell corresponds to the combination of tasks for which the bit is 1 in its binary representation. For example, if $i = 13$, then the binary representation of this position is $(1101)_2$. Specifically, this means that *partialPlan*[13] corresponds to the optimal ordering of the 1st, 2nd and 4th tasks. Analogously, the partial plan $\{1, 3, 4, 5\}$ is stored in position $2^{1-1} + 2^{3-1} + 2^{4-1} + 2^{5-1} = 29$ of the *partialPlan* matrix. The *Costs* and *Sel* vectors hold the aggregate cost and selectivity of the subplans, respectively. The last cell of *PartialPlan* and *Costs* contain the optimal plan and its total cost, respectively. A complete pseudocode is shown in Algorithm 1. For the sake

of simplicity of presentation, the algorithm is not fully optimized; e.g., in line 18, the update of vertices may occur only once after the final best plan is found.

An example: We give an example of the algorithm with a flow with $n = 5$; the task metadata are shown in Fig. 1. The *DP* example is in Fig. 4. First of all, all the subsets R of T of length $s = \{1, 2, \dots, n\}$ are found. For single task subsets, such as $\{t_1\}, \{t_2\}, \dots, \{t_n\}$, *DP* estimates their position in the *partialPlan* matrix, e.g. $\{2\}$ subset is positioned in *partialPlan*(2^{2-1}). For subsets with length greater than 1, e.g., the subset $\{1, 3, 4\}$, we examine the case that each element of that subset is placed at the end of the subset. If the precedence constraints are violated, *DP* continues to the next placement. If the precedence constraints are not violated, the algorithm estimates the cost of the valid partial plan with that element positioned at the end of the subset, reusing the results of the orderings of smaller subsets. Similarly, the cost of all orderings in the subset is estimated and the algorithm finds the ordering of the subset with the minimum cost. The optimal partial plan, its cost and the product of task selectivities are stored in the corresponding position in the *partialPlan* and *Costs* and *Sel* vertices, respectively.

DP algorithm			Subsets R of different sizes	
Costs	Sel	partialPlan		
1: 5	0,5	1: 1	4 5	3 4 5
2: 10	0,2	2: 2	3 5	2 4 5
3: 10	0,1	3: 12	3 4	2 3 5
4: 15	0,3	4: 3	2 5	2 3 4
12	0,15	13	2 4	1 4 5
13	0,06	23	2 3	1 3 5
11,5	0,03	123	1 5	1 3 4
8: 12	0,2	8: 4	1 4	1 2 5
11	0,1	14	1 3	1 2 4
12,4	0,04	24	1 2	1 2 3
11,2	0,02	124		
15	0,06	43		
12,5	0,03	143		
13	0,012	243		
11,5	0,006	1243		
16: 25	0,5	16: 5		
17,5	0,25	15		
15	0,1	25		
12,5	0,05	125		
22,5	0,15	35		
16,25	0,075	135		
14,5	0,03	235		
12,25	0,015	1235		
24: 17	0,1	24: 45		
13,5	0,05	145		
13,4	0,02	245		
11,7	0,01	1245		
16,5	0,03	435		
29: 13,25	0,015	29: 1435		
30: 13,3	0,006	30: 2435		
31: 11,65	0,003	31: 12345		

Example computations:

- For subsets of size 2, we examine the subset {4 5}:
 - 4 – 5 valid partial plan (no precedence constraints)
 - 5 – 4 not valid partial plan (precedence constraints)
- 1) Find the position of subsets {4} and {5}:
 - pos(4) = $2^{4-1} = 2^3 = 8$
 - pos(5) = $2^{5-1} = 2^4 = 16$
 - pos(4-5) = pos(4) + pos(5) = 8 + 16 = 24
- 2) Cost and selectivity estimation of partial order 4 – 5:
 - Costs(24) = Costs(8) + Sel(8) * Costs(16) = 12 + 0,2*25 = 17
 - Sel(24) = Sel(8) * Sel(16) = 0,2 * 0,5 = 0,1
- For subsets of size 3, we examine the subset {2 4 5}:
 - 2 – 4 – 5 valid partial plan (no precedence constraints)
 - 4 – 5 – 2 not valid partial plan (precedence constraints)
 - 2 – 5 – 4 not valid partial plan (precedence constraints)
- For 2 – 4 – 5:
 - 1) Find the position of subsets {2} and {45}:
 - pos(4-5) = $2^{4-1} + 2^{5-1} = 24$
 - pos(2) = $2^{2-1} = 2$
 - pos(2-4-5) = pos(4-5) + pos(2) = 24 + 2 = 26
 - 2) Cost and selectivity estimation of partial order 2 – 4 – 5:
 - Costs(26) = Costs(2) + Sel(2) * Costs(24) = 13,4
 - Sel(26) = Sel(2) * Sel(24) = 0,02

Fig. 4. Example of the DP algorithm.

Algorithm 2. TopSort

Require: 1. A set of n tasks, $T=\{t_1, \dots, t_n\}$. 2. A directed acyclic graph PC with precedence constraints.

Ensure: An ordering of the tasks P representing the optimal plan.

```

1:  $P=\{t_1, t_2, \dots, t_n\}$  { $P$  is initialized with a valid topological ordering of  $PC$ .}
2:  $i=1$ 
3:  $\text{minCost} \leftarrow \text{computeSCM}(P)$ 
4: while  $i < n$  { $n$  is the total number of tasks} do
5:    $k \leftarrow \text{location of } t_i \text{ in } P$ 
6:    $k1 \leftarrow k + 1$ 
7:   if  $P(k1)$  task has prerequisite  $t_i$  then
8:     // Rotation stage
9:     Rotate the elements of  $P$  from positions  $i$  to  $k$ 
10:     $\text{cost} \leftarrow \text{computeSCM}(P)$ 
11:     $i \leftarrow i+1$ 
12:   else
13:     // Swapping stage
14:     Swap the  $k$  and  $k1$  elements of  $P$ 
15:      $\text{cost} \leftarrow \text{computeSCM}(P)$ 
16:      $i \leftarrow 1$ 
17:   end if
18:   if  $\text{cost} < \text{minCost}$  then
19:      $\text{bestP} \leftarrow P$ 
20:      $\text{minCost} = \text{cost}$ 
21:   end if
22: end while
23:  $P \leftarrow \text{bestP}$ 

```

Correctness: If *PartialPlan* is of size $n = 1$, the optimal solution is trivial and is found by the algorithm during initialization in lines 1–3 of Algorithm 1. We assume that a *PartialPlan* of size $n - 1$ is optimal and we need to prove that *PartialPlan* of size n is also optimal. The sketch of the proof will be based on contradiction. Let us assume that the *DP* does not produce the optimal solution. Any linear solution of size n consists of a *PartialPlan* of size $n - 1$ followed by the n -th task; *DP* checks all the alternatives for the n -th task. So, there is a different optimal solution, where the *PartialPlan* of size $n - 1$ is different of *DP*'s *PartialPlan* of the same size. According to the *SCM*, the cost of the subplan of size n is computed as the sum of two components: the cost of subplan of size $n - 1$ and the cost of the n -th task times the selectivity of the first $n - 1$ tasks. The costs of the solutions of size n , which end with the same task, differ only in the first component. According to our assumptions, the cost of *DP*'s *PartialPlan* of size $n - 1$ cannot be higher than any other subplan solution of size $n - 1$ by definition. Consequently, there is no other solution different from *DP*'s solution that can yield lower cost. This completes the proof.

3.3 Topological Sorting

The *TopSort* algorithm is a topological sorting algorithm based on [19], which finds all the possible topological sortings given a partial ordering of a finite set; in our case the partial ordering is due to the precedence constraints. The reason behind using this algorithm is that it (implicitly) prunes invalid plans very efficiently and it generates a new plan based on a previous plan after performing a minimal change. For the purposes of this work, we adapted the topological sorting algorithm in order to generate all the possible execution plans of a data flow and detect the execution plan with the minimum cost. The algorithm assumes that it can receive as input a valid task permutation $t_1 \rightarrow t_2 \rightarrow t_3 \rightarrow \dots \rightarrow t_n$, which is trivial since it can be done in linear time. We generate all other valid execution plans by applying two main operations, namely, cyclic rotations and swapping adjacent tasks.

Firstly, the process of generating all the valid flow execution plans begins with the topological sorting of the $n - 1$ tasks $t_2 \rightarrow t_3 \rightarrow \dots \rightarrow t_n$ of the flow. Based on this partial sorting, we generate all the valid orderings of the $t_1 \rightarrow t_2 \rightarrow t_3 \rightarrow \dots \rightarrow t_n$ plan. Specifically, in the first stage of the algorithm the task t_1 is placed on the left part of the partial plan $t_2 \rightarrow t_3 \rightarrow \dots \rightarrow t_n$ and in the next steps of this stage, we swap it with the tasks on its right, while the tasks of the partial plan maintain their relative position. The t_1 task stops moving when such a swap violates a precedence constraint. Then, as the task t_1 cannot be further transposed, the second stage of algorithm begins with a right-cyclic rotation of another partial plan consisted of t_1 and all the tasks that precede it, which means all the tasks which are positioned to its left. In this way, t_1 is placed to its initial position. Similarly, we generate all the topological sortings of $t_2 \rightarrow t_3 \rightarrow \dots \rightarrow t_n$, $t_3 \rightarrow t_4 \rightarrow \dots \rightarrow t_n$ and so on. For each generated plan, we estimate the total execution cost. A pseudocode is presented in Algorithm 2.

Complexity: Since the algorithm checks all the permutations the time complexity is $O(n!)$ in the worst case. However, compared to other algorithms that produce all topological sortings, it is more efficient [19]. The space complexity is $O(n)$ because only one plan is stored in main memory at any point of execution.

Implementation Issues: The algorithm exhaustively checks all the permutations that satisfy the precedence constraints, and as such, it always finds the optimal solution for linear flows. No specific data structures are required. As shown in Algorithm 2, we employ a *computeSCM* function needs to be constructed in a way that does not compute the cost of each ordering from scratch, which is too naive, but leverages the computations of the previous plans taking into account the local changes in the new plan. Note that we can implement *TopSort* in a different way, where the tasks are checked from right to left. Although in [19] this flavour is claimed to be capable of yielding better performance, this has not been verified in our flows.

Example: In Fig. 5, an example of finding the optimal plan of a flow using *TopSort* is presented. In this example, the running steps of *topSort* algorithm are depicted, given as input a valid flow execution plan (*Initial plan order plan label*)

TopSort algorithm									
i	k	k1							
			Initial plan order						
			P = <table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table>	1	2	3	4	5	SCM(P) = 12,01
1	2	3	4	5					
1	1	2	Rotate from P(1) to P(1), set i=2	After rotation P = <table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table> SCM(P) = 12,01	1	2	3	4	5
1	2	3	4	5					
2	2	3	Swap P(2) and P(3), set i=1	After swap P = <table border="1"><tr><td>1</td><td>3</td><td>2</td><td>4</td><td>5</td></tr></table> SCM(P) = 14,51	1	3	2	4	5
1	3	2	4	5					
1	1	2	Rotate from P(1) to P(1), set i=2	After rotation P = <table border="1"><tr><td>1</td><td>3</td><td>2</td><td>4</td><td>5</td></tr></table> SCM(P) = 14,51	1	3	2	4	5
1	3	2	4	5					
2	3	4	Rotate from P(2) to P(3), set i=3	After rotation P = <table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table> SCM(P) = 12,01	1	2	3	4	5
1	2	3	4	5					
3	3	4	Swap P(3) and P(4), set i=1	After swap P = <table border="1"><tr><td>1</td><td>2</td><td>4</td><td>3</td><td>5</td></tr></table> SCM(P) = 11,65	1	2	4	3	5
1	2	4	3	5					
1	1	2	Rotate from P(1) to P(1), set i=2	After rotation P = <table border="1"><tr><td>1</td><td>2</td><td>4</td><td>3</td><td>5</td></tr></table> SCM(P) = 11,65	1	2	4	3	5
1	2	4	3	5					
2	2	3	Rotate from P(2) to P(2), set i=3	After rotation P = <table border="1"><tr><td>1</td><td>2</td><td>4</td><td>3</td><td>5</td></tr></table> SCM(P) = 11,65	1	2	4	3	5
1	2	4	3	5					
3	4	5	Rotate from P(3) to P(4), set i=4	After rotation P = <table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table> SCM(P) = 12,01	1	2	3	4	5
1	2	3	4	5					
4	4	5	Rotate from P(4) to P(4), set i=5	After rotation P = <table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table> SCM(P) = 12,01	1	2	3	4	5
1	2	3	4	5					
FINISH			Final plan order P = <table border="1"><tr><td>1</td><td>2</td><td>4</td><td>3</td><td>5</td></tr></table>	1	2	4	3	5	SCM(P) = 11,65
1	2	4	3	5					

Fig. 5. Example of the TopSort algorithm.

and assuming the metadata of Fig. 1. Each of the given plans describe a plan generated after either a rotation or a swap action.

4 Evaluation of the Time Overhead

In this section, we conduct a thorough evaluation of the time overhead of the accurate optimization algorithms. We use a machine with an Intel Core i5 660 CPU and 6 GB of RAM. We construct synthetic flows so that we can evaluate the algorithms in a wide range of parameter combinations. More specifically, we produce the optimal ordering of tasks to form a chain flow after we have (i) created random PC graphs with a configurable ratio of precedence constraints, and (ii) chosen task metadata randomly. The time overhead does not actually depend on the task metadata.

Backtracking cannot scale in the percentage of PCs and *DP* cannot scale in the number of tasks. So, we compare them against *TopSort* separately.

Figure 6(top-left) presents the average execution time of the *DP* algorithm compared to the *TopSort* solution for 50% precedence constraints, and

$n = 15, \dots, 20$ flow tasks. For this ratio of PCs, *Backtracking* is practically inapplicable. The main conclusions that can be drawn from this figure is that *DP* algorithm is not a practical optimization solution even for small flows that consist of 19 flow activities for this number of constraints; the execution of a flow with 20 tasks requires over 3 days using our test machine. *TopSort* runs at least 50 times faster than *DP*, but its execution follows a similar pattern to *DP*. However, this figure is indicative regarding the superiority of *TopSort* over *DP*. In the top-right part of Fig. 6, the time overhead of *Backtracking* compared to *TopSort* is presented for $n = 15, \dots, 34$ and 90% PCs. The main observation of this figure is that *Backtracking* is orders of magnitude slower than *TopSort* and cannot scale even for medium-size flows of more than 30 tasks.

We now turn our attention to *TopSort* for further investigation. Figure 6 (bottom-left) shows the average execution time of *TopSort* for flows with $n = 10, \dots, 70$ having 98% precedence constraints, which implies that the number of the possible re-orderings is quite restricted, as also observed in TPC-DI. *TopSort* does not scale well, but the important thing is that it can run in acceptable time, e.g., a minute, even in flows with 60 tasks. Additionally, Fig. 6 (bottom-right) depicts that *TopSort* cannot scale for arbitrarily few precedence constraints even for flows with 15 and 20 flow activities. But, for these flow sizes, it can tolerate percentages of precedence constraints much lower than 90%.

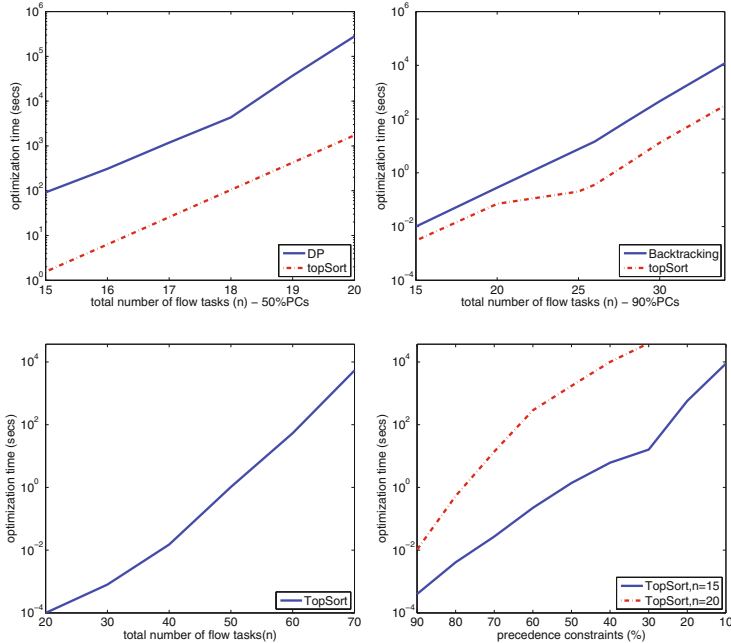


Fig. 6. Optimization overhead in several settings.

5 Related Work

Further to the discussion of related work in the introduction, task ordering for data flows is an area that has been significantly influenced by query processing techniques. In [5], an optimization algorithm for query plans with dependency constraints between algebraic operators is presented. In [20], a proposal for data integration that is also applicable to data flow task ordering is discussed. The corresponding techniques are approximate rather than exact solutions for chain flows, and they are either similar or inferior to those considered in [11]. ETL flows are analyzed in [16], where ETL execution plans are considered as states and transitions, such as swap, merge, split, distribute and so on, are employed to generate new states in order to navigate through the state space, which corresponds to the execution plan alternatives. As shown in [11], these techniques can deviate from the optimal solutions in chain dataflows by several factors.

In addition, there is a significant portion of proposals on flow optimization that proceed to flow structure optimizations but do not perform task reordering, as we do; an overview of the complete spectrum of data flow optimization techniques appears in [12]. Note that deciding the task ordering before execution in order to determine the flow structure is orthogonal to deciding the scheduling order of tasks at runtime, which is another well investigated area, e.g., see [6].

6 Conclusions

In this work, we deal with the problem of optimally ordering the constituent tasks of a chain data flow in order to minimize the sum of the task execution costs. We are motivated by the significant limitations of fully-automated optimization solutions for data flows, as, nowadays, the optimization of the complex data flows is left to the flow designers and is a manual procedure. We are also motivated by the fact that real-world flows, such as those in TPC-DI, are not very flexible, in terms of the alternative valid task orderings. As such, carefully crafted exhaustive solutions become applicable. We propose two such solutions, a dynamic programming one and one that efficiently generates all topological orderings. We explain the technical details involved and we show that the latter approach is both dominant and practical under realistic assumptions.

Our work is a big step towards optimal task ordering in generic flows instead of simple chain flows, e.g., the complete flows assumed by the TPC-DI benchmark. For this goal to be met, three further issues need to be resolved: (i) to devise solutions that, using chain optimization as a building block, apply to complete data flows; in the way we defined chains hereby, these chains are overlapping and their isolated optimizations may not be compatible among overlapping chains and do not guarantee optimality when combined; (ii) to develop efficient ways to collect the required statistical metadata and detect precedence constraints; and (iii) to take into account other metrics than the minimization of the sum of the costs with a view to considering realistic issues, such as pipelined execution and parallel execution on multi-core engines, which better reflect the running time of a flow.

References

1. Agrawal, K., Benoit, A., Dufossé, F., Robert, Y.: Mapping filtering streaming applications. *Algorithmica* **62**(1–2), 258–308 (2012)
2. Burge, J., Munagala, K., Srivastava, U.: Ordering pipelined query operators with precedence constraints. Technical report 2005–40, Stanford InfoLab (2005)
3. Chaudhuri, S.: An overview of query optimization in relational systems. In: PODS, pp. 34–43 (1998)
4. Chaudhuri, S., Dayal, U., Narasayya, V.: An overview of business intelligence technology. *Commun. ACM* **54**, 88–98 (2011)
5. Florescu, D., Levy, A., Manolescu, I., Suciu, D.: Query optimization in the presence of limited access patterns. In: SIGMOD, pp. 311–322. ACM (1999)
6. Grehant, X., Demeure, I., Jarp, S.: A survey of task mapping on production grids. *ACM Comput. Surv.* **45**(3), 37:1–37:25 (2013)
7. Halasipuram, R., Deshpande, P.M., Padmanabhan, S.: Determining essential statistics for cost based optimization of an ETL workflow. In: EDBT, pp. 307–318 (2014)
8. Hueske, F., Peters, M., Sax, M., Rheinländer, A., Bergmann, R., Krettek, A., Tzoumas, K.: Opening the black boxes in data flow optimization. *PVLDB* **5**(11), 1256–1267 (2012)
9. Ioannidis, Y.E.: Query optimization. *ACM Comput. Surv.* **28**(1), 121–123 (1996)
10. Jovanovic, P., Romero, O., Abelló, A.: A unified view of data-intensive flows in business intelligence systems: a survey. *T. Large-Scale Data Knowl. Centered Syst.* **29**, 66–107 (2016)
11. Kougka, G., Gounaris, A.: Optimization of data-intensive flows: is it needed? Is it solved? In: DOLAP, pp. 95–98 (2014)
12. Kougka, G., Gounaris, A., Simitsis, A.: The many faces of data-centric workflow optimization: a survey. *CoRR*, abs/1701.07723 (2017)
13. Ogasawara, E.S., de Oliveira, D., Valduriez, P., Dias, J., Porto, F., Mattoso, M.: An algebraic approach for data-centric scientific workflows. *PVLDB* **4**, 1328–1339 (2011)
14. Poess, M., Rabl, T., Caufield, B.: TPC-DI: the first industry benchmark for data integration. *PVLDB* **7**(13), 1367–1378 (2014)
15. Selinger, P.G., Astrahan, M.M., Chamberlin, D.D., Lorie, R.A., Price, T.G.: Access path selection in a relational database management system. In: SIGMOD, pp. 23–34 (1979)
16. Simitsis, A., Vassiliadis, P., Sellis, T.K.: State-space optimization of ETL workflows. *IEEE Trans. Knowl. Data Eng.* **17**(10), 1404–1419 (2005)
17. Simitsis, A., Wilkinson, K., Castellanos, M., Dayal, U.: Optimizing analytic data flows for multiple execution engines. In: SIGMOD, pp. 829–840 (2012)
18. Srivastava, U., Munagala, K., Widom, J., Motwani, R.: Query optimization over web services. In: Proceedings of the 32nd International Conference on Very Large Data Bases VLDB, pp. 355–366 (2006)
19. Varol, Y.L., Rotem, D.: An algorithm to generate all topological sorting arrangements. *Comput. J.* **24**(1), 83–84 (1981)
20. Yerneni, R., Li, C., Ullman, J., Garcia-Molina, H.: Optimizing large join queries in mediation systems. In: Beeri, C., Buneman, P. (eds.) ICDT 1999. LNCS, vol. 1540, pp. 348–364. Springer, Heidelberg (1999). doi:[10.1007/3-540-49257-7_22](https://doi.org/10.1007/3-540-49257-7_22)
21. Zinn, D., Bowers, S., McPhillips, T., Ludäscher, B.: Scientific workflow design with data assembly lines. In: Proceedings of the 4th Workshop on Workflows in Support of Large-Scale Science, WORKS 2009, pp. 14:1–14:10. ACM (2009)

Big Data Analytics and Knowledge Discovery
19th International Conference, DaWaK 2017, Lyon,
France, August 28–31, 2017, Proceedings
Bellatreche, L.; Chakravarthy, S. (Eds.)
2017, XIV, 488 p. 137 illus., Softcover
ISBN: 978-3-319-64282-6