

Robust Distributed Stream Processing

Chuan Lei, Elke A. Rundensteiner, Joshua D. Guttman

Computer Science Department, Worcester Polytechnic Institute
 100 Institute Road, Worcester, MA 01609, USA
 {chuanlei, guttman, rundenst}@cs.wpi.edu

Abstract—Distributed stream processing systems must function efficiently for data streams that fluctuate in their arrival rates and data distributions. Yet repeated and prohibitively expensive load re-allocation across machines may make these systems ineffective, potentially resulting in data loss or even system failure. To overcome this problem, we instead propose a load distribution (RLD) strategy that is robust to data fluctuations. RLD provides ϵ -optimal query performance under load fluctuations without suffering from the performance penalty caused by load migration. RLD is based on three key strategies. First, we model robust distributed stream processing as a parametric query optimization problem. The notions of robust logical and robust physical plans then are overlays of this parameter space. Second, our Early-terminated Robust Partitioning (ERP) finds a set of robust logical plans, covering the parameter space, while minimizing the number of prohibitively expensive optimizer calls with a probabilistic bound on the space coverage. Third, our OptPrune algorithm maps the space-covering logical solution to a single robust physical plan tolerant to deviations in data statistics that maximizes the parameter space coverage at runtime. Our experimental study using stock market and sensor networks streams demonstrates that our RLD methodology consistently outperforms state-of-the-art solutions in terms of efficiency and effectiveness in highly fluctuating data stream environments.

I. INTRODUCTION

Motivation. Distributed stream processing systems (DSPSs) are designed to execute continuous queries over streams of tuples [1], [2], [3]. Continuous queries place heavy workloads on precious system resources from CPU processing cycles, memory, and network bandwidth. Since workloads can vary in unpredictable ways, exploiting the limited resources requires robust and effective load distribution techniques.

Load distribution, the placement of the operators in a query plan to machines (nodes) in the distributed system, is an important design decision, impacting the query processing performance [2]. A carefully selected operator placement may later produce poor performance due to time-varying, unpredictable stream fluctuations. Data fluctuations may in fact necessitate repeated expensive load redistributions in such distributed systems. Such load redistribution may cause delay and potentially make the system fail to react to short-term load fluctuations in time. Example 1 illustrates this problem using a real-world application.

Example 1. Consider a query monitoring stocks that exhibit “bullish” patterns (upward price movement in the stock market) [4] in recent business news and research.

```
SELECT S.company_name, S.symbol, S.price
FROM Stock as S, News as N, Research as R
WHERE matches(S.data, BullishPatterns) /*op1*/
AND contains(S.sector, News[1 hour]) /*op2*/
AND contains(S.company_name, News[1 hour]) /*op3*/
WINDOW 60 seconds
```

The lookup table *BullishPatterns* contains “bullish” patterns of stock behavior, e.g., “symmetrical triangle”. Operator op_1 performs a similarity-based join on the latest stock data tuples from the last 60 seconds with the *Patterns* table. Operators op_2 and op_3 perform matches on the stock sector and the company name with news and research streams. Let c_i and δ_i denote the processing costs and current selectivity of op_i respectively.

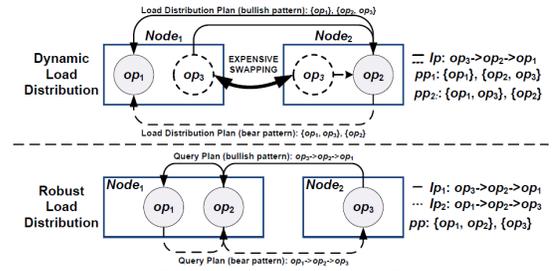


Fig. 1. Dynamic vs. Robust Load Distribution

Suppose it is a bullish market (i.e., stocks are doing well) and the following condition holds: $\delta_1 > \delta_2 > \delta_3$. Given these statistics, the best query processing order (query plan) is op_3, op_2, op_1 ($c_1 > c_2 > c_3$). Assume we have two machines, n_1 and n_2 , with resources r_1 and r_2 (i.e. CPU, memory and network bandwidth) available for processing. Then the best load distribution plan is op_1 on node n_1 , and op_2 and op_3 on node n_2 , as depicted in the upper illustration of Figure 1. Now suppose breaking news report poor stock performance. This will likely result in fewer matches with the *Patterns* table and instead more matches with news and blogs. In this case, δ_1 will be relatively lower than δ_2 and δ_3 for such data tuples. So plan op_3, op_2, op_1 is no longer the most efficient ordering. In fact, op_2 and op_3 may overload n_2 , namely $c_2 + c_3 > r_2$. Consequently, the DSPS has to relocate op_3 to n_1 . However, in the future if the market exhibits bullish pattern again, the optimizer might need to revert back to the original query processing order. Then the query executor would again need to move op_3 back to n_2 . In this particular scenario, we notice that the load distribution plan, op_1 and op_2 on n_1 , and op_3 on n_2 (lower Figure 1) would have been **robust** to support both query processing orders op_3, op_2, op_1 and op_1, op_2, op_3 in different scenarios. Proactive design of this load distribution plan achieves two objectives, i.e., robust query processing performance and tolerance to data fluctuations without dynamic load redistribution.

Insufficiency of State-of-The-Art. Traditional distributed and parallel systems [5], [6] categorize load distribution solutions as either dynamic or static. Dynamic load distribution [1] instead repeatedly improves the current load distribution plan by moving operators across machines to adapt to load changes. However, it comes with several drawbacks, including (1) amor-

tized overhead of repeated load redistributions, (2) adaptation delays, when redistribution opportunities could be missed, and (3) expensive maintenance of statistics.

Traditional static optimizers [7] determine a *single* “best” (i.e., cheapest overall execution costs) placement of query operators at compile time based on the average estimated statistics of data streams. While this approach imposes low optimization overhead, it cannot adapt the load distribution to changing statistics of data streams such as variations in input rates and selectivities. The most recent work, resilient operator distribution (ROD) [8], aims to produce a feasible physical plan to be resilient to time-varying and unpredictable workloads. However, shortcomings of ROD include: (1) it only focuses on the physical operator distribution without taking logical query plan ordering problem into consideration, (2) its operator distribution plan does not guarantee the given query processing performance, and (3) it is not proactive to changing workloads. A detailed comparison can be found in Section VII.

Our Proposed RLD Approach. Given these inevitable disadvantages of existing load distribution methods, we now propose an end-to-end solution called *Robust Load Distribution* (RLD) that exploits the key principles from load distribution methods while overcoming their respective shortcomings by integrating parametric query optimization. As foundation of RLD, we introduce the multi-dimensional parameter space model, a representation of the uncertainties in statistical estimates of streaming data. We then introduce the dual model of matching *robust logical* and *robust physical* plans which together gracefully handle the fluctuations experienced by the multi-dimensional parameter space.

Considering the intractable complexity of the search space composed of all combinations of logical plans with associated physical plan, we adopt a two-step query optimization approach. First, our proposed algorithm, Early-terminated Robust Partition (*ERP*), identifies a set of robust logical plans (i.e., *robust logical solution*) that together cover the parameter space with a *probabilistic guarantee* on the percentage coverage of the parameter space.

Next we propose a load distribution algorithm, *OptPrune* that efficiently produces an optimal *robust physical plan* (i.e., operator allocation plan) that supports the logical plans in a robust logical solution based on their probability of occurrence at runtime and their respective area of robustness in the parameter space. This results in a single robust physical plan tolerant to expected data statistic deviations at runtime.

Contributions. We introduce an end-to-end solution that overlays robust logical plans with a single physical plans in a parameter space. The contributions of this work can be summarized as follows:

1. We introduce the property of robust logical query plans (i.e., ϵ -robustness). We efficiently compute a multi-dimensional parameter space representing uncertainties in statistic estimates, including stream input rates and query operator selectivities.

2. *ERP*, our robust logical solution algorithm efficiently finds multiple robust logical plans that together assure coverage across the entire parameter space, unattainable by a single plan. *ERP* forms a *probabilistic bound* on the total uncovered area. Individual robust plans with any non-trivial area will be

found with high probability.

3. Given a robust logical solution, we design a family of algorithms that cover the spectrum from the optimization complexity to result optimality. *GreedyPhy* finds a robust physical plan in polynomial time, whereas *OptPrune* is a branch-and-bound algorithm using *GreedyPhy* as bound, which succeeds to significantly bound the search without compromising optimality.

The rest of this paper is organized as follows. We define the RLD problem in Section II, and overview our solution in Section III. Sections IV and V describe our algorithms for generating a robust logical solution and the associated robust physical solution. The experimental results are presented in Section VI. Sections VII and VIII discuss related work and the conclusion, respectively.

II. MODEL & PROBLEM STATEMENT

A. Basics of Distributed Query Plans

Common to other distributed stream work [1], [8], we assume the distributed stream processing system (DSPS) is deployed on a shared-nothing homogeneous compute cluster connected by a high bandwidth network. The DSPS accepts a continuous *query*, an expression describing the user’s information needs. Then it performs a two-step optimization [7], namely, plan generation and operator placement, to determine the most effective strategy to execute the given query. Plan generation identifies the algebra operator plan with the least estimated cost for a given input query and estimated data stream statistics. Operator placement takes the algebra plan as input and outputs a mapping of each operator in the algebra plan to a physical machine (node) in the cluster. We refer to the algebra plan and its operator placement as *logical* and *physical plans*, respectively.

B. Multi-dimensional Parameter Space

Current optimizers [1], [8] tend to use a single-point statistic estimate for plan generation. However, it is well known that estimates in streaming environments tend to fluctuate over time. We thus now model these uncertainties in estimates via a multi-dimensional space around these estimates, called *parameter space* S . This space captures all possible combinations of estimate variations. Each point *pnt* in space S is a vector $\langle d_1, \dots, d_n \rangle$, where each d_i is an estimate of the corresponding statistic modeled by that dimension of the space such as selectivity or input rate.

Different methodologies for constructing a parameter space exist, including strict upper and lower bounds [9] or levels of uncertainty to the optimizer estimates [10]. We use the latter approach in which the uncertainty level U is computed based on how statistic estimates E are derived. For example, if a value of E is available from the representative training data set, then $U = 1$ denotes low uncertainty. For simplicity we henceforth use an integer domain to denote the uncertainty levels, though other scales of uncertainty could be easily plugged in. The most crucial statistic estimates E for query optimization are the selectivities of operators and the input rates of streams [9]. We assume the statistic estimates E and the uncertainty level U correctly represent the data stream fluctuations. If suddenly some totally unexpected fluctuation

arises in the future, our current solution may not be able to handle it, and we may have to exploit operator migration to resolve such scenarios after all. The parameter space S is computed as shown in Example 2.

Example 2. Assume a simplified query $Q2$ of $Q1$:

```
Q2: SELECT *
FROM News as N, Stocks as S, Currency as C
WHERE N.subject = S.industry          /*op1*/
AND N.country = C.country             /*op2*/
WINDOW 60 seconds
```

Assume neither the selectivity for operator op_1 nor the input rate of stream News are accurate over time due to data fluctuations (e.g., breaking news in a certain industry). To capture this uncertainty, a parameter space is constructed around the single-point estimate $E = \{\delta_1 = 0.4, \lambda_N = 100 \text{ tuples/sec}\}$. First, an uncertainty level U (e.g., $U = 2$) is assigned to each estimate in E . Then the parameter space is constructed with δ_1 ranging between 0.32 and 0.48, and λ_N ranging between 80 tuples/sec and 120 tuples/sec.

Similar to parametric queries [11], in practice, each dimension of the parameter space is discretized. For ease of exposition, we henceforth work with a 2D parameter space, though the extension to higher dimensions is straightforward.

C. Notion of Plan Robustness

Let us now consider the most common queries, namely, select-project-join (SPJ) queries. Thus the cost model of a logical plan in a 2D parameter space is of the form:

$$\text{cost}(p, \text{pnt}) = c_1 \cdot \sigma_i + c_2 \cdot \sigma_j + c_3 \cdot \sigma_i \cdot \sigma_j + c_4$$

where c_1, c_2, c_3, c_4 are coefficients, and σ_i, σ_j represent the selectivities of operators op_i and op_j , respectively. Modeling a specific plan requires suitably choosing the four coefficients. This is achieved through standard surface-fitting techniques [11]. Extending the above equation to a general n -dimensional space is straightforward. Given the notations in Table I, we introduce the notion of robust query processing.

TABLE I
NOTATIONS AND DEFINITIONS

Term	Definition
n_i	The i th node
op_i	The i th operator
δ_i	The selectivity of op_i
r_i	Resource limit on node n_i
S	Parameter space
pnt	Parameter instance in the parameter space
pnt_{Hi}	Right top corner of a parameter space
pnt_{Lo}	Left bottom corner of a parameter space
lp	Logical query plan
LP_i	Robust logical solution for S , a subset of LP
pp	Robust physical plan
$\text{cost}(lp, \text{pnt})$	Cost of a query plan p at pnt
$\text{cost}(OP_i)$	Cost of a subset of OP on the i th node (i.e., the workload on the i th node)
lp_{pnt}^{OPT}	Optimal query plan at pnt

Definition 1. Given a parameter space S_i , a logical plan lp is ϵ -robust, also called *robust logical plan*, in S_i if its costs satisfy the condition (see Table 1 for notions):

$$\text{cost}(lp, \text{pnt}_{Hi}) \leq (1 + \epsilon) \times \text{cost}(lp_{\text{pnt}_{Hi}}^{OPT}, \text{pnt}_{Hi})$$

The intuition of what a robust logical plan is can be explained as follows: consider two optimal query plans lp_1 and lp_2 for two points $\text{pnt}_{Lo} = \langle d_{1,1}, \dots, d_{1,n} \rangle$ and $\text{pnt}_{Hi} = \langle d_{2,1}, \dots, d_{2,n} \rangle$ in the parameter space S_i respectively, and $\text{pnt}_{Lo} < \text{pnt}_{Hi}$, meaning $\forall i d_{1,i} \leq d_{2,i}$. If lp_1 is a robust plan in S_i , then we can provably bound the costs of plan lp_1

between the costs of plan lp_1 at pnt_{Lo} and the costs of plan lp_2 at pnt_{Hi} .

Definition 2. The *robust region* of a logical plan lp in S is the subarea S_i where lp satisfies Definition 1 at all points $\text{pnt}_i \in S_i$. For example, the robust region of lp_1 in Figure 2 (left) covers the left-top corner of the space S .

Definition 3. Given a set of robust logical plans, also called robust logical solution LP_i and resources r_i for node n_i ($\forall i : 1 \leq i \leq N$) for a DSPS, a physical plan pp is *robust*, also called *robust physical plan*, if it satisfies the following conditions: 1) $\text{cost}(OP_i) \leq r_i$, 2) $\bigcup OP_i = OP$, and 3) $\bigcap OP_i = \emptyset$ ($\forall i : 1 \leq i \leq N$), where OP_i denotes the set of operators allocated to machine n_i by pp , and OP is the full set of operators in the query.

Intuitively, a physical plan pp , shown in Figure 2 (right), is robust, if for each subset of query operators OP_i assigned to node n_i , its total cost $\text{cost}(OP_i)$ is no greater than the resource capacity r_i of n_i to execute the sub-plans of all query plans $lp_i \in LP$. Each OP_i associated with n_i also defined as a *configuration* c_i , has no overlap with any other OP_j , and the union of all OP_i forms the whole operator set OP .

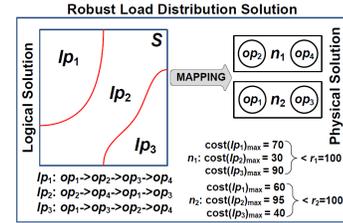


Fig. 2. Robust Load Distribution Solution

D. Problem Statement

Robust distributed query processing aims to 1) identify a robust logical solution LP_i , such that for each point in the parameter space S , there is at least one logical plan lp_i in the solution LP_i that is ϵ -robust by Def. 1 in that sub-space S_i , and to 2) produce a physical plan pp that supports the robust logical solution, i.e., pp is robust by Def. 2. The key idea is that the resulting system will be able to withstand the known data stream fluctuations, meaning, as long as the actual statistics (i.e., input stream rates and selectivities) remain within the parameter space. Our robust load distribution (RLD) problem can be formalized as follows.

Given a query q , resources r_i for node n_i ($\forall i : 1 \leq i \leq N$), statistic estimates $E = \langle e_1, \dots, e_k \rangle$, and the associated uncertainty levels $U = \langle u_1, \dots, u_k \rangle$, the **robust load distribution problem** is to find a robust physical plan pp that supports a robust logical solution LP_i in the parameter space S constructed based on E and U .

Finding the robust solution (Figure 2) for this problem requires a comparison among all possible subsets of all logical plans LP and all possible physical plans pp in PP in the worst case [8]. The above problem is prohibitively expensive as the search space of finding robust logical and physical plans are exponential in the number of query operators and the number of machines in the system, respectively [8]. Furthermore, finding a robust logical solution LP_i and a corresponding physical plan pp supporting LP_i requires the optimizer to search both logical and physical search spaces. This renders the solution intractable for large problems.

III. OVERVIEW OF RLD APPROACH

Given the intractability of RLD, we instead employ a two-step approach towards query optimization popular for both distributed and parallel database systems [7] due to its reduction of the overall complexity of the distributed query optimization problem.

In our context, the two-step optimization works as follows: 1. The first step generates a robust logical solution, in which each logical plan is designed for a particular sub-region of the parameter space. 2. The second step produces a single robust physical plan, determining the machine on which each operator is placed and supporting the given logical solution without load redistribution.

The above two steps efficiently work together to achieve an effective load distribution plan as our experiments confirm. Our overall RDL strategy conducts load distribution with full awareness of all possible scenarios at runtime, thus avoiding costly load redistributions. The RLD architecture consists of the following three components.

Robust plan optimizer. Our robust plan optimizer uses the standard query optimizer of a DSPS as a black-box to perform traditional plan optimization calls. Given a query, it estimates resource requirements of the query based on the uncertainty of the estimates. The plan optimizer then produces a *robust logical solution*, and determines a single robust physical plan to support this logical solution.

Robust load executor. Our robust load executor is built on top of a multiple-plan adaptive query executor [12] that offers the ability of switching between robust logical plans at runtime. This is accomplished by means of an online classifier that associates each computed robust logical plan with its associated specific statistics. Once a robust load distribution solution has been produced by the optimizer, the operators are instantiated by the executor on the machines accordingly. At runtime, our executor collects up-to-date statistics of running operators from statistic monitors installed on the DSPS machines. The classifier then inspects the runtime statistics to determine the best logical plan from the robust logical solution for the new batch of tuples.

Statistic monitor. The robust load executor requires runtime knowledge about the actual values of key parameters. Thus each machine in a DSPS runs a statistic monitor that periodically samples the selectivity of operators and the associated stream input rates. The monitor then transmits the statistics to the executor where all statistics are kept up to date.

IV. ROBUST LOGICAL PLAN GENERATION

A. Weighted Partition Algorithm Overview

The parameter space, defined in Section II, requires us to find the robust logical solution LP_i such that each plan $lp_i \in LP_i$ satisfies Def.1. Our approach leverages the fact that we are not required to identify the *optimality region* for each plan, but rather only its *robustness region*. As our experimental study validate, this reduces the space significantly.

Consider a two dimensional parameter space containing 6 distinct robust logical plans, with their respective robust regions depicted in Figure 3(a). In this example, the parameter space would only need to be partitioned twice (making 10 optimizer calls) to confirm the robustness of all 6 plans. The

resulting partitioned space will contain 10 sub-spaces (Figure 3(a)). On the contrary, the exhaustive approach depicted in Figure 3(b) would divide the parameter space into an 8×8 grid (making 64 optimizer calls) to discover the same 6 plans. Thus, the exhaustive approach is 6 times more expensive than needed. Worst yet, such overhead increases significantly with the number of dimensions of the parameter space.

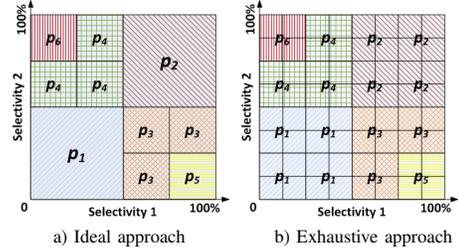


Fig. 3. Comparison of Ideal and Exhaustive Approach

Overview of the Logical RDL Steps. The first technical challenge addressed by our logical plan generation algorithm is how to efficiently find an effective partition point in the parameter space (Step 1). By the robust logical plan definition, to verify the robustness of a logical plan, we must make expensive optimizer calls in all sub-spaces of its partition. If that tested plan does not satisfy the robustness criteria, a finer-grained partition must be found. Therefore, identifying a good partitioning point is the key to avoid repeated wasteful attempts of expensive robust logical plan finding. Our insight is to leverage information about the already known query plans in the space and encode this knowledge as weights in the space. In this model, points where a new robust plan is more likely to exist are assigned higher weights.

The second issue we tackle is how to update the weights assigned to all points during the parameter space partitioning process (Step 2). Given the size of the space, we propose an efficient approach that incrementally updates the weights.

The third issue we address is when to stop partitioning the parameter space (Step 3). Fine-grained partitioning incurs significant computational overhead for query optimization. Moreover, the quality improvement of the resulting robust logical plans may no longer outweigh its growing expense of optimizer calls. Thus, we develop a termination condition that provides a probabilistic guarantee on the space coverage by the robust logical plans. We show that the possibly missed robust plans, if any, are guaranteed to not occupy a large area in the parameter space.

B. Weight Assignment in Parameter Space

We now describe our strategy of exploiting plan cost information from already identified plans to assign weights to the remaining points in the parameter space. The weight of a point should reflect the probability that the robust logical plan at that point is different from the robust plans in the bottom-left and top-right corners of the associated space. This way we would be able to exploit these weights to efficiently identify distinct robust plans in the parameter space. However, this is clearly infeasible without knowing what the costs of previously identified plans would be at all points in the space. Yet computing all these plan costs is prohibitively expensive in a large space. Thus, our goal instead is to heuristically approximate the weights without exhaustively computing all

plan costs in the space. For this, we develop a weight function based on the following principles.

Principle 1. Two points in the parameter space that are closer to each other are more likely to have the same robust plan compared to a more distant pair of points. To motivate this, consider the cost model in Section II-C. We observe that the cost of a plan is monotonically increasing along each dimension of the space [11]. Thus, in a small region of the space, the costs of a plan lp_i on two nearby points are likely to have the same robust plan by Def.1.

Principle 2. A plan is less likely to be robust at a point if the slope of the plan’s cost function, defined in Section II-C, at that point is high. This principle is based on the observation that the slope of a plan’s cost function is monotonically increasing in the parameter space. Intuitively, the closer the points move towards the margin of the plan’s robust region, the higher the slope of the plan cost function is at these points.

Based on these two principles, we now design a **weight assignment function** that increases as the ratio between the slope of a plan’s cost function and the distance of that point to the bottom left point pnt_{Lo} of that space increases. The weight function decides how quickly the weight decreases as the distance increases, and how quickly the weight increases as the slope increases. In practice, assigning weights individually to each point in the parameter space would be expensive since there are $O(n^d)$ points in a d -dimensional space assuming an n step discretization of the space along each dimension. Thus, we consider each dimension independently. A point $pnt = (x_1, \dots, x_n) \in S$ is projected onto each dimension d_i , such that $x_i \in d_i$ and the point’s weight on each dimension is assigned according to the projected distance between pnt and pnt_{Lo} (see Table I). The weight of pnt in the i -th dimension, denoted by $weight_i(pnt)$, is defined as:

$$weight_i(pnt) = \frac{\min(\text{slope}(pnt, p_{pnt_{Hi}}^{OPT}), \text{slope}(pnt, p_{pnt_{Lo}}^{OPT}))}{\text{dist}(pnt, pnt_{Lo}^i)}$$

where dist is a function that measures the distance between two points. Any distance measure such as Manhattan or Euclidean Squared Distance [13] could be plugged in.

Weight Re-Assignment Strategy. Unfortunately, the initial weight assignment will no longer be accurate after partitioning. That is partitioning produces several sub-spaces, each of which may have their own optimal plan at bottom-left $lp_{pnt_{Lo}}^{OPT}$ or top-right corner $lp_{pnt_{Hi}}^{OPT}$ of that sub-space. Recall that the optimizer finds a distinct robust plan for each sub-space. Thus, each point’s weight has to be updated accordingly in order to reflect the cost behavior of the new plans in the sub-spaces.

As described earlier, the cost of the weight assignment function largely depends on the number of points in the parameter space. Updating the entire parameter space repeatedly incurs significant overhead for the weight assignment. We now introduce a refinement of the above weight assignment approach where we update the weights of points in a sub-space as we partition the parameter space S if and only if the following condition is satisfied.

$$\overline{(lp_{pnt_{Lo}} = lp_{pnt_{Lo}}^{OPT}) \wedge (lp_{pnt_{Hi}} = lp_{pnt_{Hi}}^{OPT})} = \text{True}$$

where $weight'_i(pnt)$ denotes the updated weight assignment for the current partition point pnt ; while $S_i.pnt_{Hi}(S_i.pnt_{Lo})$ denotes the top-right (bottom-left) corner of the sub-space

$S_i \in S$. Intuitively, the above condition ensures that the update would not be triggered if the predicted logical plan at top-right (bottom-left) corner is identical to the optimal plan identified at the same point after partitioning the space.

Example 3. We illustrate the above weight re-assignment strategy using Figure 4. The weight is assigned to each point in the original parameter space in Figure 4(a). Then we partition the space into 4 sub-spaces, and compute the weights of pnt_{Hi} and pnt_{Lo} for each of the 4 sub-spaces S_i (Figure 4(b)). However, we only update the sub-spaces S_4 (Figure 4(c)) as the optimal plan is different from the predicted robust logical plan at $S_4.pnt_{Hi}$. Thus, the above condition reduces the number of sub-spaces to update.

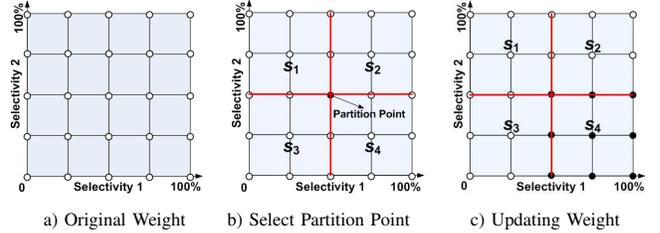


Fig. 4. Weight Assignment for Parameter Space

Example 3 shows savings by not updating weight assignments in areas where an identified robust logical plan has a "large" area of optimality. The reason is that the weight assignment in a sub-space only depend on its robust plan.

C. Parameter Space Partitioning

1) *Weight-driven Robust Partitioning Algorithm:* We now illustrate how the above weight assignment is integrated into the parameter space partitioning process. The main idea of parameter space partitioning is to incrementally find *robust* plans as the space is partitioned. A straightforward way to achieve this is using the *weight-driven partition (WRP)* procedure (Algorithm 1). First we compute the weights for all points in the discretized space S . Then, we pick the point with the highest weight as the partition point to divide the space into 2^d sub-spaces, with d denoting the number of dimensions of S . Each sub-space S_i will be further partitioned if its optimal plan at pnt_{Hi} is not *robust*. Finally, the partitioning process stops once all plans at pnt_{Hi} of all S_i are *robust*.

This partitioning process could result in unnecessarily small sub-spaces if the robustness threshold ϵ is not well chosen. On the contrary, a too large ϵ may end up with a rather small number of sub-spaces - leading to suboptimal plans. Our experiments (Section VI-C) show that relatively small increments in ϵ are sufficient to bring down the number of plans significantly, without adversely affecting the query processing quality.

Limitations of WRP. As described above, *WRP* aims to minimize the computational overhead of partitioning. Yet there is a significant explosion in costs associated with an increasing dimensionality of the parameter space. Moreover, a large percentage of the computation could be wasted as *WRP* treats all sub-spaces in the parameter space equally rather than targeting specific subareas. The intuition is that a too *strict* threshold ϵ may lead the partitioning algorithm to undertake a too fine-grained job, not merited by the coarseness of the

underlying space. In an extreme case, all optimal logical plans in the space must be identified if $\epsilon = 0$.

Algorithm 1 Weight-driven robust partitioning algorithm

Input: A parameter space S

Output: A robust logical solution LP_i

```

1: Plan  $lp \leftarrow \text{optimize}(S.pnt_{Hi})$ 
2: if  $lp$  is robust in  $S$  then
3:    $LP_i.add(lp, S)$ 
4:   return  $LP_i$ 
5: else
6:   Assign weights to points in  $S$ ;
7:   Choose appropriate point to partition  $S$  based on weights
8:   for  $i = 1 \rightarrow n(\text{number of sub-spaces})$  do
9:     Update weight assignments in  $s_i \in S$ 
10:     $\text{standardPartitioning}(s_i, LP_i)$ ;
11:   end for
12: end if
13: return  $LP_i$ 

```

2) *Early-terminated Weight Robust Partitioning*: Given the shortcomings of *WRP*, we now enhance this approach by designing an early termination strategy. This new method, called *ERP*, reduces the overhead while providing a *probabilistic guarantee that the robust plans missed cannot occupy a large area in the parameter space*. Our solution uses two key factors:

(a) **Region of robustness** for a given plan lp refers to the location where the plan lp is robust in the parameter space.

(b) **Size of robustness** for a given plan lp corresponds to the total area in the parameter space where the plan lp is robust.

We exploit the above two factors to trade off between the partitioning costs and the quality of the resulting plans. Observe that when we partition the parameter space using uniform partitioning, the probability of finding a new robust plan is proportional to its area of robustness. Given that a finite number of robust plans exists, the probability of finding a new robust plan decreases as we find more robust plans from that set of plans when partitioning. The more robust plans we have already found, the more partitioning steps it will take on average to find an additional robust plan. We propose to exploit this insight by terminating the partitioning process when we have not obtained a new robust plan for a pre-determined number of partitioning steps.

For this, we maintain an *aging counter*, which keeps track of the interval between the last two times that a new robust plan was detected in the partitioning process. Each time we call the optimizer at pnt_{Hi} of a new sub-space, if the plan at pnt_{Hi} is a new robust plan that is distinct from all robust plans observed thus far, we reset the aging counter. Otherwise, we increment the aging counter.

Assume that after partitioning t times, we found a new robust plan. Then the aging counter is reset to 0 and we start counting up again. Let c be the number of additional partitioning steps after t to find the next new robust plan. The probability of finding a robust plan is constant between t and $t + c$, since the number of missing robust plans does not change during this interval. If we denote the probability of identifying a new robust plan after t partitions by $Pr(t)$,

$$\forall 0 \leq t \leq c - 1, Pr(t + c) = \frac{n_{miss}}{n_{total}}$$

where n_{miss} is the number of unidentified robust logical plans and n_{total} is the total number of robust plans.

Theorem 1: With a probability of at least $1 - \epsilon$, if we do not find a new optimal plan in the partitioning process within c trials, where $c \leq c_0 = (1 + \epsilon^{-1/2}) / \delta$, then the total number of optimal plans not yet found is bounded by δ ($\leq \delta$).

Details of the proof can be found in our technical report [14]. Intuitively, Theorem 1 states that if the aging counter reaches a value $\geq c_0$, then with high probability the missing optimal plans cover a small area in the parameter space. We observe that the position to partition the space can have a significant impact on how quickly the aging threshold is reached. The space partitioning technique, which exploits Theorem 1 to early terminate the partitioning process, called *ERP*, is shown in Algorithm 2. Since the guarantee in Theorem 1 is probabilistic, our proposed *ERP* may miss some robust plans. The following theorem quantifies the likelihood of missing a robust plan based on its area of optimality.

Algorithm 2 Early-terminated Robust Partitioning Algorithm

Input: A parameter space S

Output: A robust logical solution LP_i

```

1:  $LP_i \leftarrow \phi$ 
2:  $counter_{miss} \leftarrow 0$ 
3:  $age\_threshold \leftarrow (1 + \epsilon^{-1/2}) / \delta$ 
4: while  $counter_{miss} \leq age\_threshold$  do
5:    $pnt \leftarrow \text{getPartitionPnt}(S)$ 
6:    $p_{pnt}^{OPT}$  is the optimal plan at  $pnt$ 
7:   if  $p_{pnt}^{OPT} \in LP_i$  then
8:      $counter_{miss}++$  and continue
9:   else
10:     $add(p_{pnt}^{OPT})$  to  $LP_i$ 
11:     $counter_{miss} = 0$ 
12:   end if
13: end while
14: return  $LP_i$ 

```

Theorem 2: Suppose we stop partitioning according to the aging threshold in Theorem 1. If the coverage of an optimal plan lp is $area(lp) \geq \gamma \cdot \delta$, for a constant $0 < \gamma \leq 1/\delta$, then the probability that the plan lp is not found is at most $e^{-\gamma(1+\epsilon^{-1/2})}$.

Details of the proof can be found in our technical report [14]. Theorem 2 states that if an optimal plan has a "large" area of optimality, then we will find it with high probability when using the stopping condition of Theorem 1. From Theorem 2 we know that the probability of missing a robust plan decreases exponentially with its area.

V. ROBUST PHYSICAL PLAN GENERATION

A. Basic Approach for Robust Physical Plan

A straightforward strategy for robust physical plan generation would entail the following steps. As input, we accept LP_i , a set of logical plans that together cover the parameter space by Def.1 produced by the logical optimizer. Then we compute all physical plans PP for each robust plan $lp_i \in LP_i$, denoted by $PP(lp_i)$. Thereafter, we find the intersection among all sets $PP(lp_i)$ for all logical plans $lp_i \in LP_i$. If the intersection is

not empty, then any physical plan in this intersection is a valid solution, which satisfies all robust logical plans in LP_i .

However, if the intersection is empty, then no physical solution supports *all* logical plans in LP_i . We would need to locate a suboptimal solution. One simple option would be to remove a logical plan lp_i from solution LP_i . Then we repeat the above procedure until the intersection is not empty and thus a valid robust physical plan can be produced.

Unfortunately, the number of physical plans for a single logical plan is $n^m/n!$ [8], where n is the number of machines and m the number of operators in the logical plan lp . Moreover, the number of different logical solutions LP_i is $2^k - 1$, where k is the cardinality of the set of all possible logical plans LP . As a result, finding the optimal solution for this problem is intractable for a large number of machines or a large number of robust logical plans.

Thus, we now propose algorithms that trade off between the optimization complexity versus the result optimality. *GreedyPhy* exploits heuristics to efficiently find a robust physical plan in polynomial time, whereas *OptPrune*, using *GreedyPhy* as baseline, is guaranteed to find the optimal physical plan to support the maximum number of logical plans, though with minimal increase in optimization time.

B. Greedy Physical Plan Generation

Given the complexity of physical plan generation, we now introduce a heuristic strategy that uses two key principles:

(a) **The area of optimality heuristic.** Intuitively, we aim to produce a robust physical plan pp that covers as much as possible of the parameter space by supporting all logical plans in LP_i . If all logical plans cannot be supported by pp , then we drop from LP_i the least important logical plan, i.e., the plan associated with the *smallest robust region* in the parameter space S from LP_i .

(b) **The probability of occurrence heuristic.** By definition, the selectivity of an operator fluctuates around its given statistic estimates. Various probability distributions could be used to model the occurrence of a point in the space S . For simplicity, we model this probability using a *normal distribution* as commonly done in the literature [13]. Therefore, the closer a point is to the given statistic estimate, the higher the possibility that the actual selectivity happens at runtime. As a result, we drop the logical plan whose optimality region is the furthest away from the given estimate point.

Weight Assignment Policy. Our strategy is to assign a weight to each robust logical plan that incorporates the above two factors. Let $area(lp_i)$ denote the robust region of plan lp_i , and $Pr(pnt_j)$ denote the probability of occurrence of a point pnt_j at run time. A robust logical plan's weight $weight(lp_i)$ is defined as:

$$weight(lp_i) = \sum_{pnt_j \in area(lp_i)} Pr(pnt_j)$$

where $Pr(pnt_j)$ can be obtained from the normal distribution.

Example 4. Consider a 2-dimensional parameter space with each dimension discretized into 16 units (see Figure 5). Suppose the space contains 5 different robust logical plans. Each plan's robust region is depicted as one or more rectangles. For example, the robust region of lp_1 includes $area_1$, $area_6$, $area_{11}$, and $area_{16}$. The probabilities of the actual run time statistics to fall within these rectangles are 2.4%, 11.7%,

11.7%, and 2.4%, respectively. Summing the values gives us the weight of lp_1 as 28.2%. The probability of occurrence with respect to a unit area is calculated as follows using $area_1$ as an example:

$$Pr(area_1) = Pr_x(area_1) \cdot Pr_y(area_1)$$

where x and y denote the x -axis and y -axis of a unit area in the 2-dimensional parameter space, respectively. The above assumes that the x and y dimensions are independent following the assumption of independence of selectivity values commonly made by query optimizers [15]. Thus the correlation between dimensions is zero.

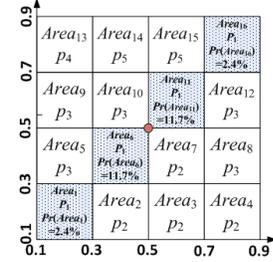


Fig. 5. Weight Assignment for Logical Plans

Algorithm 3 GreedyPhy Algorithm

Input: A robust logical solution LP_i , resources R

Output: A robust physical plan pp

```

1: terminate ← false
2:  $lp_{max} \leftarrow updateMax(LP_i)$ 
3: while terminate ≠ true do
4:    $pp \leftarrow LLF(lp_{max}, R)$ 
5:   if  $pp \neq null$  then
6:     terminate ← true
7:   else
8:      $index \leftarrow getMinWeightPlanWithMaxOp()$ 
9:      $LP_i.remove(index)$ 
10:     $lp_{max} \leftarrow updateMax(LP_i)$ 
11:   end if
12: end while
13: return  $pp$ 

```

The GreedyPhy Algorithm. After assigning weights to the logical plans, the plans are stored in a heap sorted by their weights in descending order. *GreedyPhy* algorithm exploits the above weight assignment (Algorithm 3). Given a solution LP_i produced by the logical plan optimizer, then each plan $lp_i \in LP_i$ has a weight $w(lp_i)$, *GreedyPhy* finds the physical plan pp that supports a subset of LP_i such that the sum of all weights of the supported logical plans is maximal among all possible subsets of LP_i . Intuitively, the resulting robust physical plan pp should be robust to the most frequently occurring data fluctuations. Given a robust logical solution LP_i and resource constraints r_i , *GreedyPhy* generates a logical plan lp_{max} , in which the cost of each operator is equal to its maximum cost for all logical plans $lp_i \in LP_i$ (Lines 1-2). Thereafter, in each iteration the algorithm tries to produce a physical plan by using the Largest Load First (LLF) algorithm (i.e., Longest Processing Time algorithm [8]) (Line 4). LLF orders the operators by their cost and assigns operators in descending order to the least loaded machine. If a physical plan is produced by LLF, it is thus returned as final solution.

If the algorithm fails to find a physical plan, it removes the least-weighted logical plan lp_i from the robust logical solution LP_i (Lines 8-10). After repeating lines 3-10 a maximum of $|LP_i|$ times (LP_i is empty after $|LP_i|$ times), the algorithm returns a physical plan that maximizes the total weight of the subset of logical plans selected from LP_i supported by pp .

Complexity Analysis. The worst case for *GreedyPhy* is that none of the logical plans in LP_i can be supported by the given resources. In other words, it would iterate k (the cardinality of $|LP_i|$) times before stopping. Therefore, the worst case complexity of *GreedyPhy* is $O(n)$ as the complexity of LLF is proven to be $O(n)$ [8], and the complexity of *getMinWeightPlanWithMaxOp* and *updateMax* procedures are both linear in the number of operators in lp , namely $O(n)$. Therefore, *GreedyPhy* is guaranteed to have a linear complexity taking $O(n)$ time.

C. OptPrune Physical Plan Generation

The above algorithm, being greedy, cannot always find the optimal solution. We now design a strategy that guarantees the optimal solution is found if one exists. Given the prohibitively high complexity of exhaustive search, the key idea is to devise a pruning methodology that eliminates suboptimal solutions. *OptPrune* succeeds in improving the efficiency of the search costs without compromising result optimality (see Algorithm 4).

Algorithm 4 OptPrune Algorithm

Input: A set of robust logical plans LP_i , resource limits R

Output: A robust physical plan pp

```

1:  $C \leftarrow$  all feasible configurations on a single machine
2:  $greedyPlan \leftarrow GreedyPhy(LP_i, R)$ 
3:  $bound \leftarrow score(greedyPlan)$ 
4:  $pp \leftarrow NULL$ 
5:  $sort(C)$ 
6:  $c_0 \leftarrow max(C)$ 
7:  $Search(C, c_0)\{$ 
8:  $S' \leftarrow removeConflict(C, c_0)$ 
9: for all  $c_i \in S'$  do
10:  $pp.add(c_i)$ 
11:  $score \leftarrow updateScore(pp)$ 
12: if  $completeSolution(pp)$  then
13:   return true
14: else if  $checkLimit(pp) \vee score < bound$  then
15:    $pp.remove(c_i)$ 
16:    $score \leftarrow updateScore(pp)$ 
17:   continue
18: else if  $!Search(C', c_i)$  then
19:    $pp.remove(c_i)$ 
20:    $score \leftarrow updateScore(pp)$ 
21:   continue
22: end if
23: end for
24:  $\}$ 
25: return  $pp$ 

```

In order to find the optimal solution, *OptPrune* potentially needs to examine all possible physical plans. We represent all physical plan candidates in a weighted directed graph $G = (V, E, SCORE)$ (Figure 6) where $v_i \in V$ are vertices,

$e_{ij} \in E$ are directed edges $v_i \rightarrow v_j$, and a $score \in SCORE$ is associated with each vertex v_i , denoted by $score(v_i)$. A vertex v_i represents a set of configurations c_i defined in Section II-C. The *root*, the vertex on level 0, is empty. Each *leaf*, a vertex on level N (with N the number of machines), is a *robust physical plan* (Def.3). All vertices located between level 0 and level $N - 1$ correspond to partial physical plans (i.e., $\bigcup OP_i \sqsubset OP$ but $\bigcup OP_i \neq OP$). A vertex v_j is the *child* of the vertex v_i denoted by the edge e_{ij} from v_i to v_j . The configurations in a child vertex v_i corresponds to the union of all configurations in v_i 's parents. A *score* is the minimal weights of the configurations can support in a vertex. For example, $c_2 = 0.6$ and $c_4 = 1$ in Figure 6, the score of $v_8 = \min\{c_2, c_4\} = 0.6$.

The key idea is that *OptPrune* can leverage the results generated by *GreedyPhy* as its pruning criteria for improved efficiency. For simplicity, we here assume that the machines are homogeneous. Thus a configuration such as $c_2 = \langle op_1, op_3 \rangle$ is valid, if op_1 and op_3 can fit on one machine.

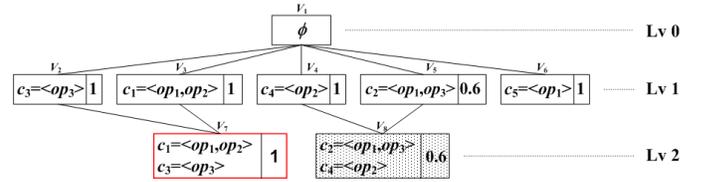


Fig. 6. Robust Physical Plan Search Graph

OptPrune traverses the above search tree in a depth-first search (DFS). *OptPrune* starts at the root, an empty plan, and continues iteratively down the graph forming a robust physical plan by adding configurations. As stated in Section V-A, the goal of *OptPrune* is to maximize the total score of the logical plans that a physical plan pp can support, denoted by $score(pp)$. The algorithm first figures out all possible configurations on a single machine (Line 1). The *score* of a physical plan is the total weight of the logical plans being supported. It sorts the configurations in decreasing order of the number of operators in each configuration. c_0 is the configuration with the most number of operators (Lines 2-6).

OptPrune starts its depth-first search (DFS) by adding one configuration to the current robust physical plan (pp) at a time (Lines 7-11). If the union of configurations in pp contain all query operators, then the algorithm terminates and returns pp (Lines 12-13). We have an effective bound that is guaranteed to only eliminate suboptimal solutions. If the current partial pp exceeds the given resource capacity or the *score* of the current pp is worse than that of the *GreedyPhy* solution, then the algorithm backtracks by removing the newly added configuration from the current pp and updating the *score* of pp accordingly (Lines 14-21).

Theorem 3: The new pp^* cannot be an optimal solution if the current pp is not an optimal one. The physical plan search graph is guaranteed to be safely pruned. *OptPrune* is guaranteed to find the optimal solution.

Proofs of the above lemma and theorem can be found in our technical report [14].

Complexity Discussion. The worst case for *OptPrune* is to have to check every configuration in the entire search space.

Therefore, the worse case complexity of *OptPrune* is the same as that of exhaustive search, namely, $O(n^m/n!)$. However, in practice our proposed pruning methods are found to be extremely effective at terminating the search much earlier, as confirmed by our experimental results (Section VI-D). The reasons are twofold. First, *GreedyPhy* produces a relatively good physical plan - this offers us an excellent bounding condition, which enables us to stop searching through most branches early on. Hence it reduces the search space significantly without affecting the search accuracy. Secondly, *OptPrune* terminates immediately if it finds the first physical plan (leaf) that supports all robust logical plans.

VI. EXPERIMENTAL STUDY

We have implemented the proposed techniques on D-CAPE [16], a distributed continuous query processing architecture employing stream query engines over a cluster of shared-nothing processors. The experiments were run on the D-CAPE system using Linux machines with Intel 2.6GHz Dual Core CPUs and 4GB memory.

TABLE II
SYSTEM PARAMETERS & DATA DISTRIBUTION

Parameter	Value	Description
<i>Data Arrival</i>	<i>Poisson</i>	Data arrival distribution
μ	500 msec	Mean inter-arrival rate
$ T_{dq} $	1,000	Maximum # of tuples dequeued by an operator at a time
<i>Ruster size</i>	100 tuples	Minimum <i>ruster</i> size
Data Distributions		
Uniform ($\alpha = 0, \beta = 100$): <i>min</i> : 0.0, <i>max</i> : 100.0, <i>med</i> : 49.0, <i>mean</i> : 49.7, <i>ave.dev</i> : 25.2, <i>st.dev</i> : 29.14, <i>var</i> : 849.18, <i>skew</i> : 0.05, <i>kurt</i> : -1.18.		
Poisson ($\lambda = 1$): <i>min</i> : 0.0, <i>max</i> : 7.0, <i>med</i> : 1.0, <i>mean</i> : 0.97, <i>ave.dev</i> : 0.74, <i>st.dev</i> : 1.01, <i>var</i> : 1.02, <i>skew</i> : 1.17, <i>kurt</i> : 1.89		

A. Data Sets and Queries

Stocks-News-Blogs-Currency data set: We have employed a data polling application (implemented in QueryMesh [12]) that collects NYSE stock prices, foreign currency exchange rates from Yahoo Finance, news and blogs via RSS feeds.

Sensor data set: This data set contains readings from sensors in the Intel Research, Berkeley Lab [9]. The sensor readings are streamed to D-CAPE in the order they are generated, as if they were submitted by sensors in real-time.

Synthetic data sets: To study the effectiveness of our strategies under data fluctuations, we design several data sets using various data distributions that model real-life phenomena. Default properties, distribution and system parameters are depicted in Table II.

Queries: We deploy N-way join queries, as those are among the core and most expensive queries in database systems. The default settings used in our experiments are listed in Table II. The queries are equi-joins of 10 streams.

B. Experimental Methodology

Our experiments are categorized into three major classes.

The first class studies the effectiveness of our *ERP* algorithm for **robust logical plan generation**. Specifically, we compare the *compile-time optimization performance* and the *quality* of the resulting robust logical plans for three alternative techniques. As baseline for the best quality robust logical solution, we employ exhaustive search (*ES*) over the discretized parameter space. We also implement a search algorithm which randomly selects points in the parameter space as plan candidates (*RS*). *RS* stops making optimizer

calls if it fails to find a distinct robust logical plan after a given number of optimizer calls. This corresponds to our partitioning technique assigning equal weights to all points in the parameter space. Finally, our weight assignment strategy with early termination is denoted as *ERP*.

The second class evaluates the effectiveness of our algorithms, *GreedyPhy* and *OptPrune*, for **robust physical plan generation**. Specifically, we compare different approaches for physical plan generation with respect to their *optimization time* to find the solution and the *space coverage* of the solution. We also measure the total weight of the area covered by the resulting physical plan. As baseline, we again choose the results from the exhaustive search over all load distribution plans, which is guaranteed to find the optimal solution.

The third class is a comparative study assessing the **runtime execution** of the overall RLD system. Specifically, we evaluate the average tuple processing time and runtime overhead of our RLD solution and compare it to state-of-the-art approaches, namely, the resilient load distribution [8] (*ROD*) and dynamic load redistribution (*DYN*) [1].

C. Effectiveness of Logical Plan Generation

Varying Robustness Thresholds and Uncertainty Levels.

This experiment assesses the impact of the robustness threshold ϵ and uncertainty level U on the effectiveness of our **logical plan finder**. The value ϵ of the robustness threshold is varied from 10% to 30% for $Q1$ (5-way join query). Figure 7 shows the number of optimization calls made by *ERP*, *RS* and *ES*. As expected, a lower value for ϵ (tighter robustness bound) results in a higher number of optimization calls, because returned plans cannot be much worse than the corresponding optimal logical plans due to the tight robustness bound ϵ . Furthermore, Figure 7 depicts the optimization efficiency under different uncertainty levels. The higher uncertainty levels result in a larger parameter space. Hence, the number of optimizer calls increases along with the increasing uncertainty level. The results, showing the parameter space S coverage of the robust logical plans identified by *ES*, *RS* and *ERP*, can be found in our technical report [14] due to the space limit.

Varying the Number of Dimensions. Our previous results are based on a fixed number of statistic estimates (i.e., dimensions). We now examine the relative efficiency of the algorithms for dimensions varying from one to ten. $Q2$ (10-way join query) is used because it has a higher number of logical plans compared to $Q1$. Thus, it is more likely to suffer from the exponential growth of complexity with a linear increase in the number of dimensions.

We consider 3 parameter configurations to evaluate the efficiency of our algorithm for finding logical plans across diverse scenarios. Our optimizer, as shown in Figure 9, is significantly more efficient than the alternative approaches. It is clear that the more dimensions the parameter space has, the more subspaces are produced by each partitioning step. That is, the number of subspaces grows exponentially with the dimensionality of the parameter space. This issue drastically affects *ES* because this approach has to check all unit subspaces in the discretized space. As depicted in Figure 9, the number of partitioning iterations increases exponentially with the number of dimensions. In contrast, our proposed

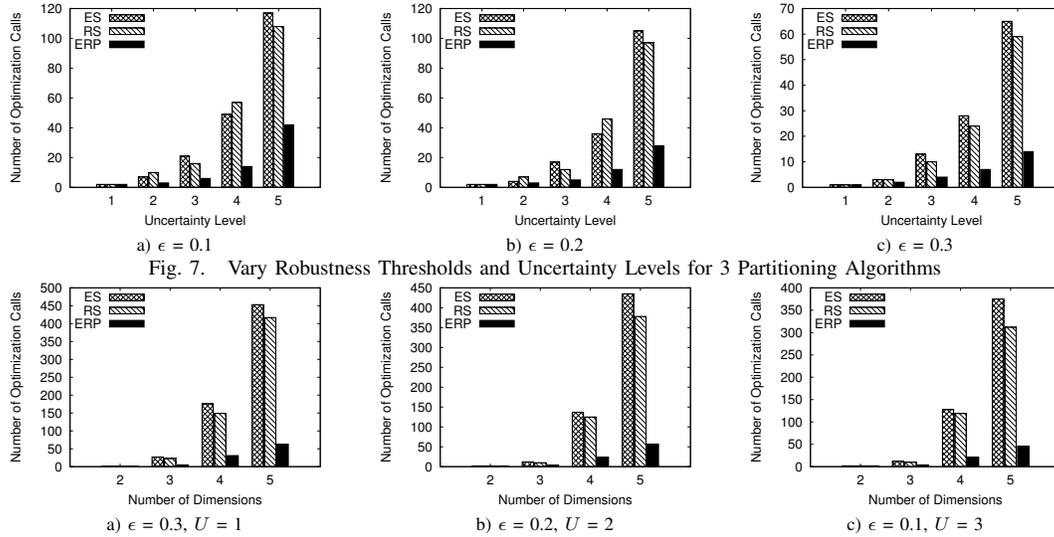


Fig. 7. Vary Robustness Thresholds and Uncertainty Levels for 3 Partitioning Algorithms

Fig. 8. Number of Optimizer Calls for 3 Partitioning Algorithms

ERP solution increases almost linearly by wisely choosing the partitioning areas. Furthermore, we benefit from our early termination mechanism that successfully avoids wasting computations on already robust areas.

D. Effectiveness of Physical Plan Generation

Next, we address the relative effectiveness of *GreedyPhy* versus *OptPrune* for **physical plan generation**. We measure the quality (i.e., space coverage and associated weight) of the respective algorithms. We vary the number of machines and also use different queries (equi-join of 10 to 20 streams).

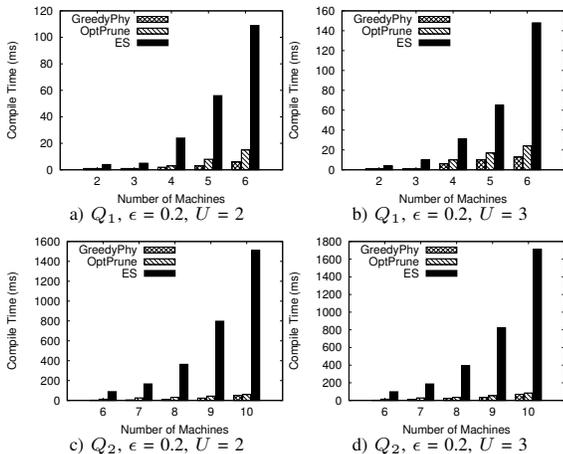


Fig. 9. Optimizer Performance for Finding Physical Plan

Figure 9 shows the average optimization time used by each algorithm for different numbers of operators. *GreedyPhy* is 12 times faster than its alternatives on average. Exhaustive search (*ES*) fares the worst because it treats the regions equally in the parameter space. This is a bad choice because *ES* wastes computations on the margin areas (i.e., less likely to actually occur at runtime) of the parameter space that are less likely to be supported by the resulting physical plan. *OptPrune* does fairly well compared to *ES*, because it tends to support the most important logical plans first. Moreover, it uses the result from *GreedyPhy* as bound for effective branch and bound search. In fact, it exhibits an optimization time similar to our *GreedyPhy* approach in most cases.

Figure 10 compares the physical plans produced by our *GreedyPhy* and *OptPrune* with the optimal solution (*ES*) for different queries (number of operators) given different system resources (number of machines). We define the average parameter coverage ratio rt_A as a metric to assess the relative effectiveness of the algorithm. The metric is defined as a ratio between the area ($Area(p.phy_A)$) covered by algorithm *A*'s physical plan and the area ($Area(p.phy_{ES})$) covered by the optimal physical plan. The physical plan generated by *OptPrune* is identical to *ES* even in the worst scenario, yet the search costs are much cheaper than those of *ES*. As for $rt_{GreedyPhy}$, the maximum rt is 0.94 and the minimum rt is 0.62 (ratio varies with different queries). Clearly, *GreedyPhy* sacrifices the quality of the robust physical plan for a reduction in compilation overhead.

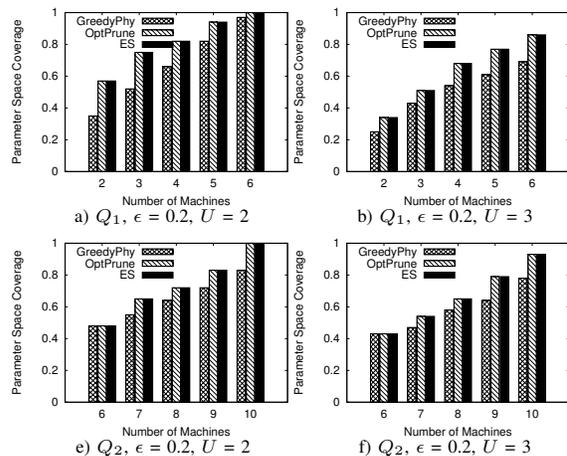


Fig. 10. Space Coverage of Physical Plan Generation

E. Measuring Runtime Performance

Average Tuple Processing Time: While our overarching goal is to achieve robust processing without load redistribution, DSPSs must process continuous queries in real-time. Thus, we now evaluate the average tuple processing performance of the RLD solution compared to the most prominent approaches, namely, resilient operator distribution [8] (ROD) and dynamic load distribution (DYN). Each query is run for 30 minutes

five times using these different solutions with the initial setup for input rates as shown in Table II. Using synthetic data, we vary stream rates by scaling the rate by a constant. The results in Figure 11(a) over such wide range of fluctuations not only show that RLD is robust in the input rate variations in most cases, but also point out where it fails. When the input rate is low (50%), ROD and DYN are almost as good as our RLD approach. When the input rate increases from 100% to 300%, neither ROD nor DYN are as robust as our RLD approach. This is likely due to ROD’s performance suffering from executing sub-optimal logical query plans once input data fluctuations arise. Without load migration, nodes may become a bottleneck. Consequently the tuple processing becomes delayed. DYN approach is also slower since moving operators may result in temporary poor performance due to the execution suspension of those operators.

However, DYN’s performance exceeds our RLD approach when the input rate fluctuation ratio is very large (400%). In this case, the load of the system cannot be well balanced by RLD since RLD only adopts one physical plan. DYN, on the other hand, performs best in such extreme fluctuations. The reason is that the computational resources are not sufficient for our RLD approach to handle such fluctuations with one single physical plan. Our RLD approach targets in particular to support fluctuations known a priori and does not optimize for operator movement at runtime.

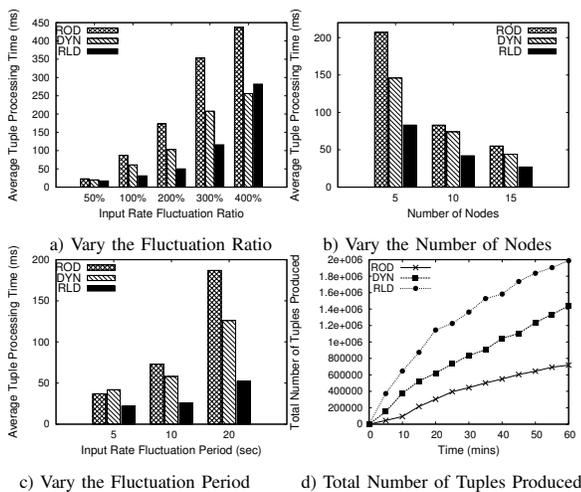


Fig. 11. RLD Runtime Performance

Our experiments further inspect whether the RLD approach is robust in the other two parameters, namely, the number of nodes and the input stream fluctuation periods. Specifically, the input stream fluctuation period alternates the input rate of each input stream periodically between a high rate and a low rate. The duration of the high rate interval equals the duration of the low rate interval. Both ROD and DYN do not exhibit the same robustness to the other two parameters as RLD does.

As indicated in Figure 11(b), when the number of machines is large, the performance difference among all three approaches is small. This is because when each machine has only a small amount of load, all 3 systems succeed to produce stable solutions that require few or no operator migrations. However, when the number of nodes is small, our RLD approach performs much better than the other two because

it pro-actively switches logical plans at runtime. Figure 11(c) shows that the average tuple processing time of RLD slightly increases while ROD and DYN suffer from the increasing fluctuation period. The reason is that ROD neither supports operator movement nor has multiple logical plans tuned to support different load fluctuations. Even though DYN supports physical plan migration, its performance is also slow due to the execution suspension caused by operator movements. Our RLD approach avoids load migration and smooths out the fluctuations by exploiting a combined solution of multiple robust logical plans and a corresponding robust physical plan.

Total Number of Tuples Produced: Figure 11(d) indicates the total number of tuples produced by the three load distributed stream models (i.e., ROD, DYN, and our RLD) over time, using the system configuration described in Table II. We ran queries for over an hour and show the average output for the first 60 minutes. Over time the total number of tuples output by RLD is significantly larger than those by either ROD or DYN solution. In ROD, the single logical plan causes a bottleneck in query processing performance, because it becomes non-optimal when the data streams fluctuate. Whereas DYN suffers from load migration overhead, whenever states of query operators get swapped among different machines [17].

Runtime Overhead: Now we compare the runtime overhead of RLD and *DYN* solutions. For both RLD and DYN, we consider any execution costs beyond the actual query processing to be runtime overhead. In RLD, tuples are grouped together into batches and assigned the appropriate logical plan based on the runtime statistics. Thus all tuples in a batch share the same plan. The only runtime overhead incurred in RLD is the initial classification to determine the execution plan for any arriving data, which was measured to be small, on average, about 2% of the query execution costs. On the contrary, DYN suffers from continuous load redistribution overhead. Multiple factors contribute to DYN’s runtime overhead, including the frequency of operator relocation, the state sizes of the moving operators, and the scale of operator relocation. The continuous re-optimization costs of DYN offset the performance gains achieved from using a better physical plan for short-term data fluctuations. In RLD, such overheads are avoided by exploiting a robust physical plan to support multiple logical plans.

VII. RELATED WORK

Robust [9], [18] and parametric query optimization [11], [19], [20], [21] are closely related areas. Robust query optimization [9], [18] aims to find one robust query plan that performs reasonably well for known uncertainties in statistics. However, if significant discrepancies exist between estimated and actual values, a single robust plan may fail to prevent performance degradation. Our work instead is unique in that it deploys multiple robust logical plans that together assure coverage across the entire parameter space, which a single plan would not be able to accomplish.

Similarly, parametric query optimization finds a set of plans that are optimal for different parameter settings. The early work [20], [21] optimized a parametric query for *all possible* values of uncertain variables, but postponed the final plan decisions to runtime once the actual statistics become known. Recent works [19], [11] proposes the concept of a plan diagram,

a pictorial enumeration of the query plans over the selectivity space. The authors propose to reduce the plan diagram for a query by merging plans whose costs are "close enough" with each other. Our problem faces different challenges compared with the above works. First, in our case the original plan diagram is not given. Instead, we compute robust logical plans based on prediction rather than observation on their cost behavior. Thus, the compile-time overhead is significantly reduced by our solution by not making traditional optimization calls repeatedly. Second, none of the above algorithms are directly applicable to our problem since assumptions made in these works do not consider the physical plan generation with resource limitations in distributed environments. Finally, unlike traditional parametric query optimization, our technique uses a probability model to capture the occurrence of points in the space at runtime, and strives to cover the most important regions in that space.

Our work is done in the context of distributed stream processing. The performance and scalability issues in centralized stream processing systems [17], [22], [23], [24] drew attention to distributed stream processing systems [2], [1]. Some works have proposed dynamic load distribution solutions for distributed stream processing [1], [25], [26]. Xing et al. [1] study how to minimize the load variances and maximize the correlations across all node pairs by *dynamically* distributing loads at runtime. SQPR [25] models query admission, allocation and reuse as a single constrained optimization formalization, and solves an approximate version. SODA [26] balances the load across all resources in the system by minimizing a weighted average of metrics that model resource utilization. However, their methodologies consider dynamic operator movements across machines as a reactive strategy when an imbalance arises. Our work instead aims to pro-actively produce a robust load distribution solution *without* runtime load redistribution and costly operator migration overhead. Given known data fluctuations as input, we pro-actively switch among pre-computed (*compile time*) robust logical plans at runtime, all of which are executed on the same physical operator allocation.

ROD [8], which we compared against in our experimental study (Section VI-E), produces a load distribution plan to keep the system feasible under workload fluctuations without load migration. However, our work has two key differences from ROD: 1. ROD only focuses on producing a feasible physical plan without exploiting multiple logical plans for a given query. Our work obtains a *many-to-one* mapping from a set of robust logical plans to a single physical plan that provides robust query processing performance. 2. The query processing performance is not guaranteed in ROD as it has no knowledge of the logical plans being executed on top of its feasible load distribution plan. On the contrary, our work uses proactive methodology to choose the best logical plan from our robust logical solution to be executed on a single physical solution. Thus, the query processing performance is further improved.

VIII. CONCLUSION

The ability to withstand stream data fluctuations is an important consideration in a distributed stream processing system. We design a scalable solution in which a DSPPS may benefit from different logical plans at runtime based on vary-

ing characteristics of the system. Our *ERP* then efficiently produces a robust logical solution that covers the parameter space. Taking the robust logical solution as input, *OptPrune* produces an optimal robust physical plan that supports the logical solution at runtime without operator relocation. Due to the effective bounding strategy, it succeeds to do so with minimal optimization time. Our experimental results on real world data show the promise of our RLD solution.

ACKNOWLEDGMENT

This work is supported by NSF grants: IIS-0917017 & 1018443 which we gratefully acknowledge.

REFERENCES

- [1] Y. Xing, S. Zdonik, and J.-H. Hwang, "Dynamic load distribution in the borealis stream processor," in *ICDE*, 2005, pp. 791–802.
- [2] M. Cherniack, H. Balakrishnan, M. Balazinska *et al.*, "Scalable distributed stream processing," in *CIDR*, 2003.
- [3] M. Balazinska, H. Balakrishnan, and M. Stonebraker, "Contract-based load management in federated distributed systems," in *Proceedings of the 1st NSDI*, 2004, pp. 15–15.
- [4] TradingMarkets, "<http://www.tradingmarkets.com/>."
- [5] B. A. Shirazi, K. M. Kavi, and A. R. Hurson, Eds., *Scheduling and Load Balancing in Parallel and Distributed Systems*. IEEE Computer Society Press, 1995.
- [6] R. Diekman and R. Preis, *Load balancing strategies for distributed memory machines*, 1999.
- [7] D. Kossmann, "The state of the art in distributed query processing," *ACM Comput. Surv.*, vol. 32, pp. 422–469, 2000.
- [8] Y. Xing, J.-H. Hwang, U. Çetintemel, and S. Zdonik, "Providing resiliency to load variations in distributed stream processing," in *VLDB*, 2006, pp. 775–786.
- [9] S. Babu, P. Bizarro, and D. DeWitt, "Proactive re-optimization," in *SIGMOD*, 2005, pp. 107–118.
- [10] N. Kabra and D. J. DeWitt, "Efficient mid-query re-optimization of sub-optimal query execution plans," in *SIGMOD*, 1998, pp. 106–117.
- [11] H. D., P. N. Darera, and J. R. Haritsa, "Identifying robust plans through plan diagram reduction," *Proc. VLDB Endow.*, pp. 1124–1140, 2008.
- [12] R. V. Nehme, E. A. Rundensteiner, and E. Bertino, "Self-tuning query mesh for adaptive multi-route query processing," in *EDBT*, 2009, pp. 803–814.
- [13] L. Peng, Y. Diao, and A. Liu, "Optimizing probabilistic query processing on continuous uncertain data," *PVLDB*, vol. 4, pp. 1169–1180, 2011.
- [14] C. Lei, E. Rundensteiner, and J. Guttman, "Robust distributed stream processing," Technical Report WPI-CS-TR-12-07, WPI, 2012.
- [15] Y. E. Ioannidis and S. Christodoulakis, "Optimal histograms for limiting worst-case error propagation in the size of join results," *ACM Trans. Database Syst.*, pp. 709–748, 1993.
- [16] T. M. Sutherland, B. Liu, M. Jbantova, and E. A. Rundensteiner, "D-cape: distributed and self-tuned continuous query processing," ser. CIKM '05, 2005, pp. 217–218.
- [17] Y. Zhu, E. A. Rundensteiner, and G. T. Heineman, "Dynamic plan migration for continuous queries over data streams," in *SIGMOD*, 2004, pp. 431–442.
- [18] V. Markl, V. Raman, D. Simmen *et al.*, "Robust query processing through progressive optimization," in *SIGMOD*, 2004, pp. 659–670.
- [19] N. Reddy and J. R. Haritsa, "Analyzing plan diagrams of database query optimizers," in *VLDB*, 2005, pp. 1228–1239.
- [20] Y. E. Ioannidis, R. T. Ng, K. Shim, and T. K. Sellis, "Parametric query optimization," in *VLDB*, 1992, p. 132151.
- [21] G. Graefe and K. Ward, "Dynamic query evaluation plans," in *SIGMOD*, 1989, pp. 358–366.
- [22] S. Chandrasekaran, O. Cooper, A. Deshpande *et al.*, "Telegraphcq: continuous dataflow processing," in *SIGMOD*, 2003, pp. 668–668.
- [23] D. J. Abadi, D. Carney *et al.*, "Aurora: a new model and architecture for data stream management," *The VLDB Journal*, pp. 120–139, 2003.
- [24] R. Motwani, J. Widom, A. Arasu *et al.*, "Query processing, approximation, and resource management in a data stream management system," in *CIDR*, 2003.
- [25] E. Kalyvianaki, W. Wiesemann, Q. H. Vu, D. Kuhn, and P. Pietzuch, "Sqpr: Stream query planning with reuse," in *ICDE*, 2011, pp. 840–851.
- [26] J. Wolf, N. Bansal, K. Hildrum *et al.*, "Soda: an optimizing scheduler for large-scale stream-based distributed computer systems," in *Middleware*, 2008, pp. 306–325.