

BOSTON UNIVERSITY  
METROPOLITAN COLLEGE

Dissertation

**QUERY PLANNING FOR STREAM PROCESSING**

by

**APOORVA TAMASKAR**

M.Sc., University of Glasgow, 2018  
B.Sc. Hons, Chennai Mathematical Institute , 2017

Submitted in partial fulfillment of the  
requirements for the degree of  
Master of Science

2020

© 2020 by  
APOORVA TAMASKAR  
All rights reserved

Approved by

First Reader

---

Kia Teymourian, PhD  
Assistant Professor of computer science

Second Reader

---

Eugene Pinsky, PhD  
Assistant Professor of Computer Science

Third Reader

---

Reza Rawassizadeh, PhD  
Associate Professor of Computer Science

## Acknowledgments

Over the course of writing this thesis, I have received lots of support and assistance. First of all I would like to thank my supervisor, Professor Kia Teymourian, who guided me and helped me formulate the research question, provided me with sources and methodology. Your insightful feedback pushed me to sharpen my thought process and helped me elevate the quality of my work.

I would like to thank the administrative staff at MET college, especially Ronette Lyle and the BU ISSO team who helped me through the administrative process.

I would like to acknowledge my friends and colleagues at Boston University for their wonderful feedback and discussions.

In addition I would like to thank my parents and my brother for their wise counsel, an empathetic ear and for always being there for me.

Apoorva Tamaskar

Metropolitan College

# QUERY PLANNING FOR STREAM PROCESSING

**APOORVA TAMASKAR**

Boston University, Metropolitan College , 2020

Major Professor: Kia Teymourian, PhD

Assistant Professor of computer science

## ABSTRACT

Data stream have become a common place scenario and they have become increasingly complex. With the network traffic volumes increasing doing meaningful calculation and inferring results from them is becoming increasingly challenging. With the recent increase in volumes of data stream lot of focus is upon cloud infrastructure and scalability. The flexibility the mathematical model of deep reinforcement learning and its wide variety of applications make it a promising approach. The continuous nature of data streams makes the conventional approach inapplicable due to reduced speed. The method applied and implemented is generic and applicable to wide range data streams. The goal of this study is to act as a proof of concept. The thesis explores the use of method on a query from Linear road benchmark data. The proposed method is demonstrated to be effective in optimizing the query execution. We hope the the work in this thesis encourages others to work on development of a robust and generalized system, capable of optimizing any query on any data stream.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Problem at hand . . . . .	2
1.3	Structure of thesis . . . . .	2
1.4	Conclusion . . . . .	3
<b>2</b>	<b>Related Work</b>	<b>4</b>
2.1	Introduction, Query optimization . . . . .	4
2.2	Converting SQL queries to parse trees . . . . .	4
2.3	Relational algebra . . . . .	5
2.3.1	Select operator $\sigma$ . . . . .	6
2.3.2	Projection operator $\pi$ . . . . .	7
2.3.3	Duplicate Elimination operator $\delta$ . . . . .	7
2.3.4	Aggregation operator $\gamma$ . . . . .	8
2.4	Converting Parse trees into logical expression . . . . .	8
2.5	Explain difficulties/ Time complexity . . . . .	10
2.5.1	Estimating size and cost . . . . .	11
2.5.2	Estimation of Projection . . . . .	12
2.5.3	Estimation of Selection . . . . .	12
2.5.4	Estimation of Join, single attribute . . . . .	12
2.5.5	Estimation of Join, multiple attribute . . . . .	13
2.5.6	Multiple Joins . . . . .	14

2.5.7	Union . . . . .	14
2.5.8	Intersection . . . . .	15
2.5.9	Difference . . . . .	15
2.5.10	Duplicate Elimination . . . . .	15
2.5.11	Grouping and Aggregation . . . . .	15
2.6	Other tools . . . . .	15
2.6.1	Histogram . . . . .	16
2.6.2	Heuristics . . . . .	16
2.7	Enumeration Methods . . . . .	17
2.7.1	Heuristic Selection . . . . .	18
2.7.2	Branch-and-Bound . . . . .	19
2.7.3	Hill Climbing . . . . .	19
2.7.4	Selinger-Style Optimization . . . . .	19
2.8	Join Order . . . . .	20
2.8.1	Join Trees . . . . .	21
2.8.2	DP to decide join order . . . . .	22
2.8.3	Greedy algorithm for join order . . . . .	22
2.9	Physical Query Plan . . . . .	23
2.9.1	Choosing a Selection Method . . . . .	24
2.9.2	Choosing a Join Method . . . . .	25
2.9.3	Pipelining Versus Materialization . . . . .	25
2.9.4	Pipelining Unary Operations . . . . .	26
2.9.5	Pipelining Binary Operations . . . . .	26
2.9.6	Notation for Physical Query Plans . . . . .	26
2.9.7	Ordering of Physical Operations . . . . .	28
2.10	Introduction to Data Streams . . . . .	29

2.11	Windowing and QoS . . . . .	30
2.12	Challenges of query optimization on data streams . . . . .	31
2.12.1	Resource Scheduling Strategies . . . . .	32
2.12.2	Load Shedding and Run-Time Optimization . . . . .	32
2.12.3	Complex Event and Rule Processing . . . . .	32
2.13	Deep Learning . . . . .	33
2.14	Reinforcement Learning . . . . .	34
2.14.1	Markov decision processes . . . . .	34
2.14.2	Finding Optima . . . . .	35
2.15	Deep Reinforcement Learning . . . . .	36
2.16	Conclusion . . . . .	36
<b>3</b>	<b>Stream Optimization</b>	<b>39</b>
3.1	Introduction . . . . .	39
3.2	Formalization Requirements . . . . .	40
3.3	Problem Statement . . . . .	40
3.4	Approach . . . . .	41
3.5	Example . . . . .	42
3.5.1	Input . . . . .	42
3.5.2	Query . . . . .	43
3.5.3	Segment Average Speed . . . . .	43
3.5.4	Toll Computation for Segments . . . . .	44
<b>4</b>	<b>Implementation</b>	<b>47</b>
4.1	Data Generation . . . . .	47
4.1.1	Schema . . . . .	48
4.2	Query Execution . . . . .	49
4.2.1	Simpler Query . . . . .	49



4.2.2	Complex Query . . . . .	50
4.3	Deep Reinforcement Learning . . . . .	54
4.3.1	MDP input . . . . .	54
4.3.2	States . . . . .	54
4.3.3	1-hot encoding . . . . .	54
4.3.4	Assumption . . . . .	55
4.4	Conclusion . . . . .	55
<b>5</b>	<b>Evaluation</b>	<b>57</b>
5.1	Optimality checking . . . . .	57
5.2	Results . . . . .	58
5.3	Discussion . . . . .	63
5.4	Conclusion . . . . .	63
<b>6</b>	<b>Conclusion and Further work</b>	<b>64</b>
6.1	Conclusion . . . . .	64
6.2	Further Work . . . . .	65
6.2.1	Data Generation . . . . .	65
6.2.2	Data Storage . . . . .	65
6.2.3	Features extraction . . . . .	65
6.2.4	Deep Reinforcement Learning . . . . .	65
6.2.5	Integration . . . . .	66
<b>A</b>	<b>Extra code</b>	<b>67</b>
	<b>References</b>	<b>75</b>

## Chapter 1

# Introduction

### 1.1 Motivation

In recent years, the ability to gather data has increased tremendously due to various sensory devices and the cheap cost of data storage. Some examples of data gathering and generating sources are,

- **Finance**, Stock market, the continuous flow includes the price of any given stock at any time.
- **Network Management**, Network administrators need to continuously monitor the traffic flow on their network as they need to help maintain a level of Quality of Service for their clients.
- **Healthcare**, Monitoring the vitals of the patients is an important and continuous process.
- **National Security** Monitoring airspace, getting continuous information from various sources and inferring results from them.

For these examples it is clear that there exists a need to continuously keep processing data stream and reporting back various statistics and any abnormality. But this continuous monitoring and reporting is difficult to do on traditional Data Base Management Systems(DBMS). A different approach has to be adapted to meet the demands in various industries.

Stream optimization in a certain sense is modification of the stream data to make the querying process faster. Few motivations for this include, to be able to make use of opportunities presented by faster calculation, mitigate risks before it is late, keeping the latest information readily available in understandable format.

## 1.2 Problem at hand

The increase in data gathering capabilities and speed need to be accompanied with increase in speed at which data can be analysed. There are various challenges that come up when addressing this, such as High frequency of the updates and reporting, buffer overflows, overheads, context-switches during processing and many more. The motivation for our approach is, most of the query optimization methods do not look at the data itself and rather try to come up with a very general optimization which can be applied to any query. Many times an in-depth look at the data might prove to be rather expensive. In this paper we try to detect trends in the data and help optimize querying using those and test our model on benchmark cases.

## 1.3 Structure of thesis

The next chapter gives an in-depth view of the pipeline used by the current state of art technology for query optimization in traditional data bases including the mathematical knowledge for simplification and the overall framework. The next chapter also introduces the reader to data stream and how data bases are used for them called DSMS and showcases an approach to optimize queries on data streams for the problem discussed above. The following chapter list out the details of implementation, challenges face, evaluation methods used, benchmark test case timings, followed by a summary of the paper.

## 1.4 Conclusion

With the help of this thesis, we hope to understand how query processing works in Data base management systems, the challenges they face, the way Data Stream Management Systems work, their challenges. We also propose a method for query optimization which will look at windows of data(defined later in the paper). The optimizations for the quering process will be tested against various test cases and compared with existing algorithms.

## Chapter 2

# Related Work

### 2.1 Introduction, Query optimization

A database can be thought of as a list of tables, where in each table itself can be considered as a list of data points ordered initially in the sequence they are entered.

There are various tools which can be used to connect to a database, here we focus on structured query languages(SQL). A simple SQL query looks like this

```
1  SELECT column_name_1 , column_name_2
2  FROM table_name
3  WHERE condition
```

This query is essentially asking to display the 2 columns from the table where the condition given is satisfied. This to particular query might be looking simple, but if the condition introduced is a complex one or if the table from which we need to return the output is complex, the question of how to execute the query optimally becomes difficult to answer.

### 2.2 Converting SQL queries to parse trees

The exact grammar for the conversion to parse tree is not listed below. But this step isn't simply conversion to a parse tree, the preprocessor which does the conversion also has several more functions.

For example, if the SQL query uses a "view" in the query as a relation, then each instance has to be replaced by the parse tree.

The preprocessor also has to conduct semantic checking, that is, check if relations used exist, check for ambiguity, and type checking. If a parse tree passes the preprocessing then it is said to be **valid**. In these parse trees, there are 2 types of nodes, one the atoms, which are essentially keywords in SQL, operators, constants and attributes. The second is Syntactic categories, these are names for families of sub-queries in triangular brackets. Each of the syntactic category has unique expansion into atoms and further syntactic categories.

## 2.3 Relational algebra

As we saw above, order of operations matters, if the order of operations is not thoughtout and done blindly alot of redundant steps are executed and memory is moved around unnecessarily. There are few ways to atleast look and analyse the operations and how they can be simplified.

Let  $R, S$  be relations. Some simple laws, associativity and commutativity can easily be verified:-

- $R \times S = S \times R$
- $(R \times S) \times T = R \times (S \times T)$
- $R \bowtie S = S \bowtie R$
- $(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$
- $R \cup S = S \cup R$
- $(R \cup S) \cup T = R \cup (S \cup T)$
- $R \cap S = S \cap R$
- $(R \cap S) \cap T = R \cap (S \cap T)$

When applying associative law on relations, need to be careful whether the conditions actually makes sense after the order is changed.

While the above identities work on both sets and bags(bags allow for repeatition). To show that laws for sets and bags do differ an easy way is to consider the distributive property.

$$A \cap_S (B \cup_S C) = (A \cap_S B) \cup_S (A \cap_S C)$$

$$A \cap_B (B \cup_B C) \neq (A \cap_B B) \cup_B (A \cap_B C)$$

We can simply show it with an example. Let  $A = \{t\}, B = \{t\}, C = \{t\}$ . The LHS comes to be  $\{t\}$ , whereas RHS is  $\{t, t\}$

### 2.3.1 Select operator $\sigma$

First we start with simple properties of the  $\sigma$  operator. Need to be careful about the attributes used in the select operator condition when pushing it down.

- $\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$
- $\sigma_{C_1 \vee C_2}(R) = (\sigma_{C_1}(R)) \cup_S (\sigma_{C_2}(R))$
- $\sigma_C(R \cup S) = \sigma_C(R) \cup \sigma_C(S)$
- $\sigma_C(R - S) = \sigma_C(R) - \sigma_C(S) = \sigma_C(R) - S$
- $\sigma_C(R \times S) = \sigma_C(R) \times S$
- $\sigma_C(R \bowtie S) = \sigma_C(R) \bowtie S$
- $\sigma_C(R \bowtie_D S) = \sigma_C(R) \bowtie_D S$
- $\sigma_C(R \cap S) = \sigma_C(R) \cap S$

### 2.3.2 Projection operator $\pi$

While for the Select operator( $\sigma$ ) the identities were quite straight forward with not many things to consider, the identities for Projection operator ( $\pi$ ) are bit more involved.

- $\pi_L(R \bowtie S) = \pi_L(\pi_M(R) \bowtie \pi_N(S))$ , where  $M, N$  are attributes required for the join or they are inputs to the projection.
- $\pi_L(R \bowtie_D S) = \pi_L(\pi_M(R) \bowtie_D \pi_N(S))$ , similar to above identity/ law.
- $\pi_L(R \times S) = \pi_L(\pi_M(R) \times \pi_N(S))$
- $\pi_L(R \cup_B S) = \pi_L(R) \cup_B \pi_L(S)$
- $\pi_L(\sigma_C(R)) = \pi_L(\sigma_C(\pi_M(R)))$

### 2.3.3 Duplicate Elimination operator $\delta$

The  $\delta$  operator eliminates duplicates from bags.

- $\delta(R) = R$ , if  $R$  does not have any duplicates.
- $\delta(R \times S) = \delta(R) \times \delta(S)$
- $\delta(R \bowtie S) = \delta(R) \bowtie \delta(S)$
- $\delta(R \bowtie_D S) = \delta(R) \bowtie_D \delta(S)$
- $\delta(\sigma_C(R)) = \sigma_C(\delta(R))$
- $\delta(R \cap_B S) = \delta(R) \cap_B S$



### 2.3.4 Aggregation operator $\gamma$

It is difficult to give identities for the aggregation operator, like done for the above operators. This is mostly due to how the details of how the aggregation operator is used.

- $\sigma(\gamma_L(R)) = \gamma_L(R)$
- $\gamma_L(R) = \gamma_L(\pi_M(R))$ , where  $M$  must at least contain the attributes used in  $L$ .

## 2.4 Converting Parse trees into logical expression

Till now, the only SQL related information presented is how to convert a Query into the parse tree, which is grammar dependent. Given the parse tree, need to substitute nodes by operators seen above, later this expression is optimized to be later converted to a physical query plan.

Now to convert the parse tree into the logical expression. First, look at the transformation of select-from-where statement.

- $\langle Query \rangle \rightarrow \langle SFW \rangle$ .
- $\langle SFW \rangle \rightarrow SELECT \langle SelList \rangle FROM \langle FromList \rangle WHERE \langle Condition \rangle$ .
- $\langle SelList \rangle \rightarrow \pi_L$ , where  $L$  is the list of attributes in  $\langle SelList \rangle$ .
- $\langle Condition \rangle \rightarrow \sigma_C$ , where  $C$  is the equivalent of  $\langle Condition \rangle$

While there is nothing inherently wrong about the last statement, need to consider the case where  $\langle Condition \rangle$  involves subqueries. A simple explanation about why it isn't allowed is a convention normally the subscript has to be a boolean condition, and if it was allowed otherwise, it would be a very expensive operation, as the subscript

in the select operator has to be evaluated at every element of the argument relation. This shows the redundancy of it. While if this is allowed, it can be simplified and made efficient, but has to be done on a case by case basis with the use of  $\bowtie$ ,  $\times$  functions.

Overall it is a good idea to not use subquerying and rather using joins.

At this point, by making the substitutions mentioned above and using the algebraic identities, we obtain a starting logical query plan. The query has to be transformed into a query which the compiler believes to be the cheapest or the optimal. But, a thing which further complicates the process is the join order.

With the current knowledge, few optimizing rules are evident.

- **Selection repositioning**, Selections should be pushed down as much as possible, but sometimes, might need to take the selection operator a level up first.
- Pushing projections down the parse tree, being careful with the new projections made in the process.
- Duplicate removal needs to be repositioned.
- $\sigma$  combined with  $\times$  below can result in equijoin, which is much more efficient.

Normally, parsers only have nodes with 0, 1, 2 children, which corresponds to unary and binary operators. But many of the operators (natural join, union, and intersection) are commutative and associative, so it helps combine them on a single level to provide the opportunity to prioritize which ones are done first. To do this step, the following guidelines are sufficient.

- We must replace the natural joins with theta-joins that equate the attributes of the same name.
- We must add a projection to eliminate duplicate copies of attributes involved in a natural join that has become a theta-join

- The theta-join conditions must be associative.

## 2.5 Explain difficulties/ Time complexity

Currently, we have taken a query as an input and converted it into a logical plan, and then applied more transformation using relational algebra to make the optimal query plan. The next step is to convert it into a physical query plan which can be executed. To complete this step, need to compare the cost of various physical query plan derived from the logical query plan. This plan with the least cost is then passed query execution engine and executed. When trying to select a query plan, need to select:-

- An order for the bracketing for associative and commutative operations like joins, unions, and intersections.
- The underlying algorithm for each operator in the logical plan, for instance, deciding whether a nested-loop join or a hash-join should be used.
- Need to consider additional operators like scanning, sorting, and so on that are needed for the physical plan but that were not present explicitly in the logical plan.
- A method for passing information/ variable values fromt the output of one operator to the next, for instance, by storing the intermediate result on disk or by using iterators and passing an argument one tuple or one main-memory buffer at a time.

Finally after selecting all these methods need to actually check the time required to excute the plan, one obvious way to calculate the time is actually executing the plan to see the cost. But if we try to find the optimal plan by this method then we are

essentially executing every plan. That is expensive and redundant. The next few section presents few methods which can help in this task.

### 2.5.1 Estimating size and cost

$B(R) :=$  is the number of blocks needed to hold all the tuples of relation  $R$ .

$T(R) :=$  is the number of tuples of relation  $R$ .

$V(R, a) :=$  is the value count for attribute  $a$  of relation  $R$ , that is, the number of distinct values relation  $R$  has in attribute  $a$ .

$V(R, [a_1, a_2, \dots, a_n]) :=$  is the number of distinct values  $R$  has when all of attributes  $a_1, a_2, \dots, a_n$  are considered together, that is, the number of tuples in  $\delta(\pi_{a_1, a_2, \dots, a_n}(R))$

Compared to algorithms learnt in an introductory algorithms class, here the memory usage is huge and need to remember that intermidate/ temporary relations calculated will also incur a cost on the system. So need to take them into account as well because the physical plan is selected to minimize the estimated cost of evaluating the query. Over all the estimation method should pass the following sanity check:-

- Give accurate estimates. No matter what method is used for executing query plans.
- Are easy to compute.
- Are logically consistent; that is, the size estimate for an intermediate relation should not depend on how that relation is computed. For instance, the size estimate for a join of several relations should not depend on the order in which we join the relations.

Even though there is no agreed upon method to do this, it does not matter because the goal is to help select a query plan and not to find the exact minimum cost. So even if the estimated cost is wrong, but if it is wrong similarly then we will still have the least costing plan.

### 2.5.2 Estimation of Projection

The projection operator is a bag operation, so it does not reduce the number of tuples, only reduces the size of each tuple.

$$T(R) = T(\pi(R))$$

$$B(R) \leq B(\pi(R))$$

### 2.5.3 Estimation of Selection

Let  $S = \sigma_{A=c}(R)$ , here  $A$  is an attribute of  $R$  and  $c$  is a constant

$T(S) = \frac{T(R)}{V(R,A)}$ , is it important to note that this is an estimate and not the actual value, this will be the actual value if all the attributes in  $A$  have equal occurrence. An even better estimate can be obtained if the DBMS stores a statistic known as histogram.

The above calculation was easy because of the equality. if  $S = \sigma_{A < c}(R)$ , we simply estimate  $T(S) = \frac{T(R)}{3}$

Now if  $S = \sigma_{A \neq c}(R)$ , can take  $T(S) = T(R)$  or  $T(S) = T(R) - \frac{T(R)}{V(R,A)}$

If  $S = \sigma_{A=c_1 \vee c_2}$ , take them to be independent conditions.

### 2.5.4 Estimation of Join, single attribute

For natural joins, we assume, the natural join of two relations involves only the equality of two attributes. That is, we study the join  $R(X,Y) \bowtie S(Y,Z)$ , but initially we assume that  $Y$  is a single attribute although  $X$  and  $Z$  can represent any set of attributes. But it is hard to find a good estimate as  $T(R \bowtie S) \in [0, T(R)*T(S)]$ , to help with this two assumptions are made.

- **Containment of Value Sets** If  $Y$  is an attribute appearing in several relations, then each relation chooses its values from the front of a fixed list of values  $y_1, y_2, y_3, \dots$  and has all the values in that prefix. As a consequence, if  $R$  and

$S$  are two relations with an attribute  $Y$ , and  $V(R, Y) \leq V(S, Y)$ , then every  $Y$ -value of  $R$  will be a  $Y$ -value of  $S$ .

- **Preservation of Value Sets** If we join a relation  $R$  with another relation, then an attribute  $A$  that is not a join attribute (i.e., not present in both relations) does not lose values from its set of possible values. More precisely, if  $A$  is an attribute of  $R$  but not of  $S$ , then  $V(R \bowtie S, A) = V(R, A)$ . Note that the order of joining  $R$  and  $S$  is not important, so we could just as well have said that  $V(S \bowtie R, A) = V(R, A)$ .

Using the above assumptions we can claim,

$$T(R \bowtie S) = \frac{T(R) * T(S)}{\max(V(R, Y), V(S, Y))}$$

Let  $V(R, Y) \leq V(S, Y)$ , then every tuple  $t$  of  $R$  has a chance of  $\frac{1}{V(S, Y)}$  of joining with a given tuple of  $S$ . Since there are  $T(S)$  tuples in  $S$ , the expected number of tuples that  $t$  joins with is  $\frac{T(S)}{V(S, Y)}$ . As there are  $T(R)$  tuples of  $R$ , the estimated size of  $R \bowtie S$  is  $\frac{T(R)*T(S)}{V(S, Y)}$ , if it was  $V(R, Y) \geq V(S, Y)$ , then  $\frac{T(R)*T(S)}{V(R, Y)}$

Guidelines for other type of joins:-

- The number of tuples in the result of an equijoin can be computed exactly as for a natural join, after accounting for the change in variable names.
- Other theta-joins can be estimated as if they were a selection following a product.

### 2.5.5 Estimation of Join, multiple attribute

Now we assume  $Y$  in  $R(X, Y) \bowtie S(Y, Z)$  represents multiple attributes. Say for example,  $R(X, y_1, y_2) \bowtie S(y_1, y_2, Z)$ . We again do a probability calculation.

Let  $r \in R, s \in S$  be a tuples, the probability that  $r, s$  agree on  $y_1$  is  $\frac{1}{\max(V(R, y_1), V(S, y_1))}$ , similarly for  $y_2$ . So we get the expected value to be

$$\frac{T(R) * T(S)}{\max(V(R, y_1), V(S, y_1)) * \max(V(R, y_2), V(S, y_2))}$$

From this the pattern is clear, need to divide  $T(R)*T(S)$  by maximum of  $V(R, y), V(S, y)$  for each attribute.

But this is only the calculation for  $T(R \bowtie S)$ , need to calculate for  $B(R \bowtie S)$  as well.

### 2.5.6 Multiple Joins

For this case we work with  $S = R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$

Here we have to make use of the containment assumption. Say the attribute  $A$  appears in  $k$  of  $R_i$ 's and the values corresponding to  $V(R_i, A)$  are  $v_1 \leq v_2 \leq \dots \leq v_k$ , need to find the probability that tuples agree on  $A$ .

Consider the tuple  $t_1$  chosen from the relation that has the smallest number of  $A$ -values,  $v_1$ . By the containment assumption, each of these  $v_1$  values is among the  $A$ -values found in the other relations that have attribute  $A$ . Consider the relation that has  $V_i$  values in attribute  $A$ . Its selected tuple  $t_i$  has probability  $\frac{1}{v_1}$  of agreeing with  $t_1$  on  $A$ . Since this claim is true for all  $i \in \{2, 3, \dots, k\}$ , the probability that all  $k$  tuples agree on  $A$  is the product  $\frac{1}{v_2 * v_3 * \dots * v_k}$ . This analysis gives us the rule for estimating the size of any join.

Start with the product of the number of tuples in each relation. Then for each attribute  $A$  appearing at least twice, divide by all but the least of the  $V(R, A)$ 's

### 2.5.7 Union

If  $U_B$  is used, it is the sum of the two individually.

If  $U_S$  is used, then number of tuples range from the max of two, to their sum. So take mean of the range.

### 2.5.8 Intersection

The number of tuples here ranges from 0 to the minimum of the two (in case of set intersection), so again can take the mean of this range.

### 2.5.9 Difference

The range for  $R - S$  is  $[T(R) - T(S), T(R)]$ , so again mean of the range.

### 2.5.10 Duplicate Elimination

The range for  $\delta(R)$  is  $[1, T(R)]$ , so again can take the mean, there can be other estimates as well.

A nice compromise is  $\min(\frac{T(R)}{2}, \prod V(R, a_i))$

### 2.5.11 Grouping and Aggregation

Same as duplicate.

## 2.6 Other tools

We assume that the "cost" of evaluating an expression is approximated well by the number of disk I/O's performed. The number of disk I/O's, in turn, is influenced by:

- The particular logical operators chosen to implement the query.
- The sizes of intermediate results, need to pass them to the next function.
- The physical operators used to implement logical operators, e.g., the choice of a one-pass or two-pass join, or the choice to sort or not sort a given relation.
- The ordering of similar operations, especially joins.



- The method of passing arguments from one physical operator to the next.

### 2.6.1 Histogram

In the earlier section the statistics were heavily used in calculations. Another statistic that can be stored is the histogram.

If  $V(R, A)$  is not too large, then the histogram may consist of the number (or fraction) of the tuples having each of the values of attribute  $A$ . If there are a great many values of this attribute, then only the most frequent values may be recorded individually, while other values are counted in groups.

**Equal width** To start, select 2 parameters,  $w$  the width and  $v_0$  a beginning point of a column in the histogram, which will initially be the considered the lower bound and if an even lower value is noticed, make a smaller column as well and update the lower bound.

**Equal height** These are the common "percentiles". A percentile  $p, 2p$ , so on.

**Most frequent values** List the most common values and their numbers of occurrences. This information may be provided along with a count of occurrences for all the other values as a group, or we may record frequent values in addition to an equal width or equal height histogram for the other values.

### 2.6.2 Heuristics

One important use of cost estimates for queries or subqueries is in the application of heuristic transformations of the query.

Heuristics applied independent of cost estimates can be expected almost certainly to improve the cost of a logical query plan

However, there are other points in the query optimization process where estimating the cost both before and after a transformation will allow us to apply a transformation where it appears to reduce cost and avoid the transformation otherwise. In particular,

when the preferred logical query plan is being generated, we may consider a number of optional transformations and the costs before and after. Because we are estimating the cost of a logical query plan, so we have not yet made decisions about the physical operators that will be used to implement the operators of relational algebra, our cost estimate cannot be based on disk I / O's. Rather, we estimate the sizes of all intermediate results their sum is our heuristic estimate for the cost of the entire logical plan.

## 2.7 Enumeration Methods

The naive method of find the least costing plan is enumerating all the possible plans and calculating the cost for them and then selecting the minimum one.

**Top-Down** Here, we work down the tree of the logical query plan from the root. For each possible implementation of the operation at the root, we consider each possible way to evaluate its argument(s) , and compute the cost of each combination, taking the best.

**Bottom-up** For each subexpression of the logical-query-plan tree, we compute the costs of all possible ways to compute that subexpression. The possibilities and costs for a subexpression E are computed by considering the options for the subexpressions for E , and combining them in all possible ways with implementations for the root operator of E.

There is actually not much difference between the two approaches in their broadest interpretations, since either way, all possible combinations of ways to implement each operator in the query tree are considered. When limiting the search, a top-down approach may allow us to eliminate certain options that could not be eliminated bottom-up. However, bottom-up strategies that limit choices effectively have also been developed. there is an apparent simplification of the bottom-up method, where

we consider only the best plan for each subexpression when we compute the plans for a larger subexpression. This approach, called **dynamic programming**, is not guaranteed to yield the best plan, although often it does. The approach called Selinger-style (or System-R-style) optimization exploits additional properties that some of the plans for a subexpression may have, in order to produce optimal overall plans from plans that are not optimal for certain subexpressions.

### 2.7.1 Heuristic Selection

One option is to use the same approach to selecting a physical plan that is generally used for selecting a logical plan: make a sequence of choices based on heuristics. Few common ones are:-

- If the logical plan calls for a selection  $\sigma_{A=c}(R)$ , and stored relation  $R$  has an index on attribute  $A$ , then perform an index-scan to obtain only the tuples of  $R$  with  $A$ -value equal to  $c$ .
- if the selection involves one condition like  $A = c$  above, and other conditions as well, we can implement the selection by an index scan followed by a further selection on the tuples, which we shall represent by the physical operator filter.
- If an argument of a join has an index on the join attribute(s), then use an index-join with that relation in the inner loop.
- If one argument of a join is sorted on the join attribute(s), then prefer a sort-join to a hash-join, although not necessarily to an index-join if one is possible.
- When computing the union or intersection of three or more relations, group the smallest relations first.

### 2.7.2 Branch-and-Bound

This approach, often used in practice, begins by using heuristics to find a good physical plan for the entire logical query plan. Let the cost of this plan be  $C$ . Then as we consider other plans for subqueries, we can eliminate any plan for a sub query that has a cost greater than  $C$ , since that plan for the sub query could not possibly participate in a plan for the complete query that is better than what we already know. Likewise, if we construct a plan for the complete query that has cost less than  $C$ , we replace  $C$  by the cost of this better plan in subsequent exploration of the space of physical query plans. An important advantage of this approach is that we can choose when to cut off the search and take the best plan found so far. For instance, if the cost  $C$  is small, then even if there are much better plans to be found, the time spent finding them may exceed  $C$ , so it does not make sense to continue the search. However, if  $C$  is large, then investing time in the hope of finding a faster plan is wise.

### 2.7.3 Hill Climbing

This approach, in which we really search for a "valley" in the space of physical plans and their costs, starts with a heuristically selected physical plan. We can then make small changes to the plan, e.g., replacing one method for an operator by another, or reordering joins by using the associative and/or commutative laws, to find "nearby" plans that have lower cost. When we find a plan such that no small modification yields a plan of lower cost, we make that plan our chosen physical query plan.

### 2.7.4 Selinger-Style Optimization

This approach improves upon the dynamic-programming approach by keeping for each subexpression not only the plan of least cost, but certain other plans that have higher cost, yet produce a result that is sorted in an order that may be useful higher up in the expression tree. Examples of such interesting orders are when the result of

the subexpression is sorted on one of:

- The attribute(s) specified in a sort operator ( $\tau$ ) at the root.
- The grouping attribute(s) of a later group-by operator ( $\gamma$ ).
- The join attribute(s) of a later join.

If we take the cost of a plan to be the sum of the sizes of the intermediate relations, then there appears to be no advantage to having an argument sorted. However, if we use the more accurate measure, disk I/O's, as the cost, then the advantage of having an argument sorted becomes clear if we can use one of the sort-based algorithms and save the work of the first pass for the argument that is sorted already.

## 2.8 Join Order

Join takes in two arguments, while the end result is independent of the order of the two arguments, the method used to compute the result may be dependent on the order. Perhaps most important, the one-pass join reads one relation preferably the smaller into main memory, creating a structure such as a hash table to facilitate matching of tuples from the other relation. It then reads the other relation, one block at a time, to join its tuples with the tuples stored in memory.

For instance, suppose that when we select a physical plan we decide to use a one-pass join. Then we shall assume the left argument of the join is the smaller relation and store it in a main-memory data structure. This relation is called the build relation. The right argument of the join, called the probe relation, is read a block at a time and its tuples are matched in main memory with those of the build relation. Other join algorithms that distinguish between their arguments include:

- Nested-loop join, where we assume the left argument is the relation of the outer loop.

- Index-join, where we assume the right argument has the index.

### 2.8.1 Join Trees

When we have the join of two relations, we need to order the arguments. We shall conventionally select the one whose estimated size is the smaller as the left argument. Notice that the algorithms mentioned above – one-pass, nested loop, and indexed – each work best if the left argument is the smaller. More precisely, one-pass and nested-loop joins each assign a special role to the smaller relation (build relation, or outer loop), and index-joins typically are reasonable choices only if one relation is small and the other has an index. It is quite common for there to be a significant and discernible difference in the sizes of arguments, because a query involving joins very often also involves a selection on at least one attribute, and that selection reduces the estimated size of one of the relations greatly.

When we need to join more than 2 relations, the order in which they are joined can be represented by a binary tree, where each node has either 0 or 2 children. A tree where the right child always a leaf node is called a left deep tree, one can similarly define a right deep tree. Any other tree is will be called **bushy**. We will stick with left deep tree due to their interaction with various common join algorithms. This introduced limitation also helps to reduce search space.

If one-pass joins are used, and the build relation is on the left, then the amount of memory needed at any one time tends to be smaller than if we used a right-deep tree or a bushy tree for the same relations.

If we use nested-loop joins, with the relation of the outer loop on the left, then we avoid constructing any intermediate relation more than once.

### 2.8.2 DP to decide join order

Suppose we need to calculate  $R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$ . For the DP algorithm construct a table with an entry for each subset of one or more of the  $n$  relations. In that table we put

- the estimated size of the join of these relations.
- the least cost of computing the join of these relations. Other, more complex estimates, such as total disk I/O's, could be used if we were willing and able to do the extra calculation involved.
- The expression that yields the least cost. This expression joins the set of relations in question, with some grouping. We can optionally restrict ourselves to left-deep expressions, in which case the expression is just an ordering of the relations.

### 2.8.3 Greedy algorithm for join order

Even the carefully limited search of dynamic programming leads to a number of calculations that is exponential in the number of relations joined. It is reasonable to use an exhaustive method like dynamic programming or branch-and-bound search to find optimal join orders of five or six relations. However, when the number of joins grows beyond that, or if we choose not to invest the time necessary for an exhaustive search, then we can use a join-order heuristic in our query optimizer.

The most common choice of heuristic is a greedy algorithm, where we make one decision at a time about the order of joins and never backtrack or reconsider decisions once made. We shall consider a greedy algorithm that only selects a left-deep tree. The "greediness" is based on the idea that we want to keep the intermediate relations as small as possible at each level of the tree.

Start with the pair of relations whose estimated join size is smallest. The join of these relations becomes the current tree. Find, among all those relations not yet included in the current tree, the relation that, when joined with the current tree, yields the relation of smallest estimated size. The new current tree has the old current tree as its left argument and the selected relation as its right argument.

## 2.9 Physical Query Plan

We have parsed the query, converted it to an initial logical query plan, and improved that logical query plan with transformations. Part of the process of selecting the physical query plan is enumeration and cost estimation for all of our options, then focused on the question of enumeration, cost estimation, and ordering for joins of several relations. By extension, we can use similar techniques to order groups of unions, intersections, or any associative/commutative operation

There are still several steps needed to turn the logical plan into a complete physical query plan.

- Selection of algorithms to implement the operations of the query plan, when algorithm-selection was not done as part of some earlier step such as selection of a join order by dynamic programming.
- Decisions regarding when intermediate results will be materialized ( created whole and stored on disk) , and when they will be pipelined (created only in main memory, and not necessarily kept in their entirety at any one time) .
- Notation for physical-query-plan operators, which must include details regarding access methods for stored relations and algorithms for implementation of relational-algebra operators.



### 2.9.1 Choosing a Selection Method

One of the important steps in choosing a physical query plan is to pick algorithms for each selection operator.

Assuming there are no multidimensional indexes on several of the attributes, then each physical plan uses some number of attributes that each have an index, and are compared to a constant in one of the terms of the selection. We then use these indexes to identify the sets of tuples that satisfy each of the conditions. For simplicity, we shall not consider the use of several indexes in this way. Rather, we limit our discussion to physical plans that:

- Use one comparison of the form  $A\theta c$ , where  $A$  is an attribute with an index,  $c$  is a constant, and  $\theta$  is a comparison operator such as  $=$  or  $<$ .
- Retrieve all tuples that satisfy the comparison, using the index scan physical operator.
- Consider each tuple selected to decide whether it satisfies the rest of the selection condition. We shall call the physical operator that performs this step Filter; it takes the condition used to select tuples as a parameter, much as the `rr` operator of relational algebra does.

In addition to physical plans of this form, we must also consider the plan that uses no index but reads the entire relation (using the table-scan physical operator) and passes each tuple to the Filter operator to check for satisfaction of the selection condition.

We decide among the physical plans with which to implement a given selection by estimating the cost of reading data for each possible option. To compare costs of alternative plans we cannot continue using the simplified cost estimate of intermediate-relation size. The reason is that we are now considering implementations of a single

step of the logical query plan, and intermediate relations are independent of implementation. Thus, we shall refocus our attention and resume counting disk I/O's.

### 2.9.2 Choosing a Join Method

On the assumption that we know (or can estimate) how many buffers are available to perform the join, we can apply the formulas for sort/ indexed/ hash join. However, if we are not sure of, or cannot know, the number of buffers that will be available during the execution of this query (because we do not know what else the DBMS is doing at the same time), or if we do not have estimates of important size parameters such as the  $V(R, a)$ 's, then there are still some principles we can apply to choosing a join method. Similar ideas apply to other binary operations such as unions, and to the full-relation, unary operators,  $\gamma, \delta$

### 2.9.3 Pipelining Versus Materialization

The last major issue we shall discuss in connection with choice of a physical query plan is pipelining of results. The naive way to execute a query plan is to order the operations appropriately (so an operation is not performed until the argument(s) below it have been performed), and store the result of each operation on disk until it is needed by another operation. This strategy is called materialization, since each intermediate relation is materialized on disk. A more subtle, and generally more efficient, way to execute a query plan is to interleave the execution of several operations. The tuples produced by one operation are passed directly to the operation that uses it, without ever storing the intermediate tuples on disk. This approach is called pipelining, and it typically is implemented by a network of iterators whose functions call each other at appropriate times. Since it saves disk I/O 's, there is an obvious advantage to pipelining, but there is a corresponding disadvantage. Since several operations must share main memory at any time, there is a chance that algorithms with higher disk-

I/O requirements must be chosen , or thrashing will occur, thus giving back all the disk-I/O savings that were gained by pipelining, and possibly more.

#### 2.9.4 Pipelining Unary Operations

Unary operations - selection and projection - are excellent candidates for pipelining. Since these operations are tuple-at-a-time, we never need to have more than one block for input, and one block for the output. We may implement a pipelined unary operation by iterators. The consumer of the pipelined result calls **GetNext()** each time another tuple is needed. In the case of a projection, it is only necessary to call **GetNext()** once on the source of tuples, project that tuple appropriately, and return the result to the consumer. For a selection  $\sigma_C$  (technically, the physical operator **Filter(C)** ) , it may be necessary to call **GetNext()** several times at the source, until one tuple that satisfies condition  $c$  is found.

#### 2.9.5 Pipelining Binary Operations

The results of binary operations can also be pipelined. We use one buffer to pass the result to its consumer, one block at a time. However, the number of other buffers needed to compute the result and to consume the result varies, depending on the size of the result and the sizes of other relations involved in the query. We shall use an extended example to illustrate the tradeoffs and opportunities.

#### 2.9.6 Notation for Physical Query Plans

##### Operators for Leaves

Each relation  $R$  that is a leaf operand of the logical-query-plan tree will be replaced by a scan operator. The options are:

- **TableScan(R)**: All blocks holding tuples of  $R$  are read in arbitrary order.

- **SortScan(R,L):** Tuples of  $R$  are read in order, sorted according to the attribute(s) on list  $L$ .
- **IndexScan(R,C):** Here,  $C$  is a condition of the form  $A\theta c$ , where  $A$  is an attribute of  $R$ ,  $\theta$  is a comparison such as  $=$  or  $<$ , and  $c$  is a constant. Tuples of  $R$  are accessed through an index on attribute  $A$ . If the comparison  $\theta$  is not  $=$ , then the index must be one, such as a B-tree, that supports range queries.
- **IndexScan(R,A):** Here  $A$  is an attribute of  $R$ . The entire relation  $R$  is retrieved via an index on  $R.A$ . This operator behaves like **TableScan**, but may be more efficient in certain circumstances, if  $R$  is not clustered and/or its blocks are not easily found.

### Physical Operators for Selection

A logical operator  $\sigma_C(R)$  is often combined, or partially combined, with the access method for relation  $R$ , when  $R$  is a stored relation. Other selections, where the argument is not a stored relation or an appropriate index is not available, will be replaced by the corresponding physical operator we have called **Filter**.

### Physical Sort Operators

Sorting of a relation can occur at any point in the physical query plan. We have already introduced the **SortScan(R, L)** operator, which reads a stored relation  $R$  and produces it sorted according to the list of attributes  $L$ . When we apply a sort-based algorithm for operations such as join or grouping, there is an initial phase in which we sort the argument according to some list of attributes. It is common to use an explicit physical operator  $Sort(L)$  to perform this sort on an operand relation that is not stored. This operator can also be used at the top of the physical-query-plan tree if the result needs to be sorted because of an **ORDER BY** clause in the original query,

### Other Relational-Algebra Operations

All other operations are replaced by a suitable physical operator. These operators can be given designations that indicate:

- The operation being performed, e.g., join or grouping.
- Necessary parameters, e.g., the condition in a theta-join or the list of elements in a grouping.
- A general strategy for the algorithm: sort-based, hash-based, or in some joins, index-based.
- A decision about the number of passes to be used: one-pass, two-pass, or multipass (recursive, using as many passes as necessary for the data at hand) . Alternatively, this choice may be left until run-time.
- An anticipated number of buffers the operation will require.

### 2.9.7 Ordering of Physical Operations

Our final topic regarding physical query plans is the matter of order of operations. The physical query plan is generally represented as a tree, and trees imply something about order of operations, since data must flow up the tree. However, since bushy trees may have interior nodes that are neither ancestors nor descendants of one another, the order of evaluation of interior nodes may not always be clear. Moreover, since iterators can be used to implement operations in a pipelined manner, it is possible that the times of execution for various nodes overlap, and the notion of "ordering" nodes makes no sense.

If materialization is implemented in the obvious store-and-later-retrieve way, and pipelining is implemented by iterators, then we may establish a fixed sequence of events whereby each operation of a physical query plan is executed. The following rules summarize the ordering of events implicit in a physical query-plan tree:

- Break the tree into subtrees at each edge that represents materialization. The subtrees will be executed one-at-a-time.
- Order the execution of the subtrees in a bottom-up, left-to-right manner. To be precise, perform a preorder traversal of the entire tree. Order the subtrees in the order in which the preorder traversal exits from the subtrees.
- Execute all nodes of each subtree using a network of iterators. Thus, all the nodes in one subtree are executed simultaneously, with GetNext calls among their operators determining the exact order of events.

## 2.10 Introduction to Data Streams

A **Data Stream** can be thought of as a continuously generated never ending sequence of data. For example, network logistics, healthcare monitoring information, security cameras, all of these continuously produce data, these data sequences/ streams are crucial to their respective users/ organisations, realtime detection of anomaly, pattern or changes thereof, or any other type of analytics is in need currently.

But there are many complications while trying to do this, the frequency of data generation isn't stable as it can keep changing throughout the cycle. Most of these service have a Quality of Service (or QoS) requirements defined which takes into account various metrics such as response time, memory usage, and throughput and more. Sometimes these requirements are so heavy and indulgent that it is infeasible to load the incoming data streams into a persistent store and process them effectively using DBMS tools. Hence special tools are required for these, similar to how DBMS has a variety of steps, Data Stream Management Systems (DSMSs) have multiple steps as well processing of sensor data, pervasive computing, situation monitoring, real-time response, approximate algorithms, on-the-fly mining, complex event processing, and more. While it may seem that optimizing Data stream applications to

their specific domains may make them difficult to study, it is possible to generalize the functions/ steps enough to have a basic architecture to study.

Queries used in traditional DBMS are called **ad hoc queries**, where as the queries used in DSMS are called **ContinuousQueries (CQ)**. Ad hoc queries are generally specified, optimized, and evaluated over a snapshot of a database. Whereas, CQs are specified once and evaluated repeatedly against new data over a specified life span or as long as there exists data in the stream. They are long-running queries that produce output continuously. The result is also assumed to be a stream possibly with differing rates and schema (as compared to the input). The difference between ad hoc queries and CQs can be best understood based on their relationship to the data over which they are processed. Different or changing ad hoc queries are processed over (relatively) static data in contrast to the same (or static) CQs that are processed repeatedly over frequently changing (or dynamic) data. It is clear that DBMS can't really cope with data streams in terms of high frequency updates, fulfilling QoS and continuous output. Hence we make use of Data Stream Management Systems (DSMSs).

## 2.11 Windowing and QoS

Most of the languages used for this are based on SQL, some of them are, Continuous Query Language (CQL), StreamSQL and ESP. Typically, continuous queries consist of relational operators such as select, project, join, and other aggregation operators. A logical query plan, analogous to a query tree used in a traditional DBMS, can be generated from the specification of a CQ. It can then be transformed into a query plan consisting of detailed algorithms used for computing each operator. Compared to DBMS where either pipelines or materialization is used, DSMS use a push or dataflow methodology. All operators in a DSMS compute on data items as they

arrive and cannot assume the stream to be finite. This has significant implications for the computation of a number of operators such as join, sort, and some aggregation operators. These operations cannot be completed without processing the entire input data set (or sets) which poses problems due to the unbounded nature of streams. As a result, these operators will block and produce no output until the stream ends. Hence, they are termed blocking operators. For an operator to output results continuously (and hopefully smoothly) and not wait until the end of the stream, it is imperative that these blocking operators be converted into non-blocking ones. The notion of a window has been introduced to overcome the blocking aspect of a number of operators. Informally, a window defines a finite portion of the stream (as a relation) for processing purposes. A window specification, added to a continuous query specification, can produce time-varying, finite relations out of a stream.

QoS management is important and critical to the success of a DSMS. Few of the popular QoS metrics are Tuple latency, Memory usage, Throughput, Smooth or bursty nature of output streams, Accuracy of results in terms of error tolerance. Most of the metrics can be specified in the query written. Overall to correlate CQs and QoS a simple question can be asked

Given a set of CQs with their QoS specifications, what resources are needed to compute the given CQs and satisfy their QoS requirements?

## 2.12 Challenges of query optimization on data streams

Picking up on the question from the last section, it is important to be able to address the inverse as well, i.e. given resources, CQs and QoS, need to verify that QoS are satisfied. Like any verification, there are two general ways of tackling this problem.

- Try to model the load generated by the input stream and get a conservative estimate on the metrics used to measure the QoS.



- Continuously monitor the metrics while the program is running, on violation identify the bottleneck and improve upon it.

As seen in previous chapters, creating a model to predict the load can be very difficult and sometimes more expensive than the best result. So we focus on the second approach.

### **2.12.1 Resource Scheduling Strategies**

Similar to the case of traditional DBMS, querying in DSMS involves multiple steps as well. The resource requirements for each step varies hence the need for a scheduling method which can determine an optimal schedule. To further complicate this, the various types of CQs have different QoS, hence the wiggle room for the performance is high, and as DSMS process queries in real time more considerations have to taken into account compared to traditional data bases, e.g. different operators require different amount of memory and time to process data.

### **2.12.2 Load Shedding and Run-Time Optimization**

It is important to note that optimizing resource allocation using any state of the art strategy does not guarantee that resources will be sufficient, as there might be some points in the data stream which overloads particular query, hence wiggle room in the performance/ accuracy, i.e. some of the data can be discarded to decrease the quality of the result and give an approximate result. This process of decreasing quality of result by not considering some tuple is called load shedding. Given that discarding tuples of the data stream is allowed, how do we decide which ones to discard?

### **2.12.3 Complex Event and Rule Processing**

In the first section we said we proceed with the second strategy, i.e. monitor the metrics. Can extend that idea further to monitor the output of CQs to detect anomalies.

Initially, Comple Event Processing was not a part of DSMS, but now many models have been developed to integrate the two. With the integration comes many additional functionalities and helps develop a theoretical base, but also puts additional burden on the system. This along with run time optimiations makes the implementation of DSMS quite complicated.

### 2.13 Deep Learning

We expand on linear regression model. Given  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$  as data points, the goal is to find a function  $h_\theta(x)$  which is "as close as possible" to  $y_i$ .

One way to do that is, define a loss function  $J$  which measures the distance between prediction and actual value of  $y^{(i)}$  as

$$J^{(i)}(\theta) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J^{(i)}(\theta)$$

The goal then reduces to finding parametrize  $h_\theta$  and finding the  $\theta$  so that  $J(\theta)$  is minimized.

Similarly in neural networks, we define a parameterization of  $h$  and then try to minimize the loss function.

Consider the simplest neural network, which consists of an input layer and a output layer with a single neuron. Say, the input is  $x$ , then the output neuron has input  $h_\theta(x)$ , and the output from the neuron is  $A(h_\theta(x))$ , where  $A$  is an activation function.

In the above example, if  $x \in \mathbb{R}^n$  then,  $h : \mathbb{R}^n \rightarrow \mathbb{R}$ , which required us to find a single  $\theta$ . Now say,  $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , then we need to find  $\theta_1, \theta_2, \dots, \theta_m$ . Which results in a matrix of dimension  $n * m$ . For each of the  $m$  neurons in the output

layer, the input will then be  $h_{\theta_i}(x)$  and the output of the  $m$  neurons will be a vector  $A(h_{\theta_1}(x)), A(h_{\theta_2}(x)), \dots, A(h_{\theta_m}(x)) \in \mathbb{R}^m$ . This is an example of **Fully connected 2-layer neural network**.

**Multi-layer fully connected neural network**, for this keep on adding one layer at a time to the neural network, and use the output from the last layer as the input.

At the end the problem is still reduced to finding parameters which minimize the loss and hence can be solved using gradient descent and similar methods, but it is much more computationally intense.

## 2.14 Reinforcement Learning

### 2.14.1 Markov decision processes

(Xu et al., 2018) Markov decision process (MDP) is used to formalize various types of stochastic processes. In MDPs, the goal of the agent is to make a sequence of actions to optimize/ maximize an objective function.

Formally a MDP is a 4-tuple

$$\langle S, A, P(s, a), R(s, a) \rangle$$

$S \rightarrow$  Set of all possible states the agent can be in.

$A \rightarrow$  Set of all possible actions the agent can take.

$P(s, a) \rightarrow$  A probability distribution of going to various states given current state and action.  $s^1 \sim P(s, a)$

$R(s, a) \rightarrow$  Reward for taking action  $a$  on state  $s$ .

To the above 4 tuple, additional parameters such as starting state and discount factor can be added.

A run of MDP looks like following:-

Start from the state  $s_0$ , choose an action  $a_0 \in A$  to take, the action will result in transition in state from  $s_0 \rightarrow s_1$ , where  $s_1$  is taken from the probability distribution  $P(s_0, a_0)$ . This change in state results in a reward  $R(s_0, a_0)$ . Then again an action  $a_1$  is taken, which then results in a transition from state  $s_1 \rightarrow s_2$ , where  $s_2$  is taken from the probability distribution  $P(s_1, a_1)$ . This change in state results in a reward  $\gamma R(s_1, a_1)$ . and so on.

This will result in a total reward of

$$\sum_{i=0}^n \gamma^i R(s_i, a_i)$$

The goal with reinforcement learning is to maximize  $\mathbb{E}[\sum_{i=0}^n \gamma^i R(s_i, a_i)]$ .

**Policy** is a function  $\pi : S \rightarrow A$ . **Executing a policy** means at state  $s$ , action  $a = \pi(s)$  is executed. This leads to the **value function** for policy  $\pi$ .

$$V^\pi(s) = \mathbb{E}\left[\sum_{i=0}^n \gamma^i R(s_i, a_i) \mid s_0 = s, \pi\right]$$

, This value function satisfies the Bellman Equation.

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P(s, \pi(s))(s') V^\pi(s')$$

Now we need to find the optimal policy, which can be formally written as

$$V^*(S) = \max_{\pi} V^\pi(S)$$

### 2.14.2 Finding Optima

```

1 for s in S:
2     V(s)=0
3 while(not converged):
4     for s in S:
5         V(s)=R(s)+max over all action [gamma*(sum(P(s,a,s')V(s')))]

```

The above algorithm is called the value iteration algorithm(Ng, 2020).

```

1 initialize random pi
2 while(not converged):
3     V=V(pi)
4     for s in S:
5         pi(s)=max over all actions[sum(P(s,a,s')V(S'))]
```

The second algorithm is called Policy iteration algorithm(Ng, 2020).

## 2.15 Deep Reinforcement Learning

Combining the two ideas as shown(Fan et al., 2020)

## 2.16 Conclusion

To summarize this chapter, we illustrated the process of converting queries into a physical query plan for DBMS. This included various steps such as,

- Convert the given SQL query into a parse tree using the language specific grammar.
- Next have to check if the given parse tree is actually a member of the grammar, that is, semantic checking.
- Substituting the nodes of the parse tree with the proper operators for conversion to a logical plan.
- Next is optimizing the logical query plan by making the algebraic transformation from relation algebra.
- To prepare for cost based search, need to have statistics ready, so have to calculate them, E.g. histogram, the tuple size and more as mentioned earlier.
- Deciding a strategy for the joining and enumeration strategy.
- Lastly for execution, decide between pipeline and materialization.

---

**Algorithm 1** Deep Q-Network (DQN)

---

*Input:* MDP  $(S, A, P, R, \gamma)$ , replay memory  $M$ , number of iterations  $T$ , minibatch size  $n$ , exploration probability  $\epsilon \in (0, 1)$ , a family of deep  $Q$ -networks  $Q_\theta : S \times A \rightarrow R$ , an integer  $T_{\text{target}}$  for updating the target network, and a sequence of stepsizes  $\{\alpha_t\}_{t \geq 0}$

Initialize the replay memory  $M$  to be empty.

Initialize the  $Q$ -network with random weights  $\theta$ .

Initialize the weights of the target network with  $\theta^* = \theta$

Initialize the initial state  $S_0$ .

**for**  $t = 0, 1, \dots, T$  **do**

With probability  $\epsilon$ , choose  $A_t$  uniformly at random from  $A$ , and with probability  $1 - \epsilon$ , choose  $A_t$  such that  $Q_\theta(S_t, A_t) = \max_{a \in A} Q_\theta(S_t, a)$ .

Execute  $A_t$  and observe reward  $R_t$  and the next state  $S_{t+1}$ .

Store transition  $(S_t, A_t, R_t, S_{t+1})$  in  $M$ .

Experience replay: Sample random minibatch of transitions  $\{(s_i, a_i, r_i, s'_i)\}_{i \in [n]}$

For each  $i \in [n]$ , compute the target  $Y_i = r_i + \gamma * \max_{a \in A} Q_{\theta^*}(s'_i, a)$

Update the  $Q$ -network: Perform a gradient descent step

$$\theta \leftarrow \theta - \alpha_t * \frac{1}{n} \sum_{i \in [n]} [Y_i - Q_\theta(s_i, a_i)] * \nabla_\theta Q_\theta(s_i, a_i)$$

Update the target network: Update  $\theta^* \leftarrow \theta$  every  $T_{\text{target}}$  steps.

Define policy  $\bar{\pi}$  as the greedy policy with respect to  $Q_\theta$

**Output** Action-value function  $Q_\theta$  and policy  $\bar{\pi}$

---

Then we looked at DSMS and saw the challenges they face and how they overcome them. This essentially is how the metrics used to measure QoS are inter-related and the windowing method.

## Chapter 3

# Stream Optimization

### 3.1 Introduction

A stream is an ordered sequence of data items, which are values that can range from simple numbers to flat tuples to more elaborate structured data that may be deeply nested and have variable size. Streams are conceptually infinite, in the sense that as the streaming computation unfolds over time, the sequence of data items is unbounded in length (Martin Hirzel, 2018).

Stream query optimization is the process of modifying a stream processing query, often by changing its graph topology and/or operators, with the aim of achieving better performance (such as higher throughput, lower latency, or reduced resource usage), while preserving the semantics of the original query.

An optimization should be both safe and profitable. An optimization is safe if it can be applied to a stream query without changing what it computes, as determined by the users requirements. An optimization is profitable if it makes the stream query faster, as measured by metrics that matter to the user, such as throughput, latency, or resource efficiency. There is a substantial literature on different stream query optimizations, with different safety and profitability characteristics. This entry lists the most common optimizations along with short descriptions.

Possible areas of optimization are, batch size, operation combining/ dividing, memory assignment, message passing, operation reordering, garbage collection. We focus mostly on operation reordering.



Last chapter showcased how Deep reinforcement learning can help solve problems by telling the best course of action. We are going to use this technique to improve our operation reordering.

### 3.2 Formalization Requirements

To even being to tackle the problem first we need to define the problem formally to be able to analyze it.

For the experiment to be reconstructible to the closest possible degree and to define a concrete problem statement for reinforcement learning, we will need to decide upon the following at least:-

- Data Soucre generator
- Query to execute
- The underlying algorithms have to be fixed
- Featurization method of DRL
- How will DRL change over time
- Output of DRL
- Evaluation

Before going on the give answers to the above, we also need to justify using this(DRL) technique.

### 3.3 Problem Statement

As discussed in the previous chapter, given a query there are multiple ways of executing the query to return the result, these ways correspond to a query plan and for

each query plan there is a cost associated to it, one component of the costs is the time required for execution.

We can represent the plans by  $P_i$  and the time corresponding to the plan as  $C_i$ . So we can represent them as  $\{(P_i, C_i)\}_{i=1}^n$ . Our goal is to find query plan  $P_i$  with the minimum  $C_i$ .

### 3.4 Approach

Currently we make the following assumptions:-

- We have our data source generator with say a fixed seed, which will ensure we have a constant source of data.
- The query to execute is given, that is the window size for data streams is given.
- The query has combination of data streams and static relations as well as operators such as selection, projections, groupby and more.

As mentioned earlier, we are going to focus on finding an order of operations to execute to reduce/ minimize the query execution time.

Say there is a data streams  $S$  with  $m$  attributes and relations  $R_1, R_2, \dots, R_n$  each with  $m_i$  attributes.

That is to say there are total of  $(m + \sum_{i=1}^n m_i) = g$  attributes. We can use 1 – 0 encoding so that any attribute or a combination in this schema can be represented by an array/vector  $v$  of length  $g$ .

$$v_i = \begin{cases} 1 & \text{if } i \in \text{combination} \\ 0 & \text{otherwise} \end{cases}$$

This vector will help represent attributes listed in the query. We can further replace 1 with the entropy of the attribute.

For our Deep reinforcement learning data we will need the rewards, actions, current and next possible state, we will need to encode these in a 4-tuple,  $(F_c, F_n, \text{Action}, \text{Reward})$ . When we have enough of these examples we are use them to train our model and try to get a better ordering.

## 3.5 Example

We look at linear road benchmark to test our hypothesis.

### 3.5.1 Input

**CarLocStr:** Stream of car location reports. This forms primary input to the system.

```

1  CarLocStr(car_id,          /* unique car identifier      */
2                      speed, /* speed of the car        */
3                      exp_way, /* expressway: 0..10       */
4                      lane,   /* lane: 0,1,2,3           */
5                      dir,    /* direction: 0(east), 1(west) */
6                      x-pos); /* coordinate in express way */

```

**AccBalQueryStr:** Stream of account-balance adhoc queries. Each query requests the current account balance of a car.

```

1  AccBalQueryStr(car_id,
2                  query_id); /* id used to associate
3                          * responses with queries */

```

**ExpQueryStr:** Stream of adhoc queries requesting the expenditure of a car for the current day.

```

1  ExpQueryStr(car_id,
2              query_id);

```

**TravelTimeQueryStr:** Stream of expected-travel-time adhoc queries.

```

1  TravelTimeQueryStr(query_id,
2                      exp_way,
3                      init_seg, /* initial segment */
4                      fin_seg, /* final segment  */
5                      time_of_day,
6                      day_of_week);

```

**CreditStr:** The stream of credits to the account corresponding to a car.

```
1 CreditStream(car_id, credit);
```

Instead of assuming multiple streams for input we will combine them and assuming a single data stream as input.

### 3.5.2 Query

In this section we showcase the query statement for a simple and a complex query, convert it into a SQL query and show the complexity within. We use the query statement from linear road itself.

- **Query 1, Segment Average Speed** : The average speed of a particular car in a segment over the last 5 minutes.
- **Query 2, Toll Computation for Segments** : The toll for each segment depends on the average speed and volume of the cars in the segment, and on the presence of accidents in downstream segments.

Now to convert the query statement to actual syntax/ query and showcase few different query plans and how the cost might change.

### 3.5.3 Segment Average Speed

First we convert the simpler query.

```
1 SELECT exp_way, dir, seg, AVG(speed) as speed,
2 FROM CarSegStr [RANGE 5 MINUTES]
3 WHERE car_id = target
4 GROUP BY exp_way, dir, seg;
```

To convert this into a query plan, as seen in the previous chapter, we use relational algebra to convert into a query plan.

$$\gamma_{\text{exp\_way, dir, seg, AVG(speed)}}(\sigma_{\text{car\_id} = \text{target}}(\text{CarSegStr}))$$

Of course, the above mentioned method/ bracketing is one way of executing the query. A small modification from the rules in previous chapter can lead to the following improvement!

$$\gamma_{\text{exp\_way,dir,seg}, \text{AVG(speed)}}(\pi_{\text{exp\_way,dir,seg}, \text{speed}}(\sigma_{\text{car\_id} = \text{target}}(\text{CarSegStr})))$$

Now here we are assuming that *CarSegStr* is a single stream, ready to be used, incase it was a more complicated, that is, it was obtained by joining 2 different streams/ relations with a condition *C* a further improvement is possible!

$$\gamma_{\text{exp\_way,dir,seg}, \text{AVG(speed)}}(\pi_{\text{exp\_way,dir,seg}, \text{speed}}(\sigma_{\text{car\_id} = \text{target}}(S_1 \bowtie_{C \wedge (\text{car\_id} = \text{target})} S_2))))$$

Of course further improvements can be made, but that will require us to look into attributes of the two streams  $S_1$  and  $S_2$

So to conclude, we made an improvement from the naive plan to a much more efficient plan

- $\gamma_{\text{exp\_way,dir,seg}, \text{AVG(speed)}}(\sigma_{\text{car\_id} = \text{target}}(\text{CarSegStr}))$
- $\gamma_{\text{exp\_way,dir,seg}, \text{AVG(speed)}}(\pi_{\text{exp\_way,dir,seg}, \text{speed}}(\sigma_{\text{car\_id} = \text{target}}(\text{CarSegStr})))$
- $\gamma_{\text{exp\_way,dir,seg}, \text{AVG(speed)}}(\pi_{\text{exp\_way,dir,seg}, \text{speed}}(\sigma_{\text{car\_id} = \text{target}}(S_1 \bowtie_{C \wedge (\text{car\_id} = \text{target})} S_2))))$

### 3.5.4 Toll Computation for Segments

Now to convert tackle the more complex query, **Query 2**,

Before going to that, note the steps to perform a query on streams requires to first generate a relation from stream, make relations from this derived relation and then output a stream from the derived relation. This is a 3 step process.

Now we look at the more complex query.

**SegAvgSpeed:** The average speed of the cars in a segment over the last 5 minutes.

```
1 SELECT exp_way, dir, seg, AVG(speed) as speed,
2 FROM CarSegStr [RANGE 5 MINUTES]
3 GROUP BY exp_way, dir, seg;
```

**SegVol:** Relation containing the number of cars currently in a segment. The relation CurCarSeg is used to determine the cars in each segment.

```
1 SELECT exp_way, dir, seg, COUNT(*) as volume
2 FROM CurCarSeg
3 GROUP BY exp_way, dir, seg;
```

**SegToll:** The toll for each segment. There are no entries in the relation for segments having no toll. A segment is tolled only if the average speed of the segment is less than 40, and if it is not affected by an accident. If a segment is tolled, its toll is  $\text{basetoll} * (\#cars - 150) * (\#cars - 150)$ . We modify this query to remove the not affected by accident condition.

```
1 SELECT S.exp_way, S.dir, S.seg, basetoll*(V.volume-150)*(V.volume
   -150)
2 FROM SegAvgSpeed as S, SegVol as V
3 WHERE S.exp_way = V.exp_way and S.dir = V.dir and S.seg = V.seg
4 and S.speed <= 40;
```

These 3 steps combined give us the desired query. Where  $S$  in turn calculated via this process

$$\pi_{\text{exp\_way, dir, seg, speed}}(\gamma_{\text{exp\_way, dir, seg}}(\text{CarSegStr}))$$

Where  $S$  in turn calculated via this process

$$\pi_{\text{exp\_way, dir, seg, column}}(\gamma_{\text{exp\_way, dir, seg}}(\text{CurCarSeg}))$$

Finally combining the two with the  $3^{rd}$  query

$$\pi_{\text{S.exp\_way, S.dir, S.seg, S.toll}}$$

$$\sigma(\text{S.exp\_way} = \text{V.exp\_way}) \wedge (\text{S.dir} = \text{V.dir}) \wedge (\text{S.seg} = \text{V.seg}) \wedge (\text{S.speed} \leq 40)(\text{S}, \text{V})$$

We can push down the projection and then break down the selection condition, for ease of notation, substitute the 4 conditions as  $C_1, C_2, C_3, C_4$  respectively.

$$\pi_{\text{S.exp\_way}, \text{S.dir}, \text{S.seg}, \text{S.toll}}(\sigma_{C_1}(\sigma_{C_2}(\sigma_{C_3}(\sigma_{C_4}(\text{S}, \text{V}))))$$

Now we can list few different query plans

$$\pi_{\text{S.exp\_way}, \text{S.dir}, \text{S.seg}, \text{S.toll}}(\sigma_{C_1}(\sigma_{C_2}(\sigma_{C_3}(\sigma_{C_4}(\text{S}, \text{V}))))$$

$$\pi_{\text{S.exp\_way}, \text{S.dir}, \text{S.seg}, \text{S.toll}}(\sigma_{C_1}(\sigma_{C_2}(\sigma_{C_4}(\sigma_{C_3}(\text{S}, \text{V}))))$$

$$\pi_{\text{S.exp\_way}, \text{S.dir}, \text{S.seg}, \text{S.toll}}(\sigma_{C_1}(\sigma_{C_3}(\sigma_{C_2}(\sigma_{C_4}(\text{S}, \text{V}))))$$

The goal is to try to figure out the best possible method of rearranging these associative operations, in this thesis we try out using deep reinforcement learning to try to predict the best possible move.

## Chapter 4

# Implementation

This chapter explores the setup required to conduct the experiments. The chapter is divided into

- Data Generation:- We generate the data stream for the Linear road benchmark.
- Query Execution:- We have written a C++ code to execute the query and extract information while executing.
- Deep Reinforcement Learning:- Written a Deep reinforcement learning agent to help optimize the order of operations.

### 4.1 Data Generation

The data is generated using the walmart linear road code generator without mistim(Cory Stegemoller, 2016).

Linear Road tests Stream Data Management Systems (SDMS) by measuring how many expressways a SDMS can support by giving accurate and timely results to four types of queries that fit two categories: continuous and historical.(Cherniack et al., 2004)

```
1 java com.walmart.linearroad.generator.LinearGen [-o <output file>]
    [-x <number of xways>] [-m <dummy value to activate multi-
    threading>]
```

**Listing 4.1:** generate linear road data

Example output of this code



```

1 0,0,13,10,8,0,0,89,469920,-1,-1,-1,-1,-1,-1
2 0,0,17,10,8,0,1,65,348479,-1,-1,-1,-1,-1,-1
3 0,0,22,10,8,0,0,12,63360,-1,-1,-1,-1,-1,-1
4 0,0,33,10,8,0,1,94,501599,-1,-1,-1,-1,-1,-1
5 0,0,42,10,8,0,0,14,73920,-1,-1,-1,-1,-1,-1
6 0,0,4,10,7,0,0,61,322080,-1,-1,-1,-1,-1,-1
7 0,0,85,10,8,0,1,30,163679,-1,-1,-1,-1,-1,-1
8 0,0,11,10,6,0,1,41,221759,-1,-1,-1,-1,-1,-1
9 0,0,23,10,7,0,1,81,432959,-1,-1,-1,-1,-1,-1
10 0,0,15,10,6,0,0,5,26400,-1,-1,-1,-1,-1,-1

```

**Listing 4.2:** linear road data example output

Each line in the example indicates 1 data entry, the schema used in the generation of the above code as given.

#### 4.1.1 Schema

The above generated data can be interpreted as follows:-

Column1	Tells the type of query. 0: position report 2: account balance request 3: daily expenditure request 4: travel time request
Column2	Timestamp position.
Column3	Vehicle identification number
Column4	Speed of the vehicle
Column5	Express way number
Column6	Lane ID (0, ..., 4)
Column7	Direction of movement(0 =East or 1 =West)
Column8	Segment ID (0, ..., 99)
Column9	Position of the vehicle.
Column10	Query identifier
Column11	Start Segment
Column12	End Segment
Column13	Day of the week
Column14	Minute of the day
Column15	Day in the past 10 weeks

## 4.2 Query Execution

After the data is prepared, we need to execute the queries listed and explained in the previous chapter. We simplified and presented multiple plans of execution and we will explore them.

As suggested, a SDMS can take in multiple data stream as input but we were going to treat it as a single input. Below is the implementation of the 2 queries in *C++*.

### 4.2.1 Simpler Query

The simpler query

```

1  SELECT exp_way, dir, seg, AVG(speed) as speed,
2  FROM CarSegStr [RANGE 5 MINUTES]
3  WHERE car_id = target
4  GROUP BY exp_way, dir, seg;
```

**Listing 4.3:** simple query

The goal of implementing this query is to get a sense of time required for executing queries, the time for taking input and time for providing an output.

The code can be broken into 3 parts:-

- Initialization, code overview
- Reading input
- executing the query

#### A.1

To start, the input file is specified, the window size is specified and the start time of the clock is initialized. The input file is treated as a stream. While many things are to be considered while taking input, as the focus of this thesis was in other area, we simply read from the file part by part.

We continue to read the file until we reach the end of it, everytime we have read the

window size amount of data points, we execute the query.

A window size was fixed which determined the number of entries to read at once, the window size was interpreted in two ways,

- The number of entries to read
- The time interval for which to read the entries

#### A.2

This information along with the value was passed to the read function A.2 everytime.

This input function treats the window size as specifying the number of entries to read.

The input was taken by simply reading a file and storing the attribute values important to the query. We only store the information necessary for our queries. This was done due to limited RAM and for expediting the process.

#### A.3

This part prepares the data for aggregation functions. The data is first inserted into a vector. It then essentially groups the rows together depending on the order provided by the means of the sort function(my\_sort in this case). The sort function helps group up entries with same express way ID, segment ID and which are going in the same direction.

#### A.4

Finally do the aggregation, simply check for the continuity of the group and update the measures of the group as you iterate. When you arrive the the end of the group simply insert the answer for the aggregate function along with the rest of required attributes into an output and return the table.

### 4.2.2 Complex Query

The complex query(Widom et al., 2004)

```

1 SELECT car_id, speed, exp_way, lane, dir, (x-pos/52800) as seg
2 FROM CarLocStr;

```

**Listing 4.4:** CarSegStr

```

1 SELECT car_id, exp_way, dir, seg
2 FROM CarSegStr [PARTITION BY car_id ROWS 1], CurActiveCars
3 WHERE CarSegStr.car_id = CurActiveCars.car_id;

```

**Listing 4.5:** CurCarSeg

```

1 SELECT exp_way, dir, seg, AVG(speed) as speed,
2 FROM CarSegStr [RANGE 5 MINUTES]
3 GROUP BY exp_way, dir, seg;

```

**Listing 4.6:** SegAvgSpeed

```

1 SELECT exp_way, dir, seg, COUNT(*) as volume
2 FROM CurCarSeg
3 GROUP BY exp_way, dir, seg;

```

**Listing 4.7:** SegVol

```

1 SELECT S.exp_way, S.dir, S.seg, basetoll*(V.volume-150)*(V.volume
   -150)
2 FROM SegAvgSpeed as S, SegVol as V
3 WHERE S.exp_way = V.exp_way and S.dir = V.dir and S.seg = V.seg
4      and S.speed <= 40;

```

**Listing 4.8:** SegToll

The goal of implementing this query is to obtain data to train a deep reinforcement learning model and to be able to test the model. We start by giving the general outline of the code.

## A.5

Similar to the simple query, we fix a window size, the data file and we start the loop

In the loop, each time we read "window size" number of data entries, calculate 3 vectors, **curcarseg**, **segavgspeed**, **segvol**. After calculating these for the the data in the data window, we execute the final query of **segtoll**

For *CurCarSeg*, simply filter the cars for the last 5 mins.

A.6

For *SegAvgSpeed*, we first filter the cars which were active in the last 5 mins. Then sort to group data entries. Then lastly apply aggregate function.

A.7

Next, calculate *SegVol*, first sort the vector by the attributes in the group by statement. Then apply the aggregation operation. After obtaining both *SegVol* and *SegAvgSpeed*, finally proceed to calculate *SegToll* and obtain data for reinforcement learning.

A.8

Given *SegVol* and *SegAvgSpeed*, need to do the following steps:-

- Calculate column wise entropy and store it.
- Store size of *SegVol* and *SegAvgSpeed*.
- For all  $4! = 24$  orderings of operations(due to them being associative) execute the *SegToll* query.
- Store the time require for execution and the number of operations required for each order.

Here is a part of output generated by the code

The first line represents the entropy of the 8 columns,

The next line represents the size of *SegAvgSpeed* and *SegVol* respectively.

Then for the next 48 lines, it alternates between the order of operations executed and then the time required in microseconds and the number of operations required.

```
1 3.32187 0.99998 6.63484 3.29218 2.92342 3.32185 0.99998 6.63482
2 1916 1914
```

```

3 0 1 2 3
4 171487 3672963
5 0 1 3 2
6 175067 3672963
7 0 2 1 3
8 176562 3672963
9 0 2 3 1
10 175096 3672963

```

**Listing 4.9:** Data for DQN

The example above can be interpreted as

The entropy of columns are 3.32187, 0.99998, 6.63484, 3.29218, 2.92342, 3.32185, 0.99998, 6.63482

The size of *SegAvgSpeed* is 1916, the size of *SegVol* is 1914.

The order line can be interpreted as

- The first operation executed is the  $0^{th}$ (S.exp\_way = V.exp\_way).
- The second operation executed is  $1^{st}$ (S.dir = V.dir).
- The third operation executed is  $2^{nd}$ (S.seg = V.seg).
- The last operation executed is  $3^{rd}$ (S.speed <= 40)

When this order is executed, the time taken is 171487 mircoseconds and the number of operations required is 3672963

Next 2 lines can be interpreted as

- The first operation executed is the  $0^{th}$ (S.exp\_way = V.exp\_way).
- The second operation executed is  $1^{st}$ (S.dir = V.dir).
- The third operation executed is  $3^{rd}$ (S.speed <= 40).
- The last operation executed is  $2^{nd}$ (S.seg = V.seg)

When this order is executed, the time taken is 175067 mircoseconds and the number of operations required is 3672963.

The fact that the number of operations remains same but the time required changes is due to the internal scheduling by the kernel.

### 4.3 Deep Reinforcement Learning

Now the data for Deep reinforcement learning is stored and ready to be used.

#### 4.3.1 MDP input

Note after selecting an order and executing the query with the selected order, we finish the "game".

Remember the input for  $MDP$  is of the format  $(S, A, S^*, R)$ . But as this game is 1 move game, all tuples will have  $S^* = s_f$

#### 4.3.2 States

The input for a DQN is of the format

$$(X_i, Y_i)_{i \in [n]}, \text{ where } X_i \in \mathbb{R}^m, Y_i \in \mathbb{R}$$

For  $Y_i$  we use the number of operations required.

For  $X_i$  we use a concatenation of entropy vector, size vector and the 1-hot encoding of the move.

#### 4.3.3 1-hot encoding

Consider the set of permutations of  $\{0, 1, 2, 3\}$ , we can order them lexicgraphically.

Given a permutation  $\mathbb{P}$ , it will have a rank  $r$  in the ordering.

To convert the permutation into a 1hot encoding, start with a 24 element vector  $v$

with all 0s. Then update

$$v[r] = 1$$

A.9

#### 4.3.4 Assumption

Note the amount of information, features extracted from the querying do not fully represent the data set our Neural network can't perform very well. To tackle this, we look at whether the shift in the predicted rewards is similar cross all the moves.

Then given data points  $(x_i, y_i), (x_j, y_j)$  and the DQN  $Q_\theta$  the condition to check becomes:-

$$y_i \leq y_j \Rightarrow Q_\theta(x_i) \leq Q_\theta(x_j)$$

If this condition holds, then the predicted optimal move is the actual optimal move.

As optimal move implies

$$y_{\text{optimal}} \leq y_i \forall i \in [24]$$

This with the previous condition will imply

$$Q_\theta(x_{\text{optimal}}) \leq Q_\theta(x_i) \forall i \in [24]$$

We train a Deep neural network with 70%, use the rest of the 30% data for testing and record the results.

## 4.4 Conclusion

In this chapter we saw :-

- The method use to generate data for the linear road benchmark test cases.
- The method implemented to take input, converted *SQL* queries to *C++* and extracted features during execution and stored them



- Showcased how the extracted features can help with deep reinforcement learning
- Showcased the method used to evaluate the performance of DQN.

In the next chapter, we

- Formalize the assumptions we made in our evaluation.
- Present the categorize/ classification of the data set. Divided into training and testing data.
- Visualize the confusion matrix.
- Give the true positive, true negative, false positive false negative as well the analysis of the confusion matrix for all 24 orderings.
- Lastly we compare how our predictions fared against the optimal and worst orderings.

present the results of our experiments, visualize them The method to compute the results are given in the next chapter.

## Chapter 5

# Evaluation

### 5.1 Optimality checking

Now we have learnt a model  $Q_\theta$ , we also have our test cases  $\{(x_i, y_i)\}_{i \in [n]}$ .

Each data window corresponds to 24 data entries. The  $t^{th}$  window corresponds to

$$\{x_{24t}, x_{24t+1}, \dots, x_{24t+23} | t \in \mathbb{I}\}$$

Each data window has 24 data points as we have tried all possible(24 of them) combinations of ordering of operations. For each data window we need to find the ordering of operations which requires the minimum moves according to our predictions.

To do this for the  $t^{th}$  window of data, find the index of the minimum of

$$Q_\theta(x_{24t}), Q_\theta(x_{24t+1}), \dots, Q_\theta(x_{24t+23})$$

, say  $k$ .

To check if the move this corresponds to is the actual optimal move, find the index of the minimum of  $y_{24t}, y_{24t+1}, \dots, y_{24t+23}$ , say  $l$ . If  $k = l$  then the move we predicted as optimal is indeed optimal.

We train the model multiple times on the training data and measure its performance on the test data each time.

To measure the performance we use the predicted optimal move vs the actual optimal

move.

## 5.2 Results

First, for each ordering of the operations we present the frequencies with which they are optimal in the dataset.

```
1 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 110, 83, 84, 79, 80, 63, 307,
   165, 2209, 1367, 43849, 43872]
```

As visible the data is highly biased towards the last 2 orderings.

We split this data into 70%, 30% for the purpose of training and testing of our reinforcement learning agent. The random seed used resulted in the following split.

```
1 #Training
2 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 81, 60, 59, 54, 57, 45, 227,
   119, 1526, 966, 30811, 30583]
3 #Testing
4 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 29, 23, 25, 25, 23, 18, 80, 46,
   683, 401, 13038, 13289]
```

The predictions of our model resulted in the following 24 true positive, False Positive, False negative and True negative results

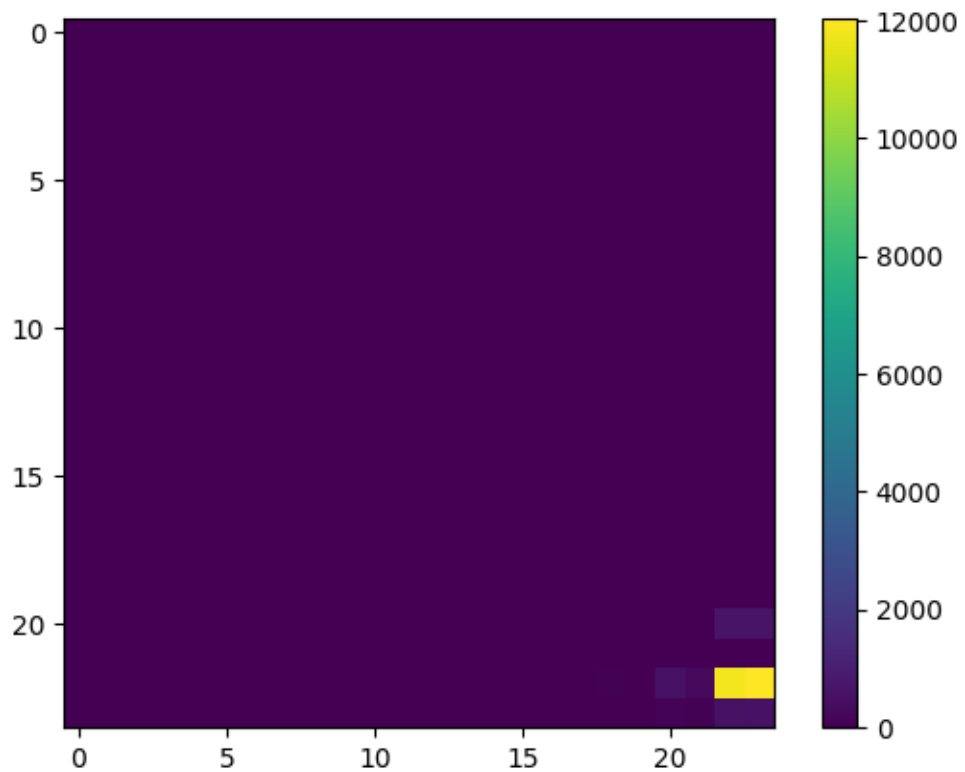
```
1 (TP, FP, FN, TN)
2 (0.0, 0.0, 0.0, 27681.0)
3 (0.0, 0.0, 0.0, 27681.0)
4 (0.0, 0.0, 0.0, 27681.0)
5 (0.0, 0.0, 0.0, 27681.0)
6 (0.0, 0.0, 0.0, 27681.0)
7 (0.0, 0.0, 0.0, 27681.0)
8 (0.0, 0.0, 0.0, 27681.0)
9 (0.0, 0.0, 0.0, 27681.0)
10 (0.0, 0.0, 0.0, 27681.0)
11 (0.0, 0.0, 1.0, 27680.0)
12 (0.0, 0.0, 0.0, 27681.0)
13 (0.0, 0.0, 0.0, 27681.0)
14 (0.0, 0.0, 29.0, 27652.0)
```

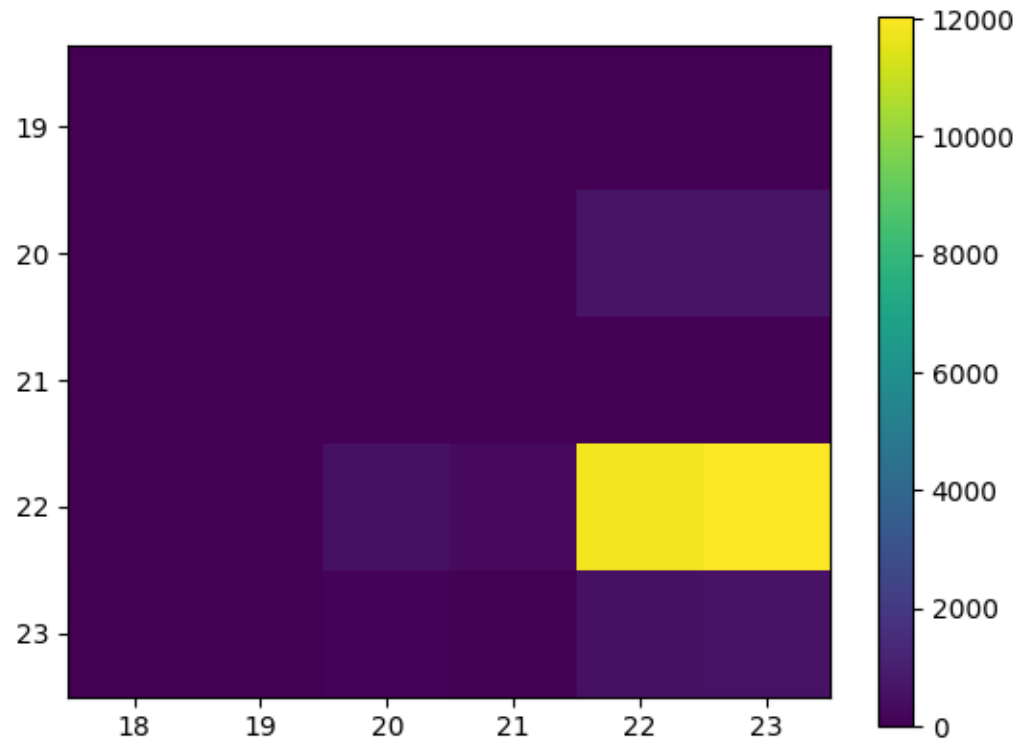
```

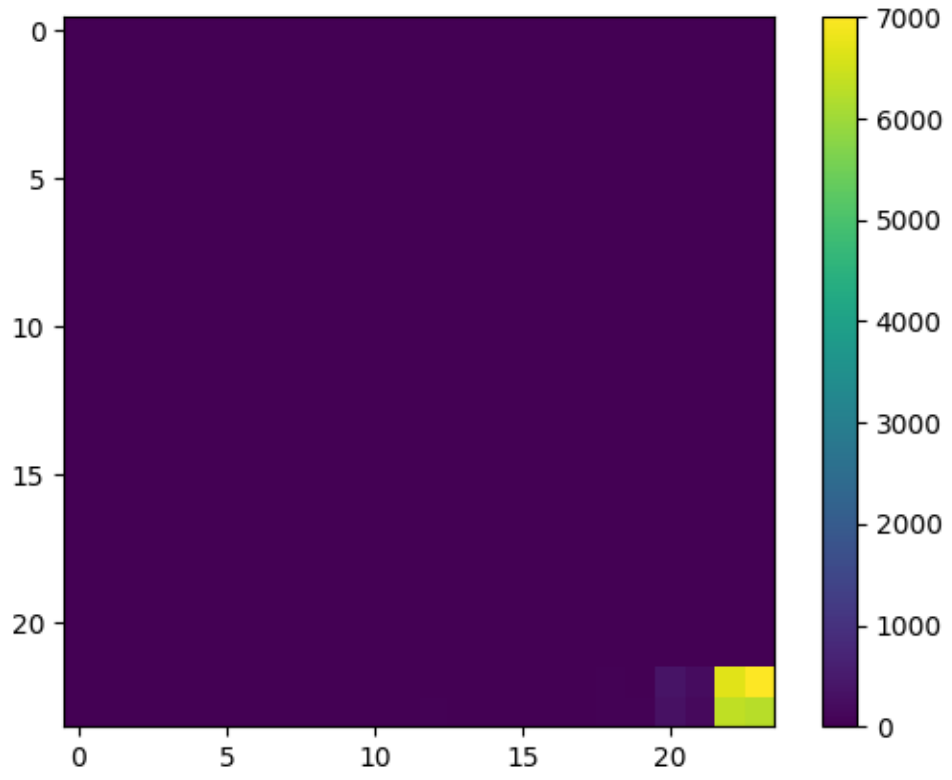
15 (0.0, 0.0, 23.0, 27658.0)
16 (0.0, 0.0, 25.0, 27656.0)
17 (0.0, 0.0, 25.0, 27656.0)
18 (0.0, 0.0, 23.0, 27658.0)
19 (0.0, 0.0, 18.0, 27663.0)
20 (0.0, 0.0, 80.0, 27601.0)
21 (0.0, 0.0, 46.0, 27635.0)
22 (0.0, 0.0, 683.0, 26998.0)
23 (0.0, 0.0, 401.0, 27280.0)
24 (6705.0, 7667.0, 6333.0, 6976.0)
25 (6273.0, 7036.0, 7016.0, 7356.0)

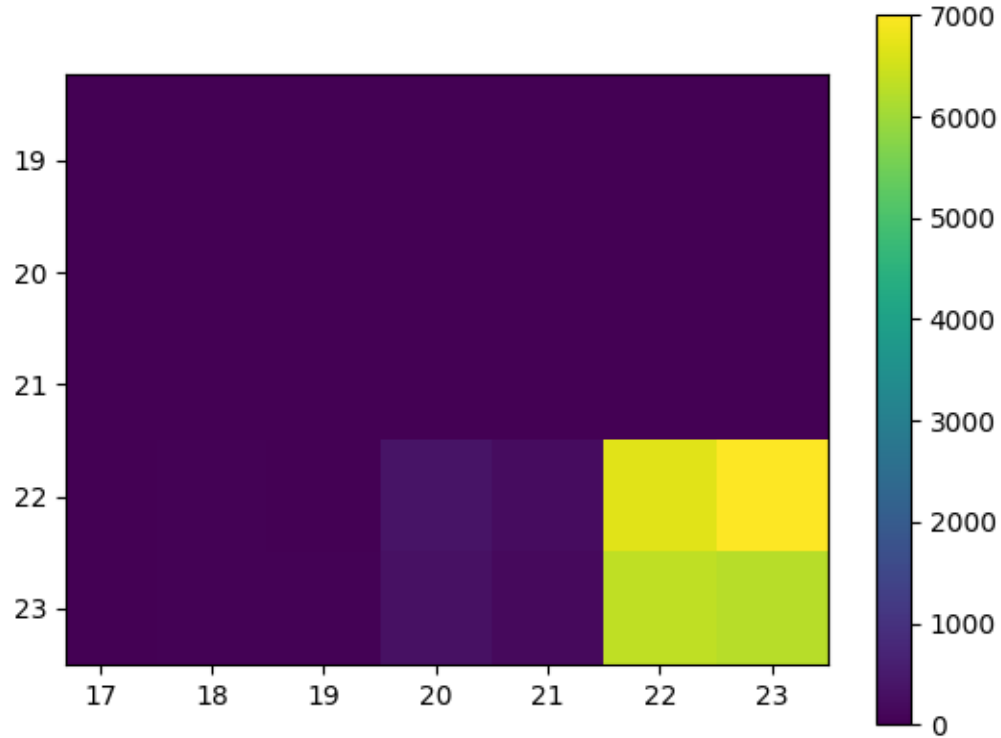
```

So we got 12978 predictions correct out of 27681 test cases that is around 47%. Below are 2 confusion matrices we obtained from two different runs of the predictions.









We have the actual optimal ordering, the actual worst ordering, the predicted best ordering from the model.

We add the number of operations required by following each of the above 3 and look at the sum of them.

The sum of operations required by the optimal ordering 890640075.0

The sum of operations required by the predicted ordering 900361559.0

The sum of operations required by the optimal ordering 2310227856.0

That is our predicted orderings are approx 39% of the worst ordering and 102% of the best ordering

A by class analysis is obtained by using `metrics.classification_report` from `scikit-learn` (Pedregosa et al., 2011).

1		precision	recall	f1-score	support
2					
3	9	0.000	0.000	0.000	1
4	12	0.000	0.000	0.000	29
5	13	0.000	0.000	0.000	23
6	14	0.000	0.000	0.000	25
7	15	0.000	0.000	0.000	25
8	16	0.000	0.000	0.000	23
9	17	0.000	0.000	0.000	18
10	18	0.000	0.000	0.000	80
11	19	0.000	0.000	0.000	46
12	20	0.000	0.000	0.000	683
13	21	0.000	0.000	0.000	401
14	22	0.467	0.514	0.489	13038
15	23	0.471	0.472	0.472	13289
16					
17	micro avg	0.469	0.469	0.469	27681
18	macro avg	0.072	0.076	0.074	27681
19	weighted avg	0.446	0.469	0.457	27681

### 5.3 Discussion

### 5.4 Conclusion



## Chapter 6

# Conclusion and Further work

### 6.1 Conclusion

The goal of this thesis is to act as a proof of concept for the application of Deep Reinforcement Learning for optimization of query processing on streams.

Say, for a particular window of data

- The size of SegVol =  $x$ .
- The size of SegAvgSpeed =  $y$ .

Due to the query used in the experiments, the number of operations required to execute the query lie between

$$[x * y, 4 * x * y]$$

The worst and best possible cases can both arise in the same window depending on the order used and the data in the query.

Of the 24 possible operation orders, we were able to predict the optimal order of operations around 48% cases, reducing the number of operations to around 39% of the worst observed(which is better than the actual worst).

$$x * y \leq \text{optimal} \leq \text{predicted} \leq \text{worst observed} \leq 4 * x * y$$

## **6.2 Further Work**

This thesis proposed a method for optimization of query execution on data stream and provided a demo implementation for experimentation. There are still many areas which can be improved and developed further.

### **6.2.1 Data Generation**

As shown in chapter 5, the data generated is highly biased. A source producing more diverse data and a query on that may lead to more interesting results. Real world data stream can be used to get better variety in the data and take input via different methods.

### **6.2.2 Data Storage**

Currently, for executing the query, we store data in vectors where as SQL uses B-trees, a change like this can impact the time required for query execution and then lead to a different learned model if the time of execution is used as reward.

On a further note data storage on systems like AWS S3 buckets, filesystems like Hadoop, hdfs and others can lead to changes in time of execution as well.

### **6.2.3 Features extraction**

The current features only include the entropy of the columns and the size of the tables. Whereas in traditional SQL with a static database(relatively) a lot more features of generally extracted. Better starting feature can perhaps lead to a better results.

### **6.2.4 Deep Reinforcement Learning**

The Deep reinforcement learning model used was rather simplified due to the less number of moves and convert the game into a single step game. Further intensive

state, action and reward tuples can be constructed to get better intermediate rewards and policies.

#### **6.2.5 Integration**

The learnt DQN can be integrated with an actual stream processing system and tested in full for scalability, reliability and more. Need to note that there is a speedup by number of operations, the runtime due to various other network constraints might have a different result.

## Appendix A

### Extra code

```

1 FILE *pFile;
2 int window_size=20000;
3 pFile = fopen (argv[1], "r");
4 int database[window_size][4];
5 std::vector<int> vect;
6 vect.clear();
7 std::vector<std::vector<int>> output;
8 auto start = std::chrono::high_resolution_clock::now();
9 while(!feof(pFile))
10 {
11     input_database(pFile, database, window_size);
12     execute(database, output, vect, window_size);
13     output.clear();
14     vect.clear();
15 }
16 auto stop = std::chrono::high_resolution_clock::now();
17 auto duration = std::chrono::duration_cast<std::chrono::microseconds>
    >(stop - start);
18 std::cout << duration.count() << std::endl;
19 fclose(pFile);

```

**Listing A.1:** code overview

```

1 void input_database(FILE * pFile, int d[][4], int window_size)
2 {
3     int i,j;
4     for(i=0;i<window_size;i++)
5     {
6         fscanf (pFile, "%i,%i,%i,%i,%i,%i,%i,%i,%i,%i,%i,%i,%i,%i,%i",
            &j,
7         &j,&j,&d[i][0],&d[i][1],&j,&d[i][2],&d[i][3],&j,&j,&j,&j,&j,
            &j,&j);
8     }

```

```
9 }
```

### Listing A.2: input function

```
1 for(i=0;i<window_size;i++)
2 {
3     if(d[i][2]==0)
4     {
5         vect.push_back(d[i][1]);
6         vect.push_back(d[i][2]);
7         vect.push_back(d[i][3]);
8         vect.push_back(d[i][0]);
9         output.push_back(vect);
10    }
11    vect.clear();
12 }
13 std::sort(output.begin(),output.end(), my_sort);
```

### Listing A.3: groupby

```
1 int count = 1;
2 std::sort(output.begin(),output.end(), my_sort);
3
4 std::vector<std::vector<int>>> op;
5 op.push_back(output[0]);
6 for(i=1;i<output.size();i++)
7 {
8     if((output[i][0]==output[i-count][0]))
9     {
10        if(output[i][1]==output[i-count][1])
11        {
12            if((output[i][2]==output[i-count][2]))
13            {
14                op.back()[3]=op.back()[3]+output[i][3];
15                count++;
16            }
17        }
18    }
19    else
20    {
21        op.back()[3]=op.back()[3]/(count);
22        op.push_back(output[i]);
23        count = 1;
24    }
```

```

25 }
26 op.back()[3]=op.back()[3]/count;
27 output.clear();
28 output=op;
29 op.clear();

```

#### Listing A.4: aggregation

```

1  int window_size=10000;
2  FILE *pFile;
3  pFile = fopen (argv[1],"r");
4
5  std::ofstream op_file;
6  op_file.open("q3out.txt");
7
8  std::vector<std::vector<int>> database, curCarSeg, segAvgSpeed,
   segVol;
9
10 int temp;
11 while(!feof(pFile))
12 {
13
14     input_database(pFile, database, window_size);
15     read_times++;
16     printf("read times is, %d\n", read_times);
17     curcarseg(database, curCarSeg);
18     temp=curCarSeg.size();
19     printf("    CurCarSeg size is, %d\n", temp);
20     segavgspeed(database, segAvgSpeed);
21     temp=segAvgSpeed.size();
22     printf("    SegAvgSpeed size is, %d\n", temp);
23     segvol(curCarSeg, segVol);
24     temp=segVol.size();
25     printf("    SegVol size is, %d\n", temp);
26     segtoll(segAvgSpeed, segVol, op_file);
27
28     curCarSeg.clear();
29     segAvgSpeed.clear();
30     segVol.clear();
31     database.clear();
32 }
33 fclose(pFile);

```

```
34 op_file.close();
```

### Listing A.5: complex query overview

```

1  std::vector<std::vector<int>>> temp;
2  int cur_time=0;
3  int i,j;
4  std::vector<int> vect;
5  vect.resize(4);
6  for(i=0;i<d.size();i++)
7  {
8      if(cur_time<d[i][0])
9      {
10         cur_time=d[i][0];
11     }
12 }
13 for(i=0;i<d.size();i++)
14 {
15     if(d[i][0]>cur_time-300)
16     {
17         vect[0]=d[i][2];
18         vect[1]=d[i][3];
19         vect[2]=d[i][4];
20         vect[3]=d[i][6];
21         temp.push_back(vect);
22     }
23 }
24 //temp=()
25 std::sort(temp.begin(),temp.end(), my_sort);
26 //(express,dir,seg, average speed)
27 int count = 1;
28 s.push_back(temp[0]);
29 for(i=1;i<temp.size();i++)
30 {
31     if(temp[i][0]==temp[i-1][0] and temp[i][1]==temp[i-1][1] and
temp[i][2]==temp[i-1][2])
32     {
33         s.back()[3]=s.back()[3]+temp[i][3];
34         count++;
35     }
36     else
37     {
38         s.back()[3]=s.back()[3]/(count);

```

```

39         s.push_back(temp[i]);
40         count = 1;
41     }
42 }
43 s.back()[3]=s.back()[3]/count;
44 temp.clear();

```

Listing A.6: SegAvgSpeed

```

1 //d=(carid,exp,dir,seg)
2 std::sort(d.begin(),d.end(), my_sort2);
3 int i,j;
4 s.push_back(d[0]);
5 s.back()[0]=1;
6 for(i=1;i<d.size();i++)
7 {
8     if(d[i][3]==d[i-1][3] and d[i][1]==d[i-1][1] and d[i][2]==d[i-1][2])
9     {
10         s.back()[0]++;
11     }
12     else
13     {
14         s.push_back(d[i]);
15         s.back()[0]=1;
16     }
17 }

```

Listing A.7: SegVol

```

1 int c1,r1,i,j,k,operation;
2 std::vector<float> entropy;
3 entropy=entropy_calc(s,v);
4 for(i=0;i<entropy.size();i++)
5 {
6     op_file<<entropy[i]<<" ";
7 }
8 op_file<<std::endl;
9 op_file<<s.size()<<" "<<v.size()<<std::endl;
10 entropy.clear();
11 //v=(count of cars,exp,dir,seg)
12 //s=(express,dir,seg, average speed)
13 int times;
14 std::vector<int> order;

```



```

15 order.push_back(0);
16 order.push_back(1);
17 order.push_back(2);
18 order.push_back(3);
19
20 int max_times=24;
21 int num_op;
22 bool b;
23 for(times=0;times<max_times;times++)
24 {
25     num_op=0;
26     auto start = std::chrono::high_resolution_clock::now();
27     for(i=0;i<s.size();i++)
28     {
29         for(j=0;j<v.size();j++)
30         {
31             b=false;
32             for(k=0;k<4;k++)
33             {
34                 operation=order[k];
35                 num_op++;
36                 switch(operation)
37                 {
38                     case 0:
39                         if(s[i][0]!=v[j][1])
40                         {
41                             b=true;
42                         }
43                     case 1:
44                         if(s[i][1]!=v[j][2])
45                         {
46                             b=true;
47                         }
48                     case 2:
49                         if(s[i][2]!=v[j][3])
50                         {
51                             b=true;
52                         }
53                     case 3:
54                         if(s[i][3]>=40)
55                         {
56                             b=true;

```

```

57         }
58     }
59     if(b)
60     {
61         break;
62     }
63 }
64 }
65 }
66 auto stop = std::chrono::high_resolution_clock::now();
67 auto duration = std::chrono::duration_cast<std::chrono::
microseconds>(stop - start);
68
69 for(i=0;i<4;i++)
70 {
71     op_file<<order[i]<<" ";
72 }
73 op_file<<std::endl<<duration.count()<<" "<<num_op<<std::endl;
74 next_perm(order);
75 }

```

Listing A.8: SegToll

```

1 def reinforcement_learning(fin_train,save_loc,epoch_number):
2     model = tf.keras.Sequential()
3     model.add(tf.keras.layers.Dense(34,activation='relu'))
4     model.add(tf.keras.layers.Dense(68,activation='relu'))
5     model.add(tf.keras.layers.Dense(136,activation='relu'))
6     model.add(tf.keras.layers.Dense(68,activation='relu'))
7     model.add(tf.keras.layers.Dense(34,activation='relu'))
8     model.add(tf.keras.layers.Dense(1))
9     model.compile(optimizer='adam', loss='mse')
10
11
12     times = int(epoch_number)
13     for i in range(times):
14         cm=np.zeros(576).reshape(24,24)
15         history=model.fit(x_tr, y2_tr, batch_size=100, epochs=1,
callbacks=[cp_callback])
16         ans=model.predict(x_te)
17         seen=0
18         correct=0
19         for j in range(len(ans)):

```

```

20         if(seen==0):
21             m1=j
22             m2=j
23         if(ans[j]<ans[m1]):
24             m1=j
25         if(y2_te[j]<y2_te[m2]):
26             m2=j
27         #m1 is the prediction
28         #m2 is the actual
29         if(seen==23):
30             seen=0
31             cm[m1%24,m2%24]+=1
32             if(m1==m2):
33                 correct+=1
34         else:
35             seen+=1
36     total=cm.sum()
37     tp = np.asarray([cm[i,i] for i in range(24)])
38     fp = cm.sum(axis=1) - tp
39     fn = cm.sum(axis=0) - tp
40     tn = np.asarray([total-(tp[i]+fp[i]+fn[i]) for i in range
(24)]
41     for i in range(24):
42         try:
43             a=(tp[i]+tn[i])/(tp[i]+tn[i]+fn[i]+fp[i])
44             p=(tp[i])/(tp[i]+fp[i])
45             r=(tp[i])/(tp[i]+fn[i])
46             f=(2*p*r)/(p+r)
47             print(a,p,r,f)
48         except:
49             print(i)
50
51     print(len(ans)/24,correct)
52     del cm

```

Listing A.9: DQN

## References

- Cherniack, M., Galvez, E., Maskey, A. S., Pedan, I., Ryvkina, E., Arasu, A., Maier, D., Stonebraker, M., and Tibbetts, R. (2004). Linear road.
- Cory Stegemoller, M. M. (2016). Linear road data generator.
- Fan, J., Yang, Z., Xie, Y., and Yang, Z. (2020). A theoretical analysis of deep q-learning.
- Martin Hirzel, Robert Soule, B. G. (2018). Stream query optimization.
- Ng, A. (2020). Reinforcement learning and control.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Widom, J., Motwani, R., Arasu, A., Babcock, B., Babu, S., Srivastava, U., Wang, Z., Cieslewicz, J., Datar, M., Ito, K., McAlister, J., Nishizawa, I., Rosenstein, J., Srinivasan, V., and Thomas, D. (2004). Stream project page for linear road benchmark.
- Xu, H., Li, Y., Tian, Y., Darrell, T., and Ma, T. (2018). Algorithmic framework for model-based reinforcement learning with theoretical guarantees. *CoRR*, abs/1807.03858.