

November 20, 2019  
DRAFT

**Replicated Training in  
Self-Driving Database Management Systems**

Gustavo E. Angulo Mezerhane

CMU-CS-19-129

December 2019

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

**Thesis Committee:**  
Andrew Pavlo, Chair  
David G. Andersen

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

Copyright © 2019 Gustavo E. Angulo Mezerhane

November 20, 2019  
DRAFT

**Keywords:** Database Systems, Replication, Machine Learning

# November 20, 2019

## DRAFT

*Para mis padres Gustavo y Claret*

November 20, 2019  
DRAFT

# November 20, 2019

# DRAFT

## Abstract

Self-driving database management systems (DBMSs) are a new family of DBMSs that can optimize themselves for better performance without human intervention. Self-driving DBMSs use machine learning (ML) models that predict system behaviors and make planning decisions based on the workload the system sees. These ML models are trained using metrics produced by different components running inside the system. Self-driving DBMSs are a challenging environment for these models, as they require a significant amount of training data that must be representative of the specific database the model is running on. To obtain such data, self-driving DBMSs must generate this training data themselves in an online setting. This data generation, however, imposes a performance overhead during query execution.

Many DBMSs operate in a distributed master-replica architecture, where the master node sends new changes to replica nodes that hold up-to-date copies of the database. We propose a novel technique named Replicated Training that utilizes existing database replicas to generate training data for models used in self-driving DBMSs. This approach load balances the expensive task of data collection across the distributed architecture, as opposed to being done entirely in the master node. It also provides the advantage of more diverse training data in the case where replicas are running in different hardware environments. Under Replicated Training, each replica can dynamically control training data collection if it needs more resources to keep up with the master node. To show the effectiveness of our technique, we implement it in NoisePage, a self-driving DBMS, and evaluate it in a distributed environment. Our experiments show that training data collection in a DBMS incurs a noticeable performance overhead in the master node, and using Replicated Training reduces this overhead while still ensuring that replicas keep up with the master. Finally, we show that Replicated Training produces ML models that have accuracies comparable to those trained solely on the master node.

November 20, 2019  
DRAFT

# November 20, 2019

# DRAFT

## Acknowledgments

I love Jello. This stuff has fed me from the day I proposed until the day I defended. Jello is my friend, my sustenance, my very being.

Oh yeah. My advisor is cool too.

In addition, **Catherine Copetas** and **Sharon Burks** should be mentioned in a **large** font in everyone's Acknowledgements section, since they said so. (It is a required part of the thesis formatting guidelines.)

November 20, 2019  
DRAFT



# November 20, 2019

# DRAFT

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	3
1.2	Contribution . . . . .	5
<b>2</b>	<b>Background</b>	<b>7</b>
2.1	Replication . . . . .	7
2.2	Training Data Collection . . . . .	8
<b>3</b>	<b>System Architecture</b>	<b>11</b>
3.1	Transactions . . . . .	11
3.2	Timestamp Manager . . . . .	13
3.3	Logging . . . . .	13
3.3.1	Log Serializer Task . . . . .	14
3.3.2	Log Consumer Tasks . . . . .	16
3.4	Recovery . . . . .	16
3.4.1	Log Record Replaying . . . . .	17
3.4.2	Transaction Replaying . . . . .	18
3.5	Internal NoisePage Protocol . . . . .	20
3.6	Replication . . . . .	21
<b>4</b>	<b>Replicated Training</b>	<b>23</b>

# November 20, 2019

## DRAFT

<b>5</b>	<b>Evaluation</b>	<b>25</b>
5.1	Replication Architecture . . . . .	26
5.1.1	Arrival Rate . . . . .	26
5.1.2	Replication delay over time . . . . .	28
5.2	Metrics Collection . . . . .	29
<b>6</b>	<b>Conclusions</b>	<b>31</b>
	<b>Bibliography</b>	<b>33</b>

# November 20, 2019

# DRAFT

## List of Figures

1.1	<b>Metrics overhead on replication delay in NoisePage</b> - Effects of replication delay when metrics collection is enabled over TPC-C with 4 warehouses and a 10k txns/sec arrival rate . . . . .	4
3.1	<b>Log Manager Architecture</b> - The log manager receives as input log record buffers from transactions, serializes these records, and sends them to different destinations (e.g. disk). . . . .	14
3.2	<b>Log record serialization formats</b> - Along with the log record, additional information must be serialized to ensure replayability. . . . .	15
3.3	<b>Schedule involving DDL changes</b> - This schedule is allowed under <code>SNAPSHOT-ISOLATION</code> , but can create problematic races involving the <code>DROP TABLE</code> command. . . . .	18
3.4	<b>Internal NoisePage Protocol packet types</b> - Packets are minimal to reduce network congestion and speed up packet processing . . . . .	20
3.5	<b>Replication Architecture</b> - While there are other components in NoisePage, this highlights the processes involved in replicating data between a master and replica . . . . .	22
5.1	<b>Sensitivity of replication delay</b> - Measuring the average replication delay with varying arrival rates in NoisePage over TPC-C with 4 warehouses on Type 2 machines. When the arrival rate exceeds the transaction replaying rate, delay sharply increases. The shaded region denotes one standard deviation from the mean for each data point. . . . .	27
5.2	<b>Replication Delay in NoisePage</b> - Measuring the average replication delay in NoisePage over TPC-C with 4 warehouses on Type 2 machines. We average the delay for each second, and plot the average over 10 benchmark runs. The red line shows the running mean of the replication delay. . . . .	28

5.3 **Metrics overhead in NoisePage** - Overhead of metrics collection as we  
scale the number of metrics exported . . . . . 29

November 20, 2019  
DRAFT

## **List of Tables**

November 20, 2019  
DRAFT

# November 20, 2019

# DRAFT

## Chapter 1

### Introduction

Database management systems (DBMSs) are notoriously hard and expensive to manage and tune. Configuration and monitoring of physical database design, knobs, and hardware provisioning can have large effects on the performance of the system. Currently, a database administrator (DBA) has the role of this configuration and monitoring. DBAs, however, are expensive to hire, do routine work, and are limited by the human knowledge they possess. **whats a better way to say that ML can see things a human cant. perhaps talk about why machines are better, and not humans suck, is there a citation for this?**

**Purpose of paragraphs: How people have used ML in databases**

Previous research has explored using machine learning (ML) to automate the runtime behavior of DBMSs. Some systems have put ML at the core, creating a new class of so called self-driving DBMSs Pavlo et al. [2017], Oracle [b]. These systems use ML for high-level decision making, such as workload prediction Ma et al. [2018], or system-wide tasks such as transaction scheduling Kraska et al. [2019]. Analogous to this, other systems use ML to solve component-scoped challenges by creating so called learned components. Examples include learned index structures Kraska et al. [2017], cardinality estimation Kipf et al. [2018], and join ordering in the query optimizer Marcus and Papaemmanouil [2018].

Other research has resulted in tools external to the DBMS that are powered by ML models. The tools output system configurations, such as knob settings Van Aken et al. [2017], Zhang et al. [2019], Duan et al. [2009], or make recommendations, such as adding or dropping indexes Ding et al. [2019]. These can then be used by the DBA to tune their system.

**Purpose of paragraph: what is this training data and why does it matter**

All these approaches face the same problem: the ML models that power these tech-

# November 20, 2019

## DRAFT

niques require a lot of training samples. These models commonly use metrics across the DBMS as training data. These metrics include query arrival rates Ma et al. [2018], latch contention in the transaction manager [citation for new paper](#), or disk write performance in the log manager [citation for new paper](#). The more diverse training data the system generates, the more models we can train. Further, it is important to generate training data in different environments to prevent the models from overfitting to one specific configuration. In cloud environments, even identical instance types can have significant variations in performance Difallah et al. [2013]. Accounting for such differences during training data generation will make the models robust and resilient to dynamic settings.

[Purpose of paragraph: explain offline training data generation and its shortcomings](#)

One possibility is to generate training data for models in an offline setting (i.e. not during production execution). Some systems Oracle [a] capture production workloads to allow DBAs to try out different system configurations on a snapshot of the database. Other tuning systems Van Aken et al. [2017], Zhang et al. [2019], Duan et al. [2009] use training data from a simulated or prior execution to train ML models that recommend configurations in new deployments. The problem with these offline approaches is that the models are trained on past workloads and environments. They do not test the effects of system configurations on the most recent workload. As workloads evolve, the efficacy of these tool’s recommendations decreases.

[Purpose of paragraph: explain online training data generation and its shortcomings](#)

In contrast, an online approach collects training data that the production system generates during live execution. Systems such as IBM’s LEO Markl et al. [2003] and Automatic Indexing Ding et al. [2019] use metrics generated from live query execution as training data. These systems can also tune their models on the fly. This technique is straightforward to instrument, but will only collect data for the specific environment and configuration the production system is currently running. As a result, the models could overfit the production setting. On the other hand, the exploration and testing of configurations in a production environment has serious performance risks. Many customers expect these systems to result in zero regressions on their performance Ding et al. [2019]. This includes both the overhead incurred from running these systems, and the potential degradations of bad recommendations or configurations.



## 1.1 Motivation

The demand for self-driving DBMS has grown with the growing size of data sets and the increased desire to run complex analytics over that data Pavlo et al. [2017]. As databases grow in size, so does their complexity, and therefore the difficulty to tune the DBMS. DBAs spend approximately a quarter of their time on tuning tasks and account for nearly 50% of a DBMS’s operating cost Debnath et al. [2008]. Self-driving DBMS alleviate this problem by using ML models to tune themselves without human intervention. This would allow DBAs to spend their valuable time on other important tasks.

In NoisePage, the training data, or input, to the ML models is metrics outputted during the execution of the system. For example, the logging component may output as a metric how long it took to write some amount of data to disk. As we discussed prior, this metrics collection comes at an overhead. To measure the effects of metrics collection in a DBMS, we execute TPC-C on PostgreSQL with and without metrics collection, taking the average over five runs. In our benchmark, the number of worker threads equal to the number of warehouses. We run this benchmark on an AWS m5d.4xlarge instance running with an Intel Xeon Platinum 8175 with eight threads (16 hyper-threads). Our results show that on average, this metrics collection overhead was around 11%. Users that may not want to pay this overhead on the master node, but still want to take advantage of the self-driving component of the DBMS can use Replicated Training technique to generate the data needed for the ML models.

Replicated Training allows offloading metrics collection to replica nodes; however, we must make sure that the replica can keep its copy of the database in sync with the master. Many database users have strict service-level agreements (SLAs) with regards to replication delay to limit the amount of data loss in the event of a node failure in an asynchronous replication environment. Allowing metrics collection to run unchecked in the replica can result in degradations in transaction replaying performance, resulting in higher replication delays. To observe the impact of these degradations on replication delay, we run the TPC-C benchmark on NoisePage with metrics collection enabled and disabled. The database is replicated asynchronously across the network on two machines with Intel(R) Xeon(R) CPU E5-2420 CPUs with six threads (12 hyper-threads). We look at the running mean of the replication delay over the benchmark entire benchmark run. The results in figure 1.1 show that metrics collection results in a significant increase in replication delay. We note that the benchmark run when metrics collection is enabled lasts longer because it takes longer for the replica to become fully in-sync with the master. In our first data point, the replication delay with metrics collection enabled and disabled is 1.48ms and 1.33ms, respectively. This difference is an 11.6% overhead, supporting the behavior we saw in

# November 20, 2019

## DRAFT

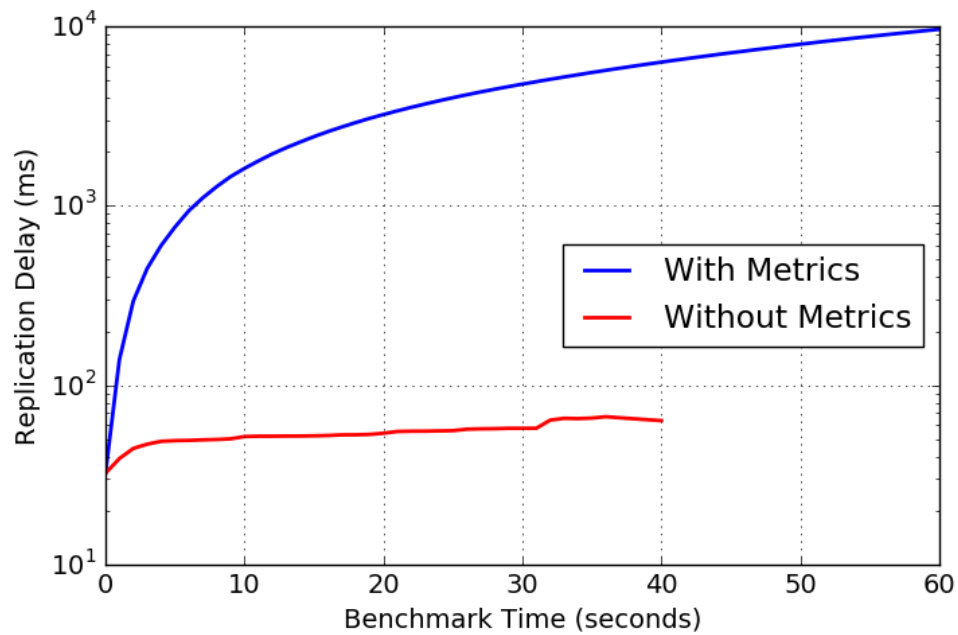


Figure 1.1: **Metrics overhead on replication delay in NoisePage** - Effects of replication delay when metrics collection is enabled over TPC-C with 4 warehouses and a 10k txns/sec arrival rate

# November 20, 2019

## DRAFT

PostgreSQL. For our last data point for each benchmark run, the replication delay with metrics enabled is  $152\times$  worse compared to when replication is disabled. When metrics collection is enabled without controls, the replica falls behind to the master. It is never able to catch up, resulting in this significant difference in the replication delay at the end of each benchmark run.

## 1.2 Contribution

Rewrite for thesis style

In this paper we present an online method for using replicas to generate training data for models used by self-driving DBMSs. Replicas hold the most up-to-date snapshot of the data, and experience the same write workload as the master node. By selecting a sample of the read workload from the master and executing it on the replica, the system can generate meaningful training data in an environment that is similar to the master node, for only a small cost of overhead on the replica. To account for the overhead required for approach, we allow the user to bound the percentage increase in replication delay incurred by the replica. If the standby falls to far behind, our system will halt to allow it to catch up to the master.

November 20, 2019  
DRAFT

## Chapter 2

## Background

### 2.1 Replication

Purpose of paragraph: explain why replication is done

Practically every production system will support some form of replication. Most systems will employ a hot-standby approach to replication, where a different database instance is running on a separate machine. This machine will copy changes from the master to create a consistent snapshot of the data, and will be ready to become the primary database in the case of a failure of the current master.

Purpose of paragraph: explain different ways replication is done

Production systems handle data replication between the master and its replicas in different ways. Most systems will do a form of log shipping, where the primary node will ship its change log to the replicas IBM, Microsoft [2017], MongoDB, PostgreSQL, Oracle [c], MySQL, which they will then replay. Other systems Google, Amazon, Snowflake abstract away how replication is done by letting an underlying cloud object storage handle the replication. To the system programmers, this appears as if all replicas are reading from the same logical copy of the data.

Regardless of how replication is done, all DBMSs that support replication will suffer from replication delay: the delay in time between when an update is done on the node that received the update request, and when a replica is able to see that update. Many DBMS users have strict SLAs that they must adhere to. A delimited replication delay is a common SLA users expect their DBMS to support in order to ensure a bounded degree of consistency between master and replica nodes. Low replication delays also reduces the

amount of data loss in the event of a crash.

Purpose of paragraph: explain tradeoffs of different log shipping techniques

If systems employ replication using log shipping, they will have to make the important decision of what type of logging to use: physical or logical logging. Physical logging involves the recording the physical changes that are made to the storage layer of the DBMS. Logical logging on the other hand logically describes the high level changes made to the data Yao et al. [2016]. Most commercial systems that employ logical logging will implement command logging, where the transactions themselves (or commands) are logged Malviya et al. [2014]. Each logging type has its own tradeoffs. Physical logging carries more overhead during execution because it produces more logs than logical logging, but is faster and more parallelizable during recovery. Logical logging on the other hand has less overhead during execution (typically a single log per transaction), but is more expensive to replay during recovery.

Purpose of paragraph: explain what our system did, maybe this does not belong here

Our system had the same design decision to make. Physical logs allow for faster replaying, allowing for a more up-to-date replica, but only exercise the lower levels of the system. Logical logs on the other hand are slower to replay, but will exercise more layers of the DBMS.

## 2.2 Training Data Collection

TODO: have citations for this section

Purpose of paragraph: introduction to area and challenges of training data collection

Effective training data generation is a hard and expensive task for production ML models. Companies have entire teams of Data Analysts and spend thousands in compute resources to generate enough data for their models to produce reasonable outputs. write some more, this paragraph doesn't feel right/finished

Purpose of paragraph: describe what aspects of training data collection we care about, and how we will describe those

There are two important facets of training data generation that we will concern ourselves with: performing actions to generate this training data, and choosing what data to generate. To illustrate the challenges in these two areas, we will discuss them in the context of self-driving cars. Self-driving cars serve as a good comparison to self-driving DBMSs as they are both expensive to deploy, and must be able to handle completely new

# November 20, 2019

## DRAFT

environments or workloads.

Purpose of paragraph: explain how other domains sample and label training data, and why its expensive

Performing actions to generate training data can be broken down into two parts: sampling and labeling. Self-driving cars sample new data by driving through streets and recording its environment through use of high-tech cameras. Apart from the difficulty of building all the technology to capture the car's environment, this process is extremely expensive. In the case of supervised learning models (should i explain supervised vs unsupervised), the sampled data must then be labeled. Often times this requires a large amount of humans manually labeling objects in the car's film. For other domains, such as ML for medicine, labeling can be incredibly expensive as often times few people are qualified enough to accurately produce labels.

Purpose of paragraph: explain how its not pheasable to sample everything

The other aspect one should consider is choosing what to sample. As we have discussed, the process of sampling and labeling can be extremely resource intensive. On one hand, the training data should be diverse in order to allow the ML models to generalize. Capturing every single environment and data point, however, is often not pheasable. For example, we can choose to have our self-driving car prioritized learning how to driving down streets we've never seen before. It's not realistic, however, to drive down every possible street during every different environment, such as daytime, nighttime, rain, rush hour, etc.

Purpose of paragraph: explain how expensive/impossible it is to build an offline simulator

A different approach is to build an environment where the workload can be simulated to generate training data. For example, one could construct a test track where a self-driving car could drive around and learn. There are some obvious limitations to this approach. The first is cost, it may be impossible to build such an environment depending on the resources required. Many production database systems are simply too large and expensive to duplicate in a simulated environment. Secondly, these simulated environments approximate a real environment. A car test track would not perfectly replicate the chaotic and variable reality of driving. Similarly, a simulated environment would not be able to perfectly represent the variabilities of a mission critical production system, such as workload spikes or machine failures. Should I discuss gyms? or is that another thing that isn't relevant here

November 20, 2019  
DRAFT



## Chapter 3

### System Architecture

NoisePage is an in-memory database management system (DBMS) that supports ACID transactions and `SNAPSHOT-ISOLATION`. The system is built in C++, and supports the PostgreSQL Wire Protocol for communicating with the database server. NoisePage was originally built to be a single node system, and we will discuss in Section 3.6 the addition of hot-standby replication. Transactions can be executed synchronously or asynchronously, and crash recovery is possible. We now discuss the architectural components of NoisePage relevant to supporting replicated training.

#### 3.1 Transactions

Transactions are the core atomic unit of work in a DBMS. The transactional approach taken by NoisePage is extremely conducive to a streamlined logging, recovery, and physical replication scheme. NoisePage’s transactional engine is a multi-versioned delta store Reed [1978] that supports `SNAPSHOT-ISOLATION` Berenson et al. [1995]. The transactional engine is coordinated by the `TransactionManager` component. In our transaction engine, readers do not block writers and vice versa, however, write-write conflicts on a per tuple basis are not allowed. When logging (Section 3.3) is enabled, transactions are considered active between when they begin, and when they commit or abort. When logging is disabled, transactions are active between when they begin, and when their changes have been serialized by the log manager (described in Section 3.3). This can be further extended by enabling synchronized commit, which guarantees that a transaction is active until its changes have been persisted in disk.

Tuples are uniquely identified by a *TupleSlot* object that stores offset information about

# November 20, 2019

## DRAFT

the tuple. The `TupleSlot` is created when a tuple is first inserted into a table. Each table has an additional, invisible column reserved for storing the head of this version chain.

Transactions update the database by adding their changes to a version chain for a specific tuple. To read a tuple, a transaction traverses the version chain until it sees the first change visible to that transaction. Changes are stored newest-to-oldest in the version chain to facilitate fast reads.

These changes, or delta records, are not an updated copy of the tuple, but rather the after-image changes made by the transaction on the tuple. There are four types of records: Redo, Delete, Commit, and Abort. Aside from the changes or transactional action, the delta records hold additional information such as a transaction *start* timestamp and what database and table were modified. These pieces of information are necessary for ensuring correct replaying of the record during recovery or replication. It is important to note, however, that delta records do not hold all the information needed to accomplish this. It is the role of the log manager, discussed in section 3.3, to serialize any additional information needed along with the record.

Each transaction has its own redo buffer to store these records. Rather than using an extendible buffer, redo buffers are fixed size (4096 bytes). Upon creation, transactions are given a buffer from a pre-allocated, centralized buffer pool. To record a change, a transaction reserves space in its redo buffer, and writes in the new change. In the case that there is not enough space left in the buffer, the transaction will hand off the buffer to the log manager, and receive a new one from the buffer pool. This allows downstream consumers, such as logging and replication to process these changes before a transaction has completed.

Committing a transaction is straightforward: the transaction will write in a commit record to the redo buffer, and hand it off to the log manager. During commit time, the oldest active transaction timestamp is also polled from the timestamp manager, and included in the commit record. Aborts, on the other hand, require additional logic to ensure correct behavior during recovery and replication. We discussed before that a transaction hands off its buffer to the log manager once it is full. Due to this, it is possible for an aborting transaction to have already persisted records. In order for correct behavior during recovery and replication, an aborting transaction that has previously flushed records must also flush an abort record. Without this abort record, a downstream consumer would be unable to differentiate such changes from a long running transaction. From observation, however, it is rare that an aborting transaction will ever flush a buffer, as the amount of data generated by an OLTP transaction rarely fills up an entire redo buffer.

### 3.2 Timestamp Manager

The timestamp manager is a central component that is in charge of providing atomic timestamps to running transactions. Timestamps are globally unique 64-bit unsigned integers. A transaction is given a *start* timestamp when it begins, and a *commit* timestamp when it successfully commits.

The timestamp manager also maintains an *active transaction set* that contains the *start* timestamp of every active transaction currently running in the system. Using this, the timestamp manager can be polled for the oldest active transaction timestamp (smallest *start* timestamp in the set). If there is no oldest active transaction, we indicate so using a special value. The importance of this information is discussed in section 3.4. A transaction is not removed from this set until its contents have been serialized by the log manager (discussed in section 3.3). This prevents background garbage collection (GC) from cleaning up the transaction before its changes have been persisted.

Fast performance of polling for the oldest active transaction timestamp is a crucial because it must be done during the critical section of every committing transaction. When logging is disabled, the number of active transactions is bounded by the number of worker threads available to the system. When logging is disabled, however, there is no bound on the number of active transactions, since transactions are active until they are at least serialized. Due to this, polling for the oldest active transaction, which requires scanning the entire active transaction set, can be a costly operation. We instead keep a cached oldest active transaction timestamp in the timestamp manager. During commit, transactions atomically read this value instead of scanning the active transaction set. The cached value is periodically refreshed by background Garbage Collection (default to every 5 ms). While the cached value may be a stale view of the system for a committing transaction, it still maintains correctness, as the transaction associated with the cached value is guaranteed to have been active at some point and older relative to any transaction which reads the cached value.

### 3.3 Logging

Changes to NoisePage are persisted on disk using write-ahead logging Mohan et al. [1992] through a dedicated component called the log manager. The log manager is in charge of serializing delta records such that they are entirely replayable on their own without any additional in-memory metadata, such as in the case of recovery after system crash or replication. To do so, the log manager coordinates multiple parallel tasks structured in a

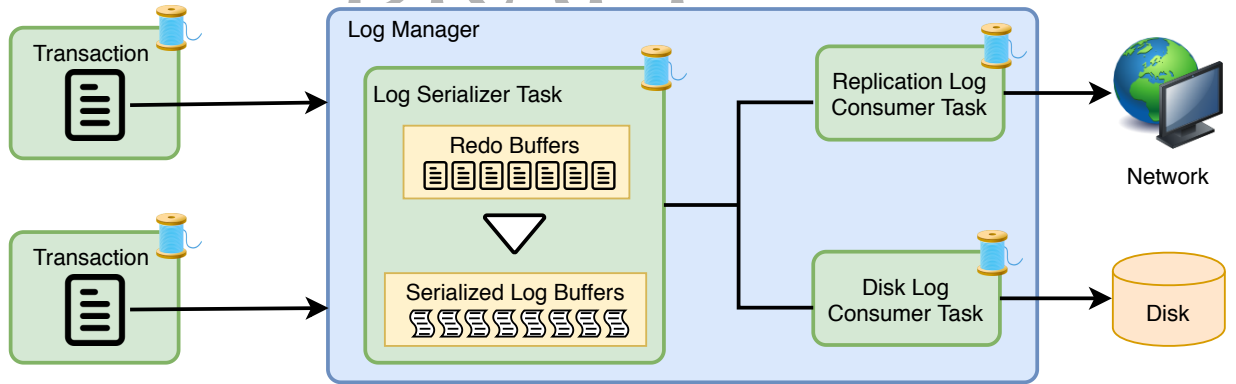


Figure 3.1: **Log Manager Architecture** - The log manager receives as input log record buffers from transactions, serializes these records, and sends them to different destinations (e.g. disk).

Producer-Consumer architecture shown in Figure 3.1. The producer, a log serializer task, feeds serialized logs to multiple log consumer tasks (e.g. disk log consumer task).

### 3.3.1 Log Serializer Task

The log serializer tasks receives redo buffers from transactions, and serializes their contents into fixed-sized buffers (4096 bytes) to be processed by the log consumers. The serializer task ensures that all the information needed to replay the delta record is included alongside it. For example, large varlens in NoisePage are not inlined in the delta record, but rather stored in a separate memory location, with just a pointer to the varlen entry stored in the delta record. The serializer task must thus fetch the varlen entry and serialize it inline. Serialization also removes any padding used in the in-memory log record.

The serialized format of the log records is shown in figure 3.2. The delete, commit, and abort records have fixed serialized sizes of 29, 29, and 13 bytes respectively. Due to the variability in data being updated, the size of serialized `RedoRecord`s varies depending on the delta record's contents. Over a run of the TPC-C benchmark with 4 warehouses, the average `RedoRecord` size is 100 bytes. For reference, the same execution of TPC-C will generate approximately 1 GB of total log data.

The serialized ordering of the records is important to ensure correct replayability. Recall from Section 3.1 that transactions can hand off record buffers as soon as they fill them. This means that records of different transactions can be interleaved in the log. Additionally, transactions can appear in non-serial order. The only guarantee that is made is that

# November 20, 2019

## DRAFT

record len (uint32)	record type (uint8)	txn start timestamp (uint64)	
database ID (uint32)		table ID (uint32)	
TupleSlot (uint64)		num columns (uint16)	
col ID 1 (uint32)	col ID 2 (uint32)	col ID 3 (uint32)	...
col 1 attr size (uint8)	col 2 attr size (uint8)	col 3 attr size (uint8)	...
null bitmap (variable)		val 1	val 2
val 3 varlen size (uint32)	val 3 varlen content		...

(a) Redo Record

record len (uint32)	record type (uint8)	txn start timestamp (uint64)
database ID (uint32)	table ID (uint32)	TupleSlot (uint64)

(b) Delete Record

record len (uint32)	record type (uint8)	txn start timestamp (uint64)
txn commit timestamp (uint64)		oldest active txn timestamp (uint64)

(c) Commit Record

record len (uint32)	record type (uint8)	txn start timestamp (uint64)
---------------------	---------------------	------------------------------

(d) Abort Record

Figure 3.2: **Log record serialization formats** - Along with the log record, additional information must be serialized to ensure replayability.

## DRAFT

records for an individual transactions appear in order that they were created, with a commit or abort record always being the last record to appear. As we will discuss in Section 3.4, this guarantee, along with the oldest active transaction timestamp discussed in Section 3.1, is enough to achieve a consistent snapshot of the database after replaying the log.

### 3.3.2 Log Consumer Tasks

The buffers generated by the log serializer are handed off to an arbitrary number of log consumers. Each consumer is given a copy of the serialized buffer so they can each work independently of each other, preventing a slow consumer from slowing down the others. Currently, NoisePage supports two consumers: (1) a disk consumer that writes logs to a file on disk, and (2) a replication consumer that sends logs over network to replicas.

The disk consumer task waits until it receives buffers from the serializer task, and writes them to a log file specified by the user. Writing to the log file is fast because we are always doing sequential writes. For good performance in persisting the log file to the disc, we take advantage of batch commit, which is configurable by the user with a combination of time and data size settings. For example, under the default settings, the disk consumer task will persist the log file every 10 milliseconds or if more than one megabyte of data has been written since the last persist. The configurability of this process will allow the self-driving infrastructure to find the optimal combination of settings for a given workload.

The replication consumer task also receives buffers from the serializer task and sends them over the network to any replicas listening to the master. In order to minimize the replication delay, there is no batch commit done. Instead, serialized logs are sent as soon as they are handed off to the replication consumer task. The network protocol used for sending logs between nodes is described in detail in Section 3.5.

## 3.4 Recovery

Recovery in NoisePage is performed by a component called the recovery manager. The recovery manager receives serialized log records from an arbitrary source, and replays them to produce a consistent view of the database.

A `LogProvider` class will deserialize data into `RedoRecords`, and hand them off to the recovery manager. The `LogProvider` class provides an abstraction to the recovery manager as to what the source of the records are. This way, log records can come from any source, such as a log file or over network, without any changes needed to the log

replaying logic of the recovery manager.

Abstracting away the source of log records gives us the great advantage that we can implement the replaying component of replication for "free". By simply having the source of records be a stream of logs over network, a standby replica can instantiate a recovery manager to replay the log records arriving from the master node. This approach for replication is also done by other systems like PostgreSQL. This, however, requires the processing model of the recovery manager to be a streaming model. We can not take advantage of cases when all the log data is available apriori, as is the case during single node crash recovery. Other recovery algorithms, such as ARIES Mohan et al. [1992], take advantage of having access to all the data from the start and trim out unnecessary processing.

### 3.4.1 Log Record Replaying

The API for updating tables in the system (`SqlTable`) allows for easy replaying of records. Recall from Section 3.1 that transactions must write their changes as delta records in their private buffers. The system takes advantage of this by passing `SqlTables` a pointer to these records in the buffer. This prevents having to make an additional copy of the data. Since recovery deserializes `RedoRecords`, there is no need to transform the data, the `SqlTable` API simply accepts these `RedoRecords` directly.

The `TupleSlot` contained in the replayed record is no longer valid during recovery, as it represented a unique memory location prior to recovery. Instead, we use it as an internal mapping from the original `TupleSlot` to the new one when the insert is replayed. Using the mapping, we can correctly identify what `TupleSlot` updates should be applied to after recovery.

Processing records that modify the catalog require additional logic. The metadata stored in the catalog is kept in tables. These tables are the same structure used for user tables in the system. This makes recovering the catalog metadata largely easy, as they simply appear as updates to the catalog tables. Despite this, additional logic is needed to reinstantiate certain in-memory objects in the system, such as indexes, views, or user tables.

While its is possible that we could replay changes as we see them in the log and roll them back in the case of aborts, we will see in Section 3.4.2 that we must defer all updates until we see a commit or abort record anyway. Buffering also gives us the added advantage that we prevent the unnecessary of replaying records for aborted transactions. When we see an abort record, we clean up and discard any records buferred for that transaction. When we see a commit record, we process the transactions as described in Section 3.4.2.

DRAFT

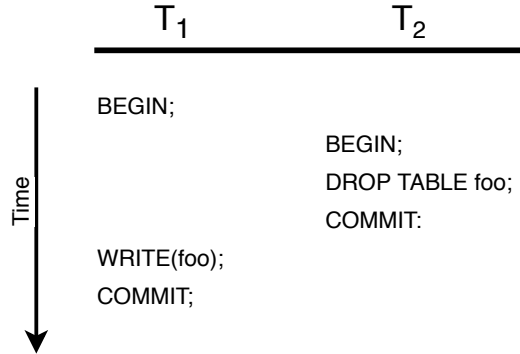


Figure 3.3: **Schedule involving DDL changes** - This schedule is allowed under `SNAPSHOT-ISOLATION`, but can create problematic races involving the `DROP TABLE` command.

The transaction's changes are always applied in the order they were made prior to recovery.

### 3.4.2 Transaction Replaying

Special considerations must be made when replaying transactions because of the streaming processing model of recovery and the design of our transactional engine. Recall from Section 3.3 that the only guarantee we have about the ordering of logs is that changes for an individual transaction are in the order they occurred. There are no guarantees about how transactions are ordered relative to each other in the log, so it is the responsibility of the recovery manager to execute them in an ordering that results in a consistent view of the database.

Recall from Section 3.1 that NoisePage supports `SNAPSHOT-ISOLATION`, which means that each transaction operates on a "snapshot" of the database taken when the transaction begins. Additionally, any committed transactions which executed concurrently are guaranteed to not have any write-write conflicts with each other on a per tuple basis. Thus, we are guaranteed that there are no dependencies between transactions that executed concurrently, and all committed transactions that are replayed will successfully commit. Further, because we use physical logging, all the values being written during log replaying are predetermined - i.e. no writes are based on reads or randomization. Based on these two guarantees, we are able to replay transactions sequentially (i.e. a transaction commits before the next one is allowed to begin).

We have determined that we can execute transactions sequentially, but the order in



## DRAFT

which we execute them is also important. Consider the schedule in figure 3.3 which is allowed under `SNAPSHOT-ISOLATION`. While there is no write-write conflict in this schedule, there is an implicit conflict due to the DDL change (`DROP TABLE`). During replaying, transaction  $T_1$  must be replayed before transaction  $T_2$  in order to ensure that  $T_1$ 's write to `foo` occurs before  $T_2$  deletes `foo`. There are no guarantees about ordering of logs between transactions, so its possible for  $T_2$ 's changes to appear before  $T_1$ 's changes in the log. Even worse, if  $T_1$  is a long running transaction, it's changes may not appear until much further along in the log. This raises the issue of when is it safe to execute a transaction.

Executing transactions in the order in that they appear in the log could violate `SNAPSHOT-ISOLATION`, since executing a newer transaction first would create a different snapshot than what an older transaction saw when it was executed before recovery. Thus, we must execute the transactions in the order that they were created i.e. ordered by their *start* timestamp.

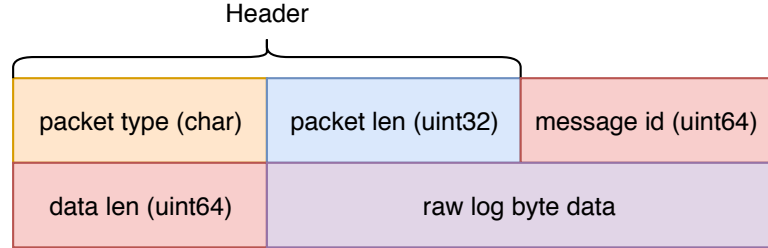
A simple aproach for accomplishing this would be as follows: A transaction  $T_i$  with *start* timestamp  $i$  is safe to replay after we have replayed transaction  $T_{i-1}$ . The transactional engine, however, does not guarantee that consecutive transactions have consecutive *start* timestamps. Additional processes, such as garabage collection or assigning commit timestamps, also receive timestamps from the timestamp manager.

The solution to this problem is by using the oldest active transaction timestamp described in Section 3.1. When a transaction commits, the oldest active transaction timestamp at the time of commit is included in the commit record (shown in figure 3.2(c)). When a transaction is entirely deserialized, rather than executing it right away, we defer its execution. Using the oldest active transaction timestamp  $i$  stored in the commit record, we then execute, in sorted order oldest-to-newest, all deferred transactions with *start* timestamps  $j$  where  $j \leq i$ . If  $i$  is the special value reserved for indicating there are no active transactions, then we execute all deferred transactions.

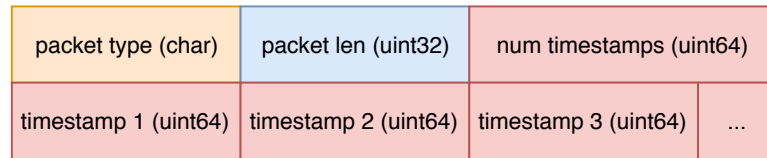
Once again, consider the schedule in figure 3.3. Suppose the *start* timestamps of  $T_1$  and  $T_2$  are 1 and 2 respectively. The commit record of  $T_2$  will indicate that the oldest active transaction timestamp at commit time is 1. If  $T_2$  is serialized **before**  $T_1$ , it will be deferred because  $1 < 2$ . Eventually  $T_1$  is deserialized and executed, followed immedietly by  $T_2$ , because there were no older transactions at the time  $T_1$  committed.

# November 20, 2019

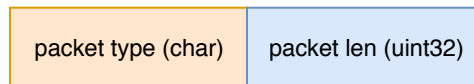
## DRAFT



(a) ReplicationDataPacket



(b) CommittedTransactionsPacket



(c) EndReplicationPacket and  
ReplicaSyncedPacket

Figure 3.4: **Internal NoisePage Protocol packet types** - Packets are minimal to reduce network congestion and speed up packet processing

## 3.5 Internal NoisePage Protocol

Internal NoisePage Protocol is the network protocol used to communicate between nodes of NoisePage. Internal NoisePage Protocol uses TCP/IP sockets to send data between the master and replica node(s). TCP is used over UDP for its message delivery guarantee; it's important that no log data is lost over network. NoisePage has a network layer which sits at the top of the system to handle client and other NoisePage node connections. The network layer is designed to support multiple protocols on separate ports. Currently it supports Internal NoisePage Protocol and the PostgreSQL Wire Protocol (described in Li [2019]).

Messages are packets consisting of a header and payload. The header is used by NoisePage's network layer to parse the packet. It consists of single `char` to identify the packet type, and a `uint32_t` value for the entire size of the packet. This portion of

## DRAFT

the protocol resembles the PostgreSQL Wire Protocol. To support a new protocol, users write a `PacketWriter` class, which contains logic to write the payload of the packets. Finally, users write a `ProtocolInterpreter` class which parses and processes the payload based on the packet type.

The current packet types for Internal NoisePage Protocol are shown in figure 3.4. The `ReplicationDataPacket` (figure 3.4(a)) holds variable-length portions of the serialized log data from the master node's log manager to the replica. The `CommittedTransactionsPacket` (figure 3.4(b)) is sent from the replica to the master to notify when the replica has successfully replayed and committed transactions. This provides support for synchronous replication. The *start* timestamp of the replayed transactions is stored in the packet's payload. The `EndReplicationPacket` is sent by the master and tells the replica to end replication. The `ReplicaSyncedPacket` is sent from the replica and notifies the master that the data in both nodes is in sync. Both these packets (figure 3.4(c)) require no payload as the type in the header entirely describes the purpose of the packet.

### 3.6 Replication

An overview of the replication architecture in NoisePage is shown in figure 3.5. Replication can be done by two NoisePage instances, a master and a replica, running on different machines connected to the internet. Replication is established by a connection between the master's log manager and the replica's network layer. When the logs are serialized by the log manager, they are placed into a `ReplicationDataPacket` and shipped over the network by an Internal NoisePage Protocol specific `PacketWriter` described in Section 3.5. The packets reach the replica's network layer and are handed off to the recovery manager running in the system. We implement a `LogProvider` class called the `ReplicationLogProvider` to parse the log data from these arriving packets into log records. For replaying these logs, we discussed in Section 3.4 that we can accomplish replication using the same recovery manager logic used for crash recovery.

If synchronous replication is enabled, the replica will send the master a `CommittedTransactionsPacket` when a transaction is replayed. If the replica is in sync with the master (i.e. there are no more logs left to replay), it will also send a `ReplicaSyncedPacket`.

The recovery manager running on the replica will sit in a loop, continually processing log records as they arrive over the network. This is unlike in crash recovery where the recovery manager will terminate when it reaches the end of the log. Instead, the master can terminate replication with a replica by sending it an `EndReplicationPacket`.

# November 20, 2019

## DRAFT

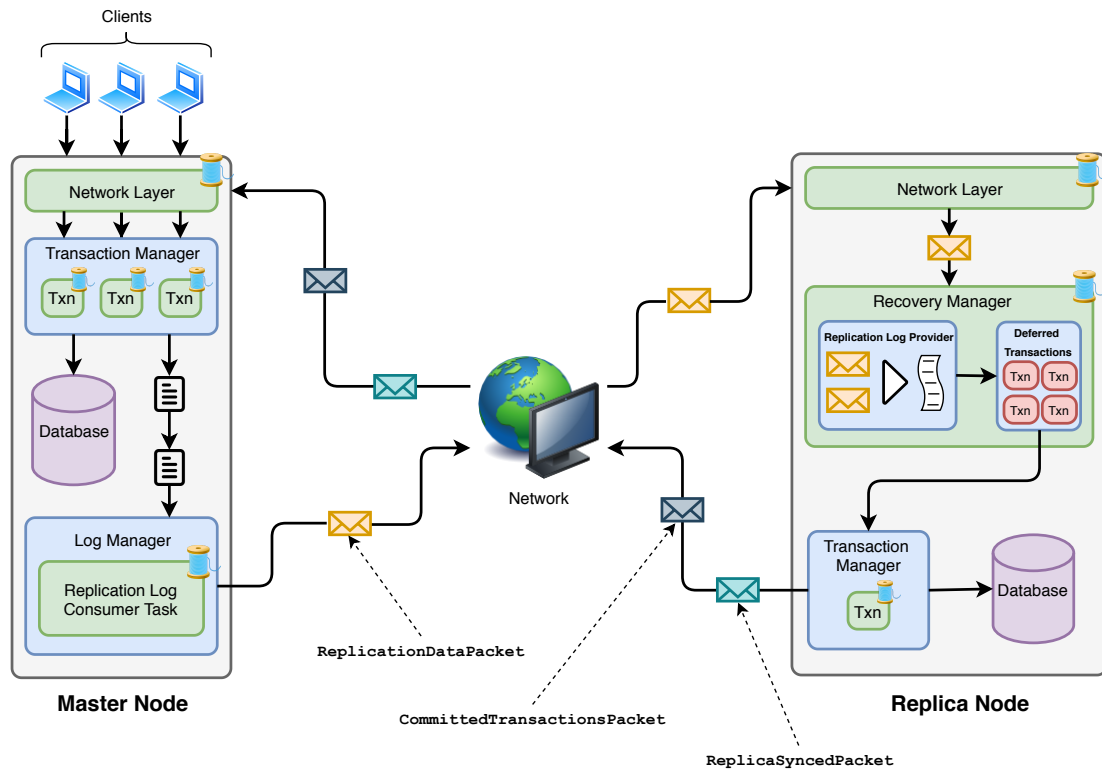


Figure 3.5: **Replication Architecture** - While there are other components in NoisePage, this highlights the processes involved in replicating data between a master and replica

## **Chapter 4**

### **Replicated Training**

Discuss different dynamic metrics collection control policies

Discuss how we can turn off areas of metrics collection, rather than just the entire thing. We could even take an active learning approach for choosing what areas to turn off.

November 20, 2019  
DRAFT

## Chapter 5

### Evaluation

We now evaluate and analyze our replication architecture and Replicated Training technique. We build all the infrastructure within NoisePage. We use the following two types of machines for our experimental evaluation:

- Type 1: Dual-socket 10-core Intel Xeon E5-2630v4 CPU, 128 GB of DRAM, and a 500 GB Samsung 970 EVO Plus SSD. This machine is used for single node microbenchmarking
- Type 2: Single-socket 6-core Intel Xeon CPU E5-2420 CPU and 32GB of DRAM. These machines are used for replication between two nodes experiments, and are wired together using [add info about network connectivity here](#)

We first give an analysis of metrics collection in NoisePage. We next evaluate the OLTP performance of our replication architecture described in Chapter 3. We then observe the behavior of dynamically controlling metrics exporting. Finally, we analyze the effectiveness of our Replicated Training technique to build accurate ML models.

**REMOVE LATER** benchmarks for metrics exporting:

- How much data is outputted with the following heuristics
  - No limiting
  - Hardcoded replication delay
  - Delay exceeds some standard deviation of running mean

benchmarks for ML model

- Use the previous heuristics and evaluate how they affect model performance
  - Make sure to get baseline value

## 5.1 Replication Architecture

To evaluate our system architecture, we want an OLTP workload. Even though reads are not replicated, a write-only workload is not representative of a real OLTP workload. Towards this goal, we use the TPC-C benchmark, which simulates a delivery system, with orders being placed and received.

To simulate a distributed environment, we execute our benchmarks between two NoisePage instances running on Type 2 machines. We use a simple master-replica architecture where one instance is the master and serves requests, and the other machine is a replica node. We replicate data asynchronously across the two machines. For synchronizing clocks in our machines to get measurements for replication delay, we use the Network Time Protocol (NTP). We do not require extreme clock precision because we are not doing any logic based on the times, we are simply trying to estimate delay.

We now evaluate our replication architecture using microbenchmarks, and define some important test configurations and baselines used moving forward.

### 5.1.1 Arrival Rate

One important consideration to make with any benchmark is the arrival rate. The arrival rate is defined as the frequency the database is queried across all worker threads. For example, if there are four worker threads, each executing 2,500 transactions per second (txns/sec), then the arrival rate is 10,000 txns/sec.

It is important to pick a good arrival rate for measuring the replication delay. If the arrival rate exceeds the rate at which replication is able to replay transactions, then the replication delay will grow unboundedly because the replica is not able to keep up with the master node. Figure 5.1 showcases this effect in NoisePage. We can see how the replica remains in sync with the master with a reasonable delay until the arrival rate reaches 14,000 txns/sec. After that, the replica is not able to keep up with the arrival rate of transactions, and accordingly the delay sharply increases. This is a natural limitation in any system, although they may vary in the arrival rate they are able to tolerate. For measurements,



# November 20, 2019

## DRAFT

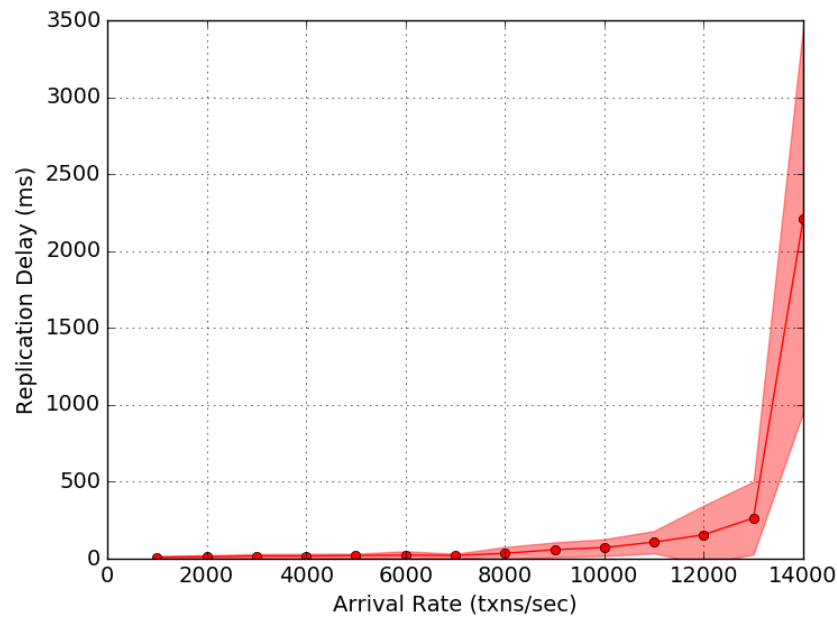


Figure 5.1: **Sensitivity of replication delay** - Measuring the average replication delay with varying arrival rates in NoisePage over TPC-C with 4 warehouses on Type 2 machines. When the arrival rate exceeds the transaction replaying rate, delay sharply increases. The shaded region denotes one standard deviation from the mean for each data point.

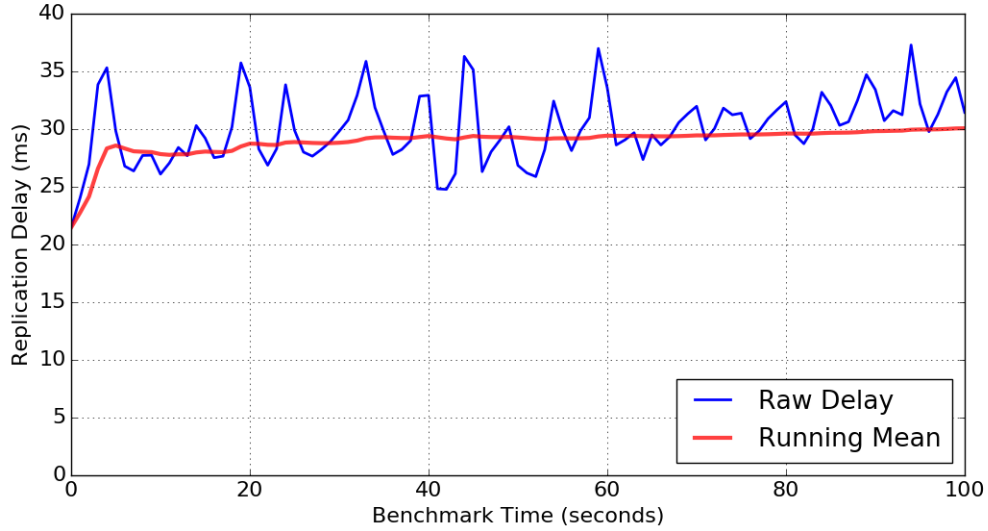


Figure 5.2: **Replication Delay in NoisePage** - Measuring the average replication delay in NoisePage over TPC-C with 4 warehouses on Type 2 machines. We average the delay for each second, and plot the average over 10 benchmark runs. The red line shows the running mean of the replication delay.

we assume an arrival rate of 10,000 txns/sec for future experiments to get stable delay readings.

### 5.1.2 Replication delay over time

As we discussed in section 2.1, many DBMS users have replication delay SLAs that they expect the DBMS to support. To get an idea of replication delay in NoisePage, we execute TPC-C using asynchronous replication, measuring the replication delay over the span of the benchmark in figure 5.2. The sharp spikes in replication are as a result of the delivery transaction of TPC-C that takes longer to replay relative to the other transactions. Despite the spikes, we see from the running mean (red line) that the replication delay remains fairly stable throughout the benchmark execution. Over the entire benchmark, the average replication delay is approximately 30ms.

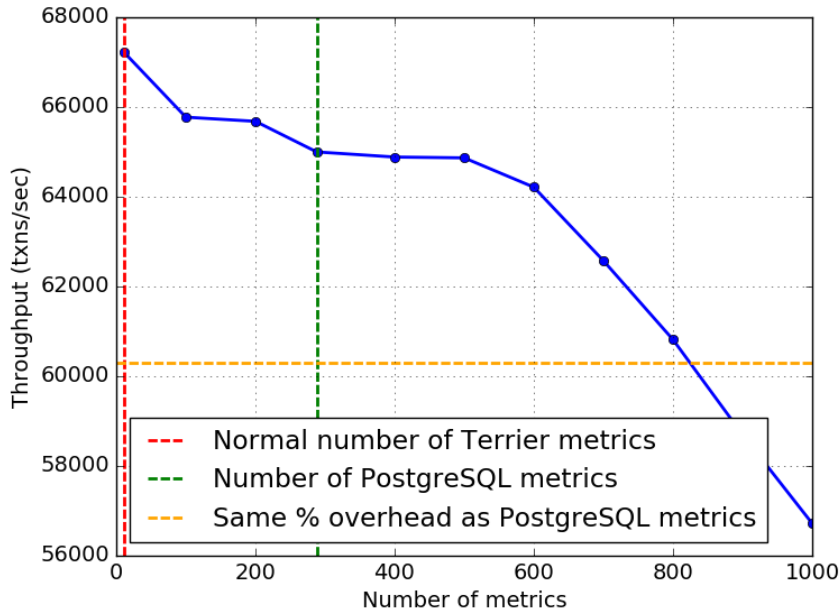


Figure 5.3: **Metrics overhead in NoisePage** - Overhead of metrics collection as we scale the number of metrics exported

## 5.2 Metrics Collection

In section 1.1, we motivated our decision to use database replicas by discussing the metrics collection overhead, and observed it is approximately 11% in PostgreSQL. For our this analysis, we used PostgreSQL instead of NoisePage because NoisePage is still a relatively new, unfinished system. In particular, NoisePage still has an immature metrics collection and does not yet output as many metrics as a system of its size should. Due to this, we decide simulate NoisePage having similar metrics overheads to PostgreSQL.

NoisePage currently has two (out of 10) high level components, the `TransactionManager` and log manager that export metrics. Collectively, these two components out 11 unique metrics (e.g. disk write speed, transaction latch wait time). For reference, we calculated PostgreSQL to output approximately 300 unique metrics by looking at its internal metadata tables.

The approach we took to simulate a realistic DBMS metrics overhead was to scale up the amount of metrics data exported by each component i.e when a metric is generated,

# November 20, 2019

## DRAFT

it is exported multiple times. To show the effect of this approach, we execute the TPC-C benchmark with 6 warehouses on machine Type 1. We compare the transaction throughput while scaling up the number of metrics exported using the approach described. Figure 5.3 shows the effects of this technique on transactional throughput. With the current metrics in the system (red line), NoisePage executes  $\sim 67,000$  txns/sec. If we scale number of metrics to the number of metrics estimate PostgreSQL exports (green line), we see a throughput of  $\sim 65,000$  txns/sec, only a 3% overhead. This does not equate the 11% we expect to see because throughout the lifecycle of a transaction, different metrics are exported at different frequencies. Therefore, scaling NoisePage's metrics to the same number of metrics that PostgreSQL exports is an insufficient comparison, as NoisePage exports at different frequencies than PostgreSQL. Instead, we scale up the number of metrics until we see the 11% overhead (yellow line), which occurs at approximately 800 metrics. We use this scale of metrics collection for future experiments.

November 20, 2019  
DRAFT

## **Chapter 6**

### **Conclusions**

Our goal is to parallelize individual transactions in a database, using FOO.

November 20, 2019  
DRAFT

## Bibliography

Amazon. Amazon aurora documentation: Replication with amazon aurora. 2.1

Hal Berenson, Phil Bernstein, Jim Gray, Jim Melton, Elizabeth O’Neil, and Patrick O’Neil. A critique of ansi sql isolation levels. In *Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data*, SIGMOD ’95, pages 1–10, New York, NY, USA, 1995. ACM. ISBN 0-89791-731-6. doi: 10.1145/223784.223785. URL <http://doi.acm.org/10.1145/223784.223785>. 3.1

B.K. Debnath, D.J. Lilja, and M.F. Mokbel. SARD: A statistical approach for ranking database tuning parameters. In *ICDEW*, pages 11–18, 2008. 1.1

Djellel Eddine Difallah, Andrew Pavlo, Carlo Curino, and Philippe Cudre-Mauroux. Oltpbench: An extensible testbed for benchmarking relational databases. *Proc. VLDB Endow.*, 7(4):277–288, December 2013. ISSN 2150-8097. doi: 10.14778/2732240.2732246. URL <http://dx.doi.org/10.14778/2732240.2732246>. 1

Bailu Ding, Sudipto Das, Ryan Marcus, Wentao Wu, Surajit Chaudhuri, and Vivek R. Narasayya. Ai meets ai: Leveraging query executions to improve index recommendations. In *Proceedings of the 2019 International Conference on Management of Data*, SIGMOD ’19, pages 1241–1258, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5643-5. doi: 10.1145/3299869.3324957. URL <http://doi.acm.org/10.1145/3299869.3324957>. 1

Songyun Duan, Vamsidhar Thummala, and Shivnath Babu. Tuning database configuration parameters with ituned. *PVLDB*, 2:1246–1257, 2009. 1

Google. Cloud spanner documentation: Replication. 2.1

IBM. High availability through log shipping. 2.1

- Andreas Kipf, Thomas Kipf, Bernhard Radke, Viktor Leis, Peter A. Boncz, and Alfons Kemper. Learned cardinalities: Estimating correlated joins with deep learning. *ArXiv*, abs/1809.00677, 2018. 1
- Tim Kraska, Alex Beutel, Ed Huai hsin Chi, Jeffrey Dean, and Neoklis Polyzotis. The case for learned index structures. In *SIGMOD Conference*, 2017. 1
- Tim Kraska, Mohammad Alizadeh, Alex Beutel, Ed H. Chi, Jialin Ding, Ani Kristo, Guillaume Leclerc, Samuel Madden, Hongzi Mao, and Vikram Nathan. Sagedb: A learned database system. 2019. 1
- Tianyu Li. Supporting Hybrid Workloads for In-Memory Database Management Systems via a Universal Columnar Storage Format. Master’s thesis, May 2019. 3.5
- Lin Ma, Dana Van Aken, Ahmed Hefny, Gustavo Mezerhane, Andrew Pavlo, and Geoffrey J Gordon. Query-based workload forecasting for self-driving database management systems. In *Proceedings of the 2018 ACM International Conference on Management of Data*, SIGMOD ’18, 2018. 1
- Nirmesh Malviya, Ariel Weisberg, Samuel Madden, and Michael Stonebraker. Rethinking main memory oltp recovery. *2014 IEEE 30th International Conference on Data Engineering*, pages 604–615, 2014. 2.1
- Ryan Marcus and Olga Papaemmanouil. Deep reinforcement learning for join order enumeration. In *Proceedings of the First International Workshop on Exploiting Artificial Intelligence Techniques for Data Management*, aiDM’18, pages 3:1–3:4, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5851-4. doi: 10.1145/3211954.3211957. URL <http://doi.acm.org/10.1145/3211954.3211957>. 1
- V. Markl, G. M. Lohman, and V. Raman. Leo: An autonomic query optimizer for db2. *IBM Systems Journal*, 42(1):98–106, 2003. ISSN 0018-8670. doi: 10.1147/sj.421.0098. 1
- Microsoft. Sql server: Transactional replication. 2017. 2.1
- C. Mohan, Don Haderle, Bruce Lindsay, Hamid Pirahesh, and Peter Schwarz. Aries: A transaction recovery method supporting fine-granularity locking and partial rollbacks using write-ahead logging. *ACM Trans. Database Syst.*, 17(1):94–162, March 1992. ISSN 0362-5915. doi: 10.1145/128765.128770. URL <http://doi.acm.org/10.1145/128765.128770>. 3.3, 3.4
- MongoDB. Mongoddb documentation: Replication. 2.1



MySQL. Mysql documentation: Replication. 2.1

Oracle. Database real application testing user’s guide. a. 1

Oracle. Database: Oracle autonomous database. b. 1

Oracle. Timesten in-memory database replication guide: Overview of timesten replication. c. 2.1

Andrew Pavlo, Gustavo Angulo, Joy Arulraj, Haibin Lin, Jiexi Lin, Lin Ma, Prashanth Menon, Todd Mowry, Matthew Perron, Ian Quah, Siddharth Santurkar, Anthony Tomasic, Skye Toor, Dana Van Aken, Ziqi Wang, Yingjun Wu, Ran Xian, and Tieying Zhang. Self-driving database management systems. In *Conference on Innovative Data Systems Research*, 2017. 1, 1.1

PostgreSQL. High availability, load balancing, and replication. 2.1

D. P. Reed. Naming and synchronization in a decentralized computer system. Technical report, Cambridge, MA, USA, 1978. 3.1

Snowflake. Snowflake documentation: Key concepts & architecture. 2.1

Dana Van Aken, Andrew Pavlo, Geoffrey J. Gordon, and Bohan Zhang. Automatic database management system tuning through large-scale machine learning. In *Proceedings of the 2017 ACM International Conference on Management of Data*, SIGMOD ’17, pages 1009–1024, 2017. 1

Chang Yao, Divyakant Agrawal, Gang Chen, Beng Chin Ooi, and Sai Wu. Adaptive logging: Optimizing logging and recovery costs in distributed in-memory databases. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD ’16, pages 1119–1134, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-3531-7. doi: 10.1145/2882903.2915208. URL <http://doi.acm.org/10.1145/2882903.2915208>. 2.1

Ji Zhang, Yu Liu, Ke Zhou, Guoliang Li, Zhili Xiao, Bin Cheng, Jiashu Xing, Yangtao Wang, Tianheng Cheng, Li Liu, Minwei Ran, and Zekang Li. An end-to-end automatic cloud database tuning system using deep reinforcement learning. In *Proceedings of the 2019 International Conference on Management of Data*, SIGMOD ’19, pages 415–432, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5643-5. doi: 10.1145/3299869.3300085. URL <http://doi.acm.org/10.1145/3299869.3300085>. 1