



Self-Driving Database Management Systems

Database meets deep learning



16 September 2017
Sung-Soo Kim
sungsoo@etri.re.kr

Smart Data Research Group

ETRI



References & Slide Credits

2

- Andrew Pavlo and *et al*, *Self-Driving Database Management Systems*, CIDR '17.
- Joy Arulraj and *et al*, *Bridging the archipelago between row-stores and column-stores for hybrid workloads*, SIGMOD'16.
- Joy Arulraj and *et al*, *Write-Behind Logging*, VLDB'16.
- Sung-Soo Kim and *et al*, *Sweet KIWI: Statistics-Driven OLAP Acceleration using Query Column Sets*, EDBT'16.
- Sung-Soo Kim and *et al*, *Flying KIWI: Design of Approximate Query Processing Engine for Interactive Data Analytics at Scale*, BigDAS'15.

Outline

3

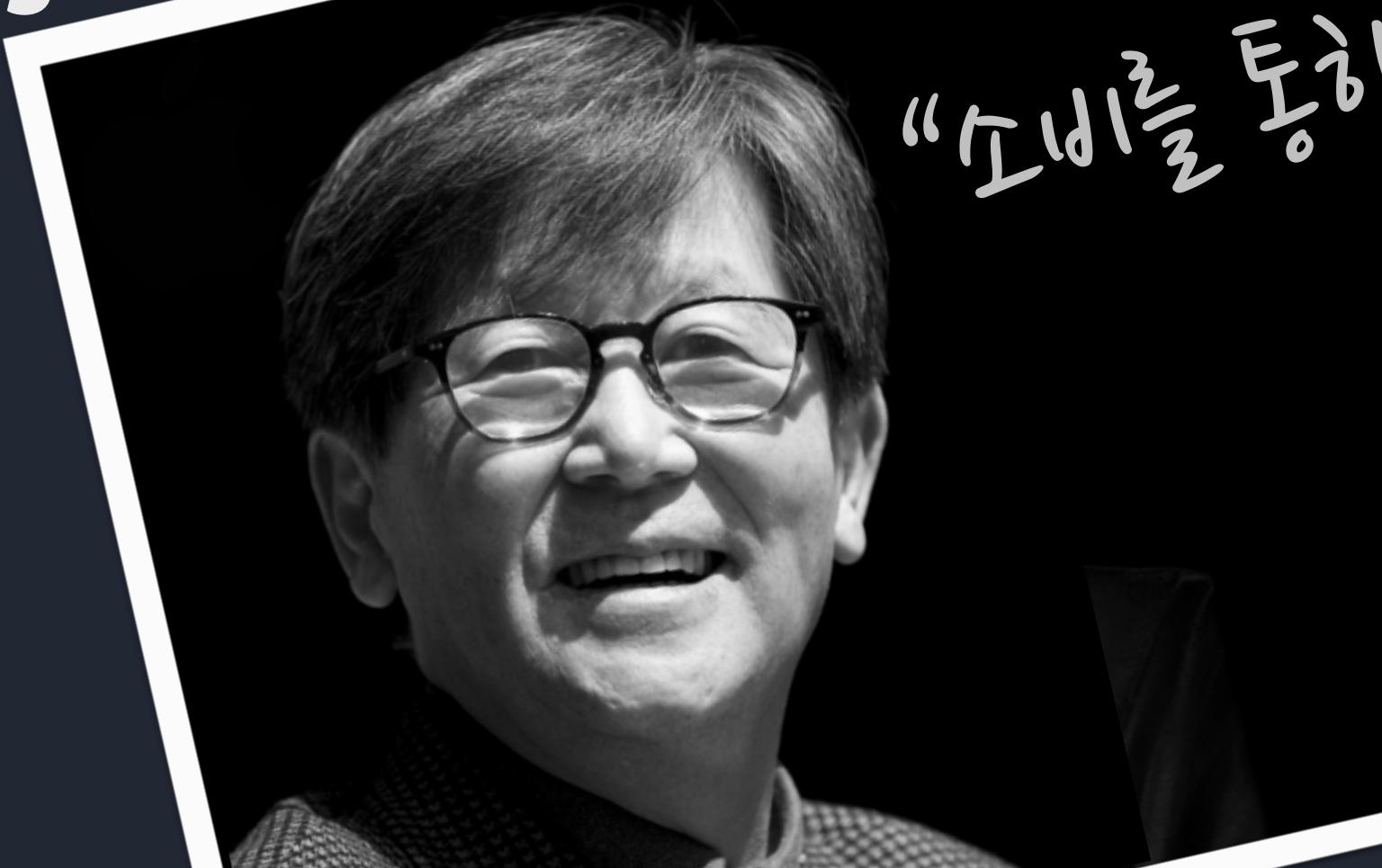
- Introduction
- Problem Overview
- Self-Driving Architecture
 - Workload Classification
 - Workload Forecasting
 - Action Planning & Execution
- Experimental Results
- Related Work
- Conclusion



Self-Driving Database Management Systems CIDR'17



<http://pelotondb.org>



“소비자를 통하여 차기 경제성을 만들어 낼 수는 없다.
인간의 경제성을 생산을 통해 형성된다.”

- 신영복 <담론>

Sweet KIWI: Statistics-Driven OLAP Acceleration using Query Column Sets

Sung-Soo Kim, Taewhi Lee, Moonyoung Chung and Jongho Won
Electronics and Telecommunications Research Institute (ETRI)
218 Gajeong-ro, Yuseong-gu, Daejeon
South Korea
{sungsoo, taewhi, mchung, jhwon}@etri.re.kr

ABSTRACT

KIWI is a SQL-on-Hadoop system enabling batch and interactive analytics for big data. In database systems, materialized views, stored pre-computed results for queries, are one of the most commonly used techniques to improve the query processing performance. However, the key challenge is using materialized views to maintain their freshness as base data changes. We propose a new approach for accelerating OLAP queries by using workload statistics and *query column* views. We present an architecture for generating query column sets of original tables. Experimental results demonstrate that our system improves performance by 1.77x on average in

Keyword

Column Sets, SQL-on-Hadoop

1. INTRODUCTION

Data warehouse (DW) on Hadoop is now being used intensively by users in enterprises as well as scientific (SoH) is a class of "Big Data" analytic applications. We introduce a column-oriented acceleration mechanism for deep analytics at scale.

SYSTEM OF FV

architecture for query processing and describes two main sections of the system.

an assumption about query workloads is similar to historical queries

This component is responsible for analyzing a set of historical query workloads to identify frequently used queries in the past. In order to construct the query column sets, we extract the metadata for query column sets over the entire original tables. The set of query column sets are updated both with the arrival of new data, and when the query workloads changes.

C_1, \dots, C_m , in horizontal sampling, let $S_h = \{R_i, R_{i+1}, \dots, R_{i+l}\}$, where $i \leq i + l \leq r$, denote a *row set* that consists of l rows in T . In vertical sampling, let $S_v = \{C_j, C_{j+1}, \dots, C_{j+k}\}$, where

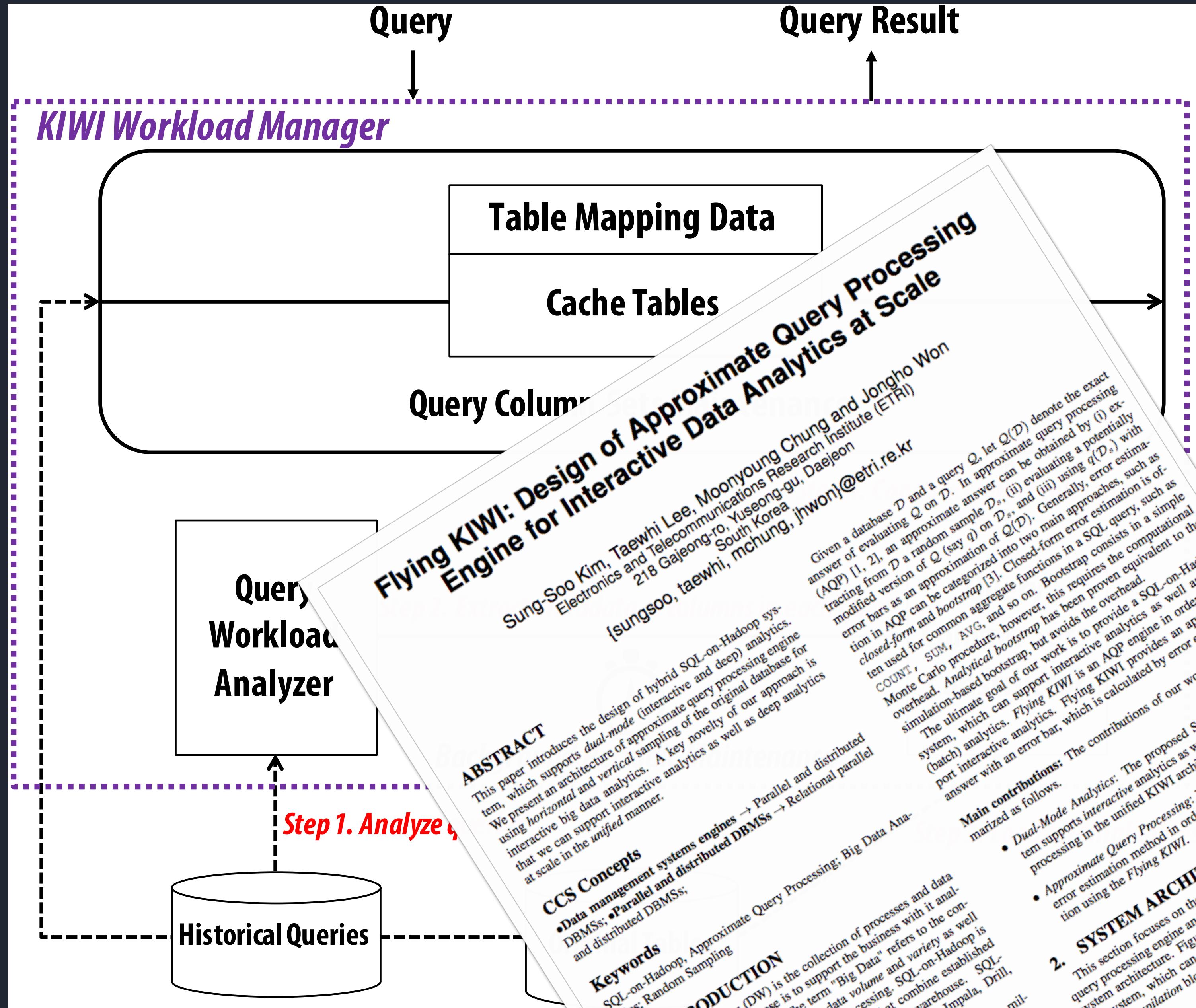
©2016, Copyright is with the authors. Published in Proc. 19th International Conference on Extending Database Technology (EDBT), March 15-18, 2016 - Bordeaux, France: ISBN 978-3-89318-070-7, on OpenProceedings.org. Distribution of this paper is permitted under the terms of the Creative Commons license CC-by-nc-nd 4.0

Series ISSN: 2367-2005

680

10.5441/002/edbt.2016.84

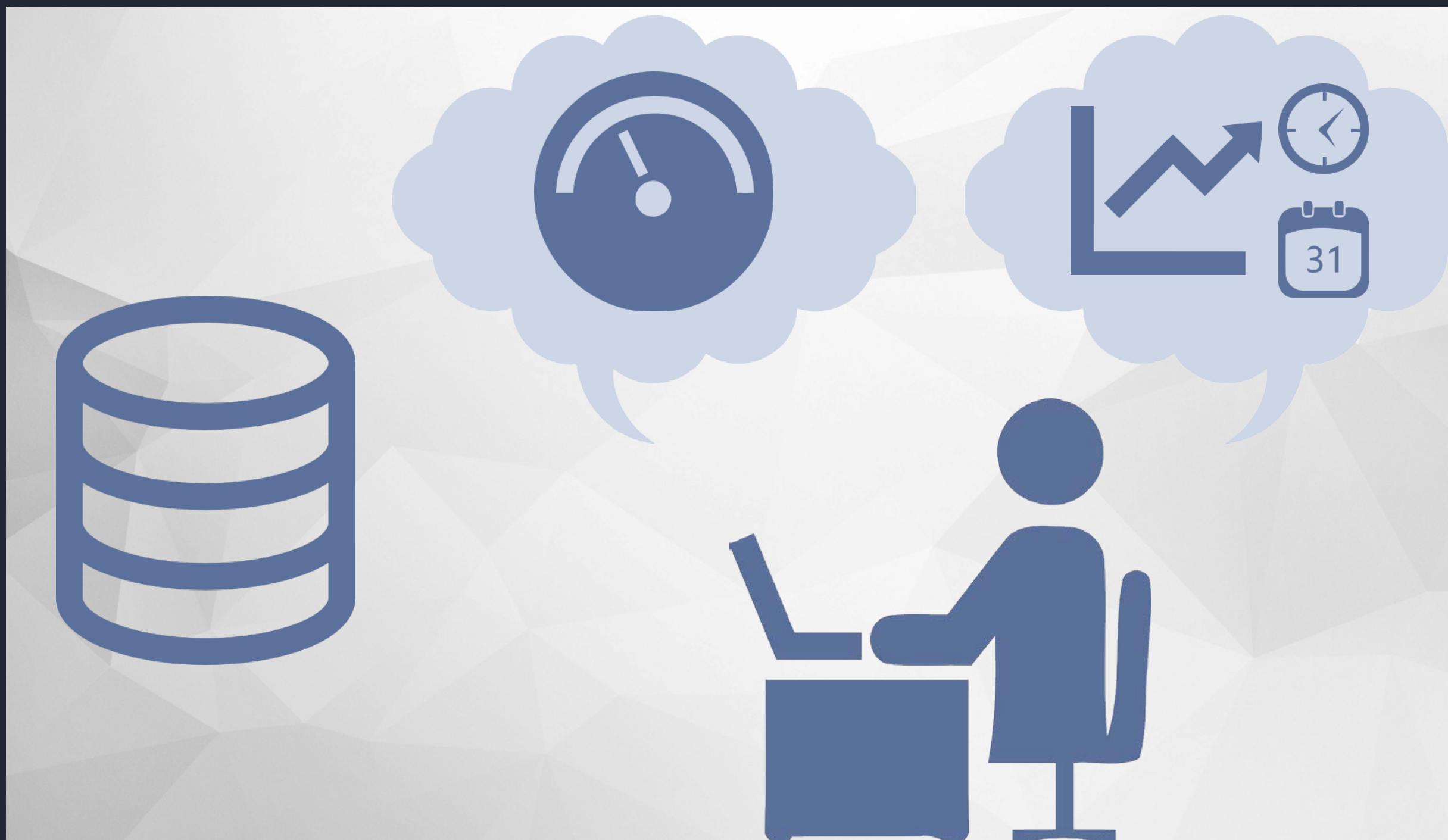
Sweet KIWI: Statistics-Driven OLAP Acceleration using Query Column Sets



Why?

5

- DBMSs are now the critical part of every *data-intensive* application.
- Existing *automated tuning tools* require too much human intervention.
- Average salary for database administrators (DBAs) in 2015 was \$81,710.



Follow Us | What's New | Release Calendar | Blog
Search BLS.gov

BUREAU OF LABOR STATISTICS

Home ▾ Subjects ▾ Data Tools ▾ Publications ▾ Economic Releases ▾ Students ▾ Beta ▾

OOH HOME | OCCUPATION FINDER | OOH FAQ | OOH GLOSSARY | A-Z INDEX | OOH SITE MAP | EN ESPAÑOL

OCCUPATIONAL OUTLOOK HANDBOOK

Computer and Information Technology > Database Administrators

EN ESPAÑOL

Summary [What They Do](#) [Work Environment](#) [How to Become One](#) [Pay](#) [Job Outlook](#) [State & Area Data](#) [Similar Occupations](#) [More Info](#)

Summary

Quick Facts: Database Administrators	
2015 Median Pay	\$81,710 per year \$39.29 per hour
Typical Entry-Level Education	Bachelor's degree
Work Experience in a Related Occupation	Less than 5 years
On-the-job Training	None
Number of Jobs, 2014	120,000
Job Outlook, 2014-24	11% (Faster than average)
Employment Change, 2014-24	13,400

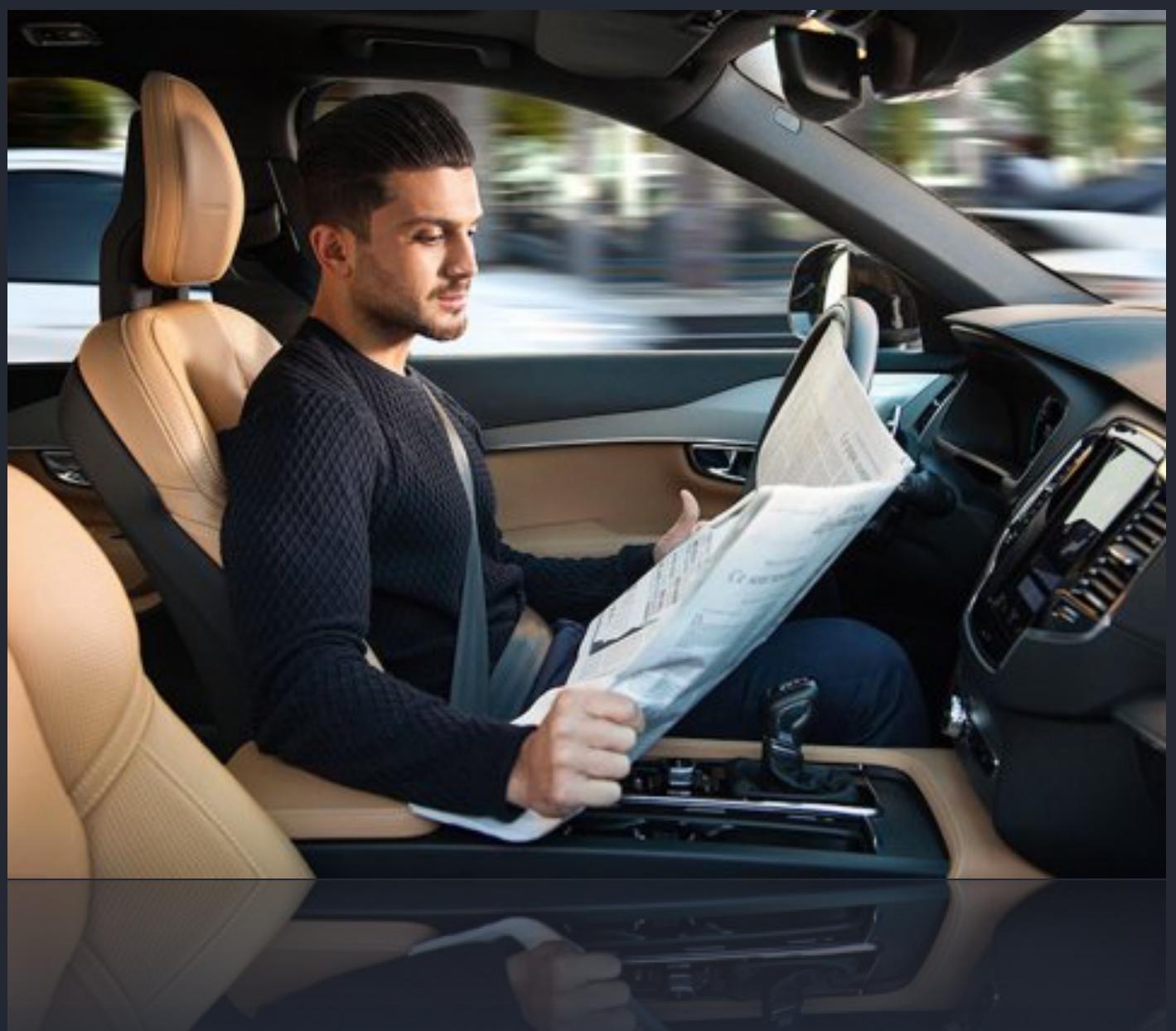
What Database Administrators Do

Database administrators (DBAs) use specialized software to store and organize data, such as financial information and customer shipping records. They make sure that data are available to users and are secure from unauthorized access.

Database administrators ensure that data are available to many different users.

SELF-DRIVING

A DBMS THAT CAN CONFIGURE,
TUNE, AND OPTIMIZE ITSELF
WITHOUT ANY HUMAN
INTERVENTION.



Can We Automate This?

7

YES

DATABASE DESIGN

DATA PLACEMENT

QUERY OPTIMIZATION

KNOB CONFIGURATION

BACK-UP & RECOVERY

PROVISIONING

NO

SECURITY & ACLS

DATA INTEGRATION

UNPLANNED HALTS

VERSION CONTROL

WHAT'S NEW?

PREVIOUS EFFORTS ARE
REACTIVE & HUMAN-DRIVEN.

A SELF-DRIVING DBMS HAS TO
BE PREDICTIVE.

WHY NOW?

RECENT ADVANCEMENTS IN
HARDWARE AND DEEP NEURAL
NETWORKS MAKE AUTONOMOUS
OPERATION NOW POSSIBLE.



Peloton

Non-volatile ↗

I N - M E M O R Y

O L T P + O L A P

L L V M E X E C

A U T O N O M O U S

THE BRAIN

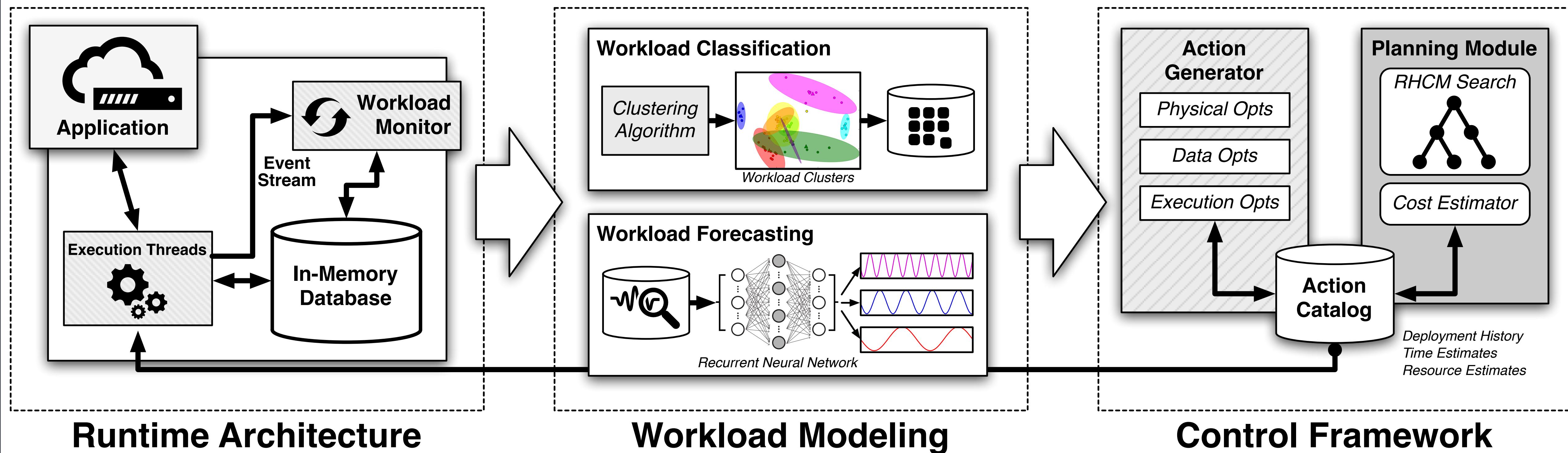
INTEGRATED DEEP LEARNING
FRAMEWORK TO MODEL,
PREDICT, AND OPTIMIZE HTAP
DATABASE WORKLOADS.



Self-Driving Database Management Systems
CIDR 2017

- Features**
- (1) *a query's runtime metrics*
 - (2) *a query's logical semantics*

DBSCAN algorithm



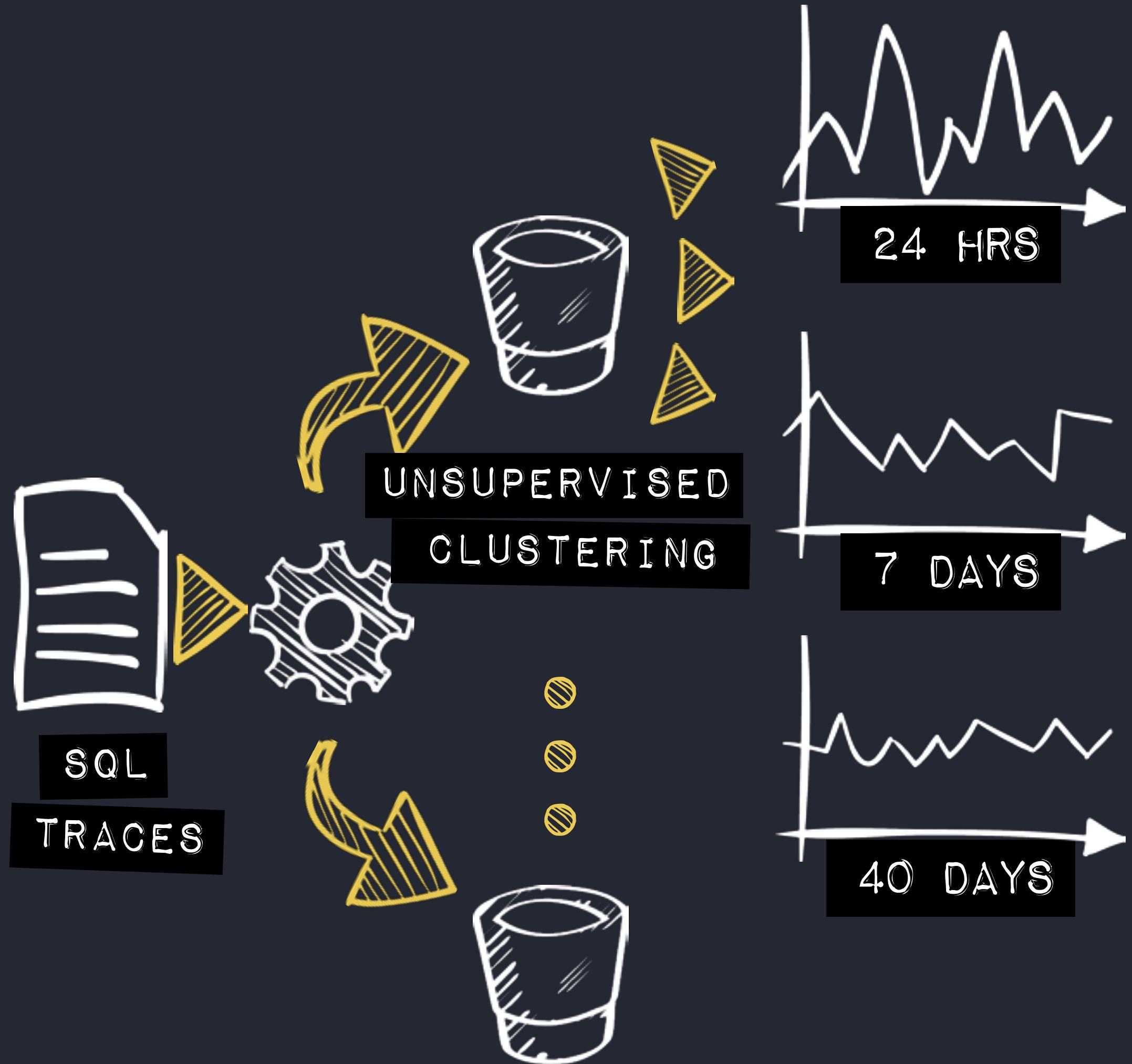
Recurrent neural networks (RNNs)

Previous attempts at autonomous systems:
Auto-regressive-moving average model (ARMA)

Long short-term memory (LSTM)

WORKLOAD CATEGORIZATION

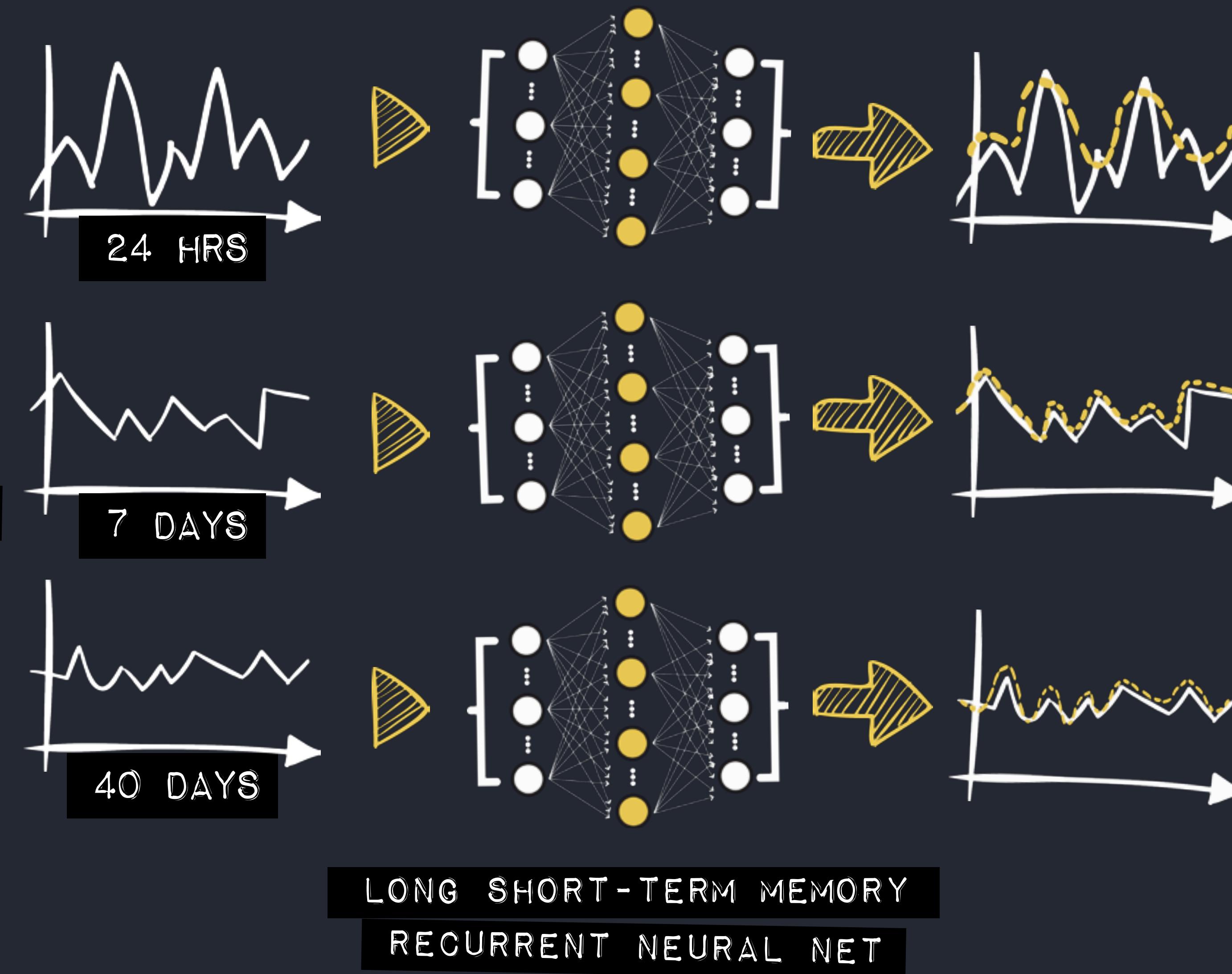
13



WORKLOAD CATEGORIZATION



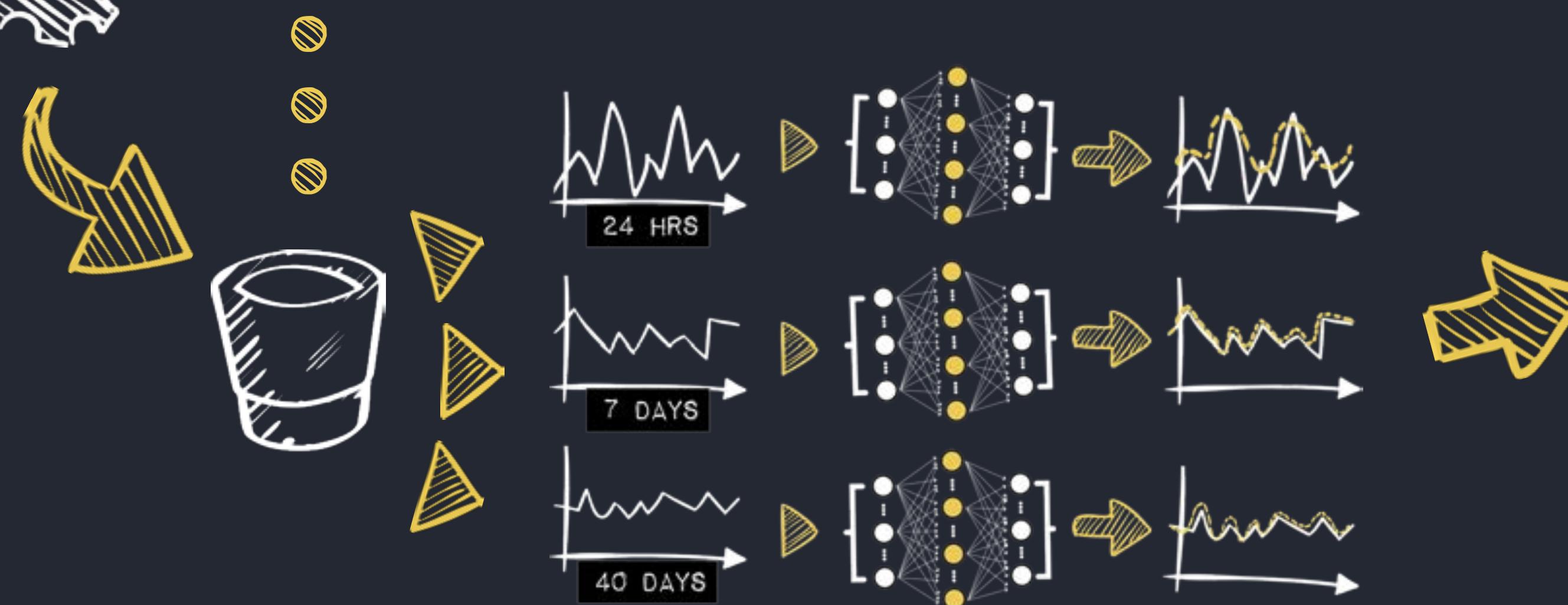
WORKLOAD FORECASTING



WORKLOAD CATEGORIZATION



WORKLOAD FORECASTING

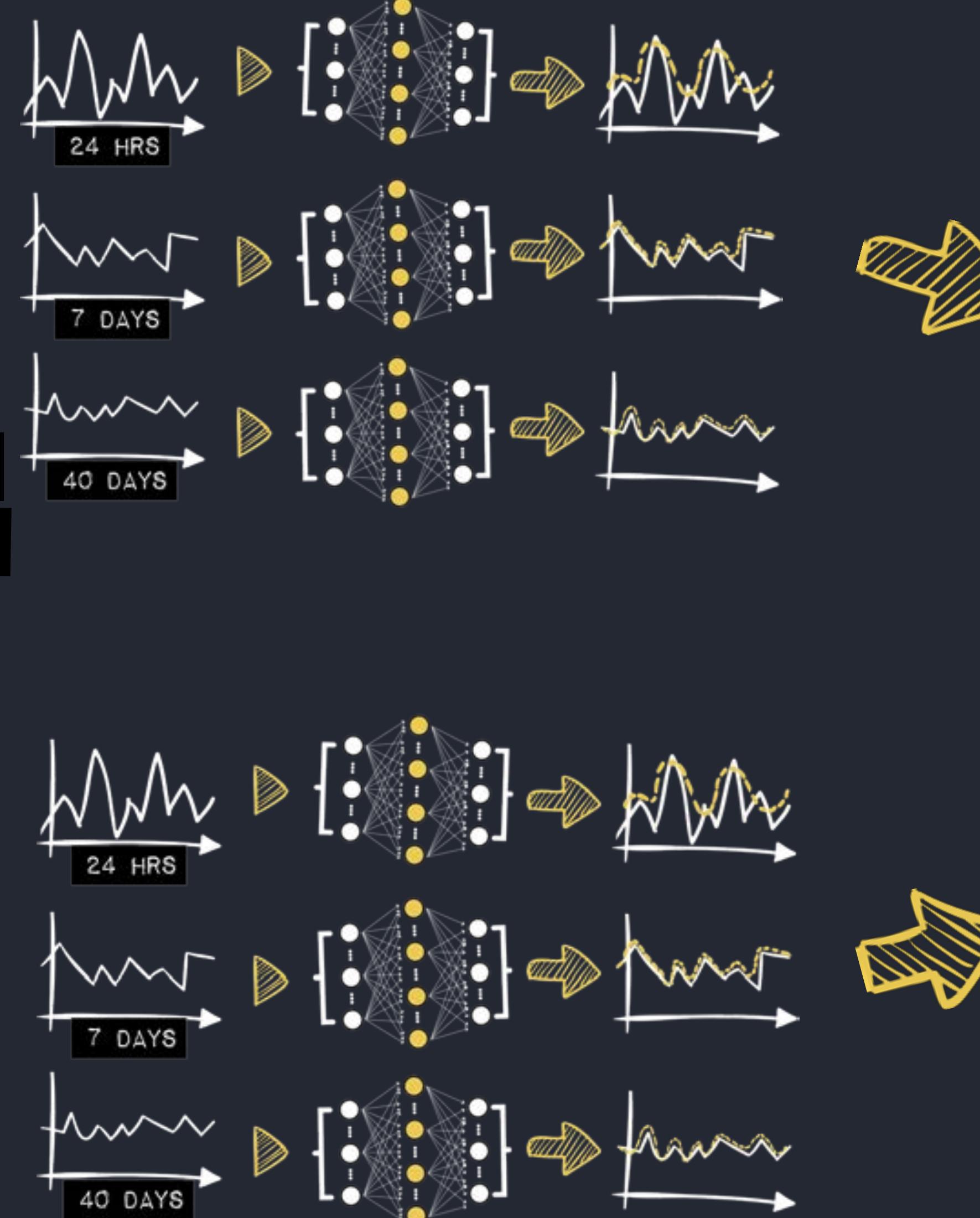


OPTIMIZATION PLANNING

WORKLOAD CATEGORIZATION



WORKLOAD FORECASTING



OPTIMIZATION PLANNING

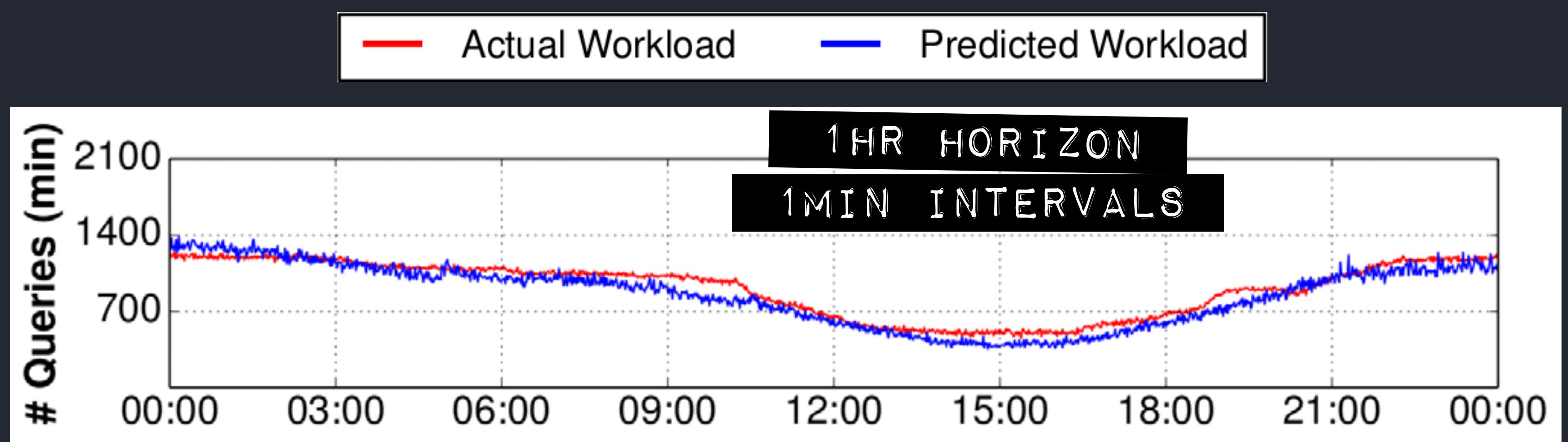


EVALUATION

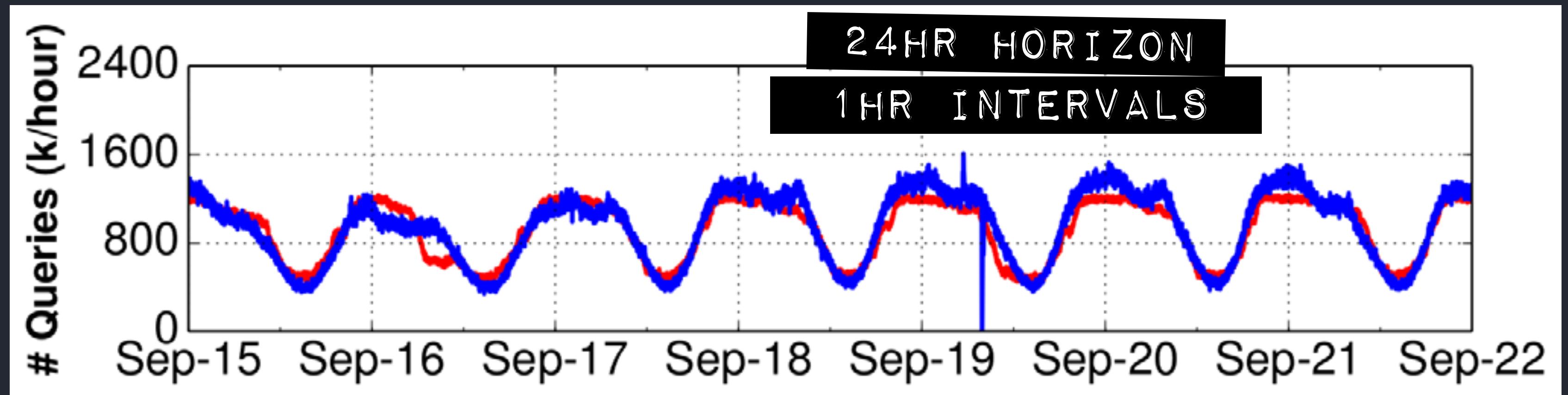
SYNTHETIC WORKLOAD BASED
ON REDDIT TRAFFIC DATA .

FORECAST WITH TENSORFLOW .
ADAPTIVE STORAGE .





ERROR RATE: 14.7%
CPU TRAINING: 25MIN
PROBE: 2MS
UPDATE: 5MS
SIZE: 2MB



ERROR RATE: 17.9%
CPU TRAINING: 18MIN
PROBE: 2MS
UPDATE: 5MS
SIZE: 2MB

ADAPTIVE STORAGE

CHANGE THE LAYOUT OF DATA
OVER TIME BASED ON HOW IT
IS ACCESSED.



Bridging the Archipelago Between Row-Stores and
Column-Stores for Hybrid Workloads
SIGMOD 2016

O L T P



```
UPDATE myTable  
SET A = 123,  
     B = 456,  
     C = 789  
WHERE D = "xxx"
```

Hot

ORIGINAL TABLE

Cold

OLAP



```
SELECT AVG(B)  
FROM myTable  
WHERE C < "yyy"
```

OLTP



```
UPDATE myTable
SET A = 123,
      B = 456,
      C = 789
      D = "xxx"
WHERE
```

Hot

OLAP



```
SELECT AVG(B)
FROM myTable
WHERE C < "yyy"
```

ORIGINAL TABLE

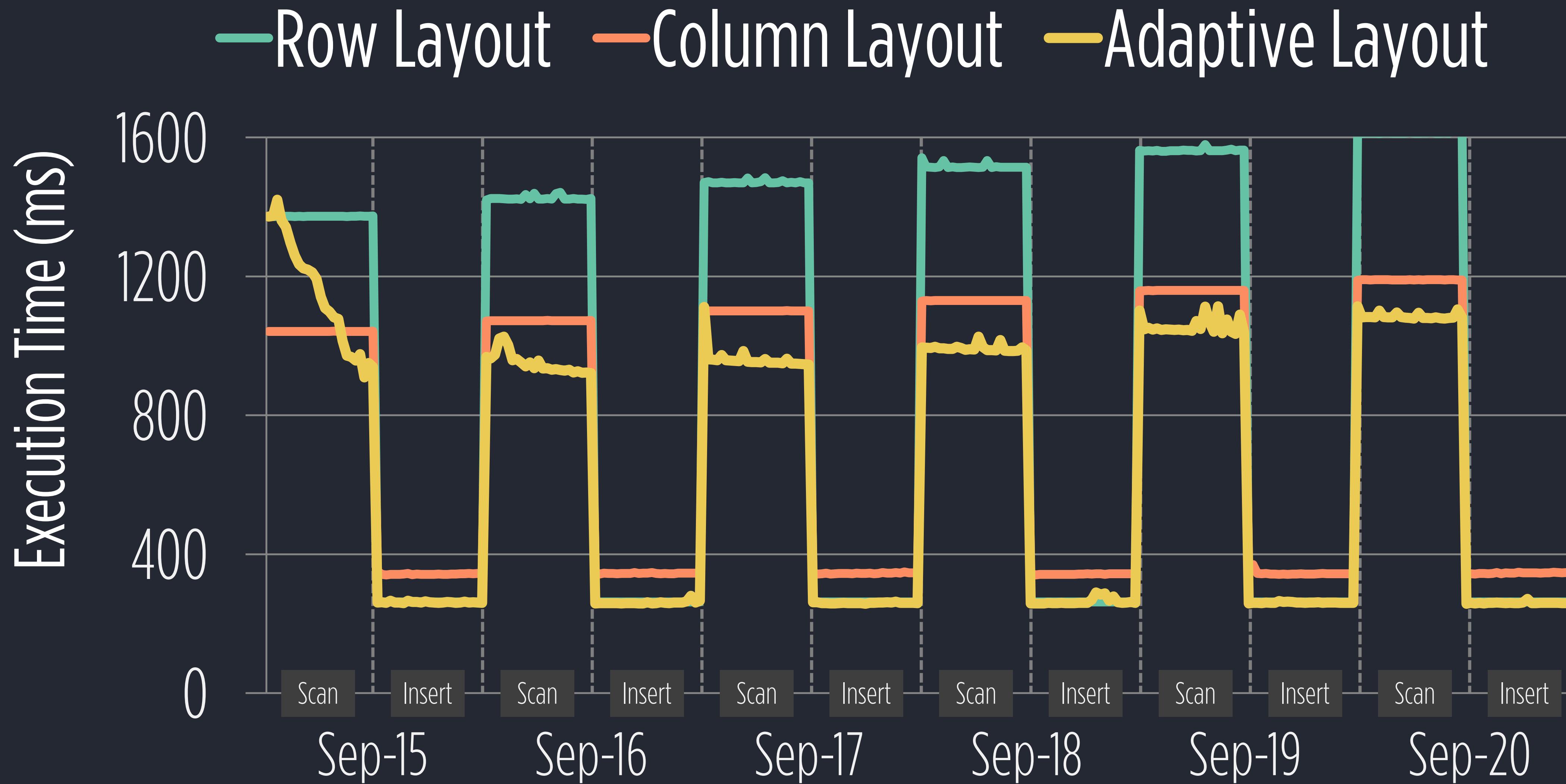
A	B	C	D

Cold

ADAPTED TABLE

A	B	C	D

A	B	C	D



NVM Optimizations

- » Avoid the overhead of DBMS recovery
- » Larger-than-Memory databases
- » Write-Behind Logging
- » Data Tiering

Write-Behind Logging
VLDB2016

Write-Behind Logging

Joy Arulraj
Carnegie Mellon University
jarulraj@cs.cmu.edu

Matthew Perron
Carnegie Mellon University
mperron@cmu.edu

Andrew Pavlo
Carnegie Mellon University
pavlo@cs.cmu.edu

ABSTRACT

The design of the logging and recovery components of database management systems (DBMSs) has always been influenced by the difference in the performance characteristics of volatile (DRAM) and non-volatile storage devices (e.g., SSDs). Until now, the key assumption has been that non-volatile storage is much slower than DRAM and can only support block-oriented reads/writes. But the arrival of non-volatile memory (NVM) storage that is almost as fast as DRAM invalidates these previous design choices.

This paper explores the unique properties of these NVM devices in systems that still include volatile DRAM. We make the case for a recovery approach, called write-behind logging, that enables a DBMS to recover almost instantaneously from system failures. Our evaluation of this recovery algorithm in an in-memory DBMS shows that across different workloads it reduces recovery time by 20x and improves the lifetime of an NVM device by 2x.

1. INTRODUCTION

The data storage and retrieval requirements of modern service-oriented applications has given rise to a new class of on-line transaction processing (OLTP) applications that have strict demands than application from previous years. The core component of these applications are high-performance DBMSs. The data processing capacity of the OLTP applications therefore depends on the performance of the DBMS. This has always been constrained by the overhead associated with storing the data on a non-volatile storage device and its durability.

The DBMS is responsible for ensuring that the database state is not corrupted due to hardware, system, or device failure [23]. It ensures the durability of all updates to the database by writing out changes to a non-volatile device such as a SSD, before committing the transaction and returning an acknowledgement back to the application. These devices can only support slow bulk data transfers as blocks. In contrast, a DBMS can quickly read and write a single byte from a volatile DRAM device, where all data is lost once power is lost. The difference in the characteristics of volatile and non-volatile storage devices has always influenced the design of the recovery component of the DBMS.

In-memory DBMSs avoid the overhead of managing disk-resident data by storing the entire database in main memory [19]. These systems outperform disk-oriented DBMSs by virtue of the memory-oriented design. The fundamental problem with these memory-oriented DBMSs, however, is that they assume that the size of the database is smaller than the size of DRAM in the system. DRAM has inherent physical limitations that limit its scaling capabilities beyond today's levels [34]. Given these constraints, we anticipate that large OLTP databases will no longer fit within DRAM. As a consequence, running OLTP DBMSs on a large amount of SSDs incurs a lot of energy since it requires periodic refreshing to preserve data that is not even actively used. Studies have shown that DRAM can account for 40% of the total power consumption in commercial servers [30]. Flash-based SSDs have larger storage capacities and expand less than DRAM. But they are much slower than DRAM and can only support block-based access methods. Even if a transaction only updates a single byte of a tuple stored on SSD, the DBMS must write out the changes in a block (typically 4 KB). OLTP applications exacerbate this fundamental constraint of SSDs as they inherently perform many small changes to a database. Stop-gap solutions, such as battery-backed DRAM caches, can help shrink the performance gap but do not resolve these other challenges.

Non-volatile memory (NVM) is a broad class of technologies, including phase-change memory [41], memristors [43], and STT-MRAM [20] but provide low latency reads and writes on the same order of magnitude as DRAM. They support persistent writes and large storage capacity like a SSD [10]. Table 1 provides a comparison of the characteristics of NVM against other storage technologies.

The unique characteristics of NVM devices necessitate changes in the design of existing DBMS architectures [11, 17]. Disk-oriented DBMSs (e.g., Oracle RDBMS, MySQL DB2, MySQL) are predicted on using block-oriented devices for data storage that are slow at random access. To avoid expensive retrievals from disk within the critical path of the transaction, they maintain a memory cache for blocks of tuples and try to maximize the amount of sequential reads and writes to storage. Memory-oriented DBMSes (e.g., MySQL) also contain special components to overcome the volatility of DRAM and ensure the durability of data. The design of these with non-addressable NVM that supports fast sequential access.

In this paper, we present a recovery algorithm that is designed for a hybrid storage stack comprising of DRAM and NVM, implemented two different recovery protocols in Peloton [21], a hybrid in-memory DBMS: (1) the write-ahead logging protocol, and (2) the write-behind logging protocol. We then analyzed the availability, storage footprints, and computational overhead of the DBMS while

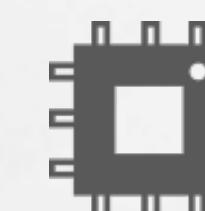
Write-Behind Logging



```
UPDATE myTable  
SET A = 123,  
WHERE C = "xxx"
```

In-Memory Heap

ver	A	B	C
x1			
x2			
x3			



PCommit

NVM Heap

A	B	C	D
x1			
x2			
x3			

NVM Log

001:Txn1-0x0001

Conclusion

25

- The demand for an *autonomous DBMS*
- The *self-driving architecture* of the Peloton DBMS
- They argued that the ambitious system that they proposed is *now possible*
 - due to *advancements in deep neural networks,*
 - *improved hardware,*
 - and *high-performance database architectures.*

