

Database Learning

Building databases that become smarter over time

Yongjoo Park

Ahmad Shahab Tajik

Michael Cafarella

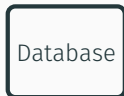
Barzan Mozafari

University of Michigan, Ann Arbor

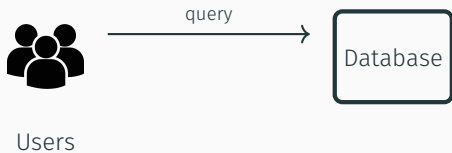
Today's Databases



Users



Today's Databases



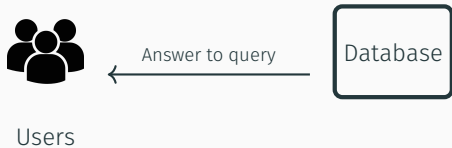
Today's Databases



Users



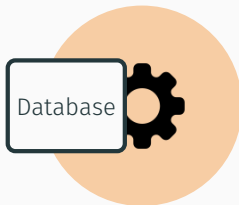
Today's Databases



Today's Databases



Users

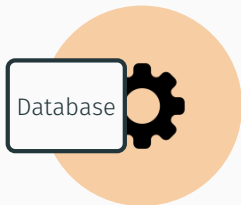


After answering queries,
THE WORK is almost completely **WASTED**.

Today's Databases



Users



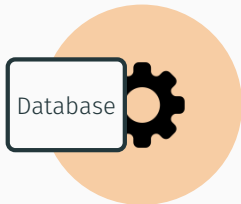
After answering queries,
THE WORK is almost completely **WASTED**.

Small exceptions:

Today's Databases



Users



After answering queries,
THE WORK is almost completely **WASTED**.

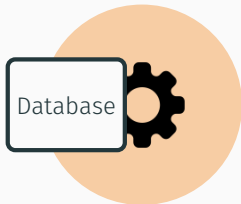
Small exceptions:

- Caching

Today's Databases



Users



After answering queries,
THE WORK is almost completely **WASTED**.

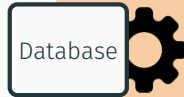
Small exceptions:

- Caching
- Identical queries

Today's Databases



Users



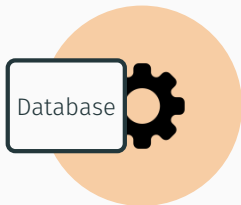
After answering queries,
THE WORK is almost completely **WASTED**.

Small exceptions:

- Caching
- Identical queries
- Indexing/Materialization hints



Users



After answering queries,

THE WORK is almost completely **WASTED**.

Small exceptions:

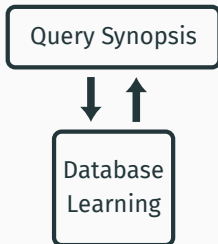
- Caching
- Identical queries
- Indexing/Materialization hints

Our Goal: reuse **the work.**

A New Paradigm in AQP Setting



Users

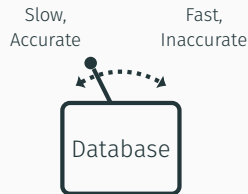
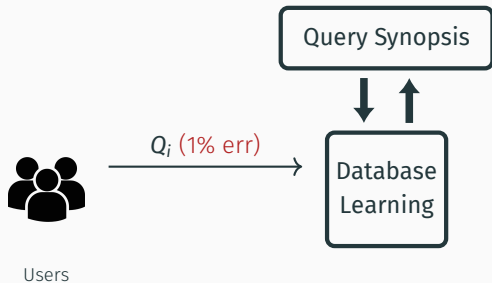


Slow,
Accurate

Fast,
Inaccurate



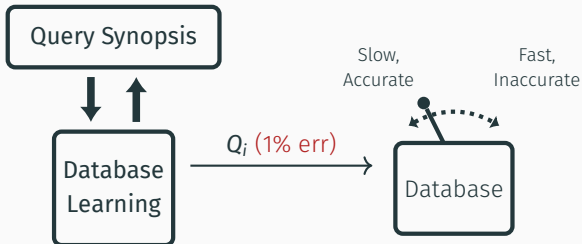
A New Paradigm in AQP Setting



A New Paradigm in AQP Setting



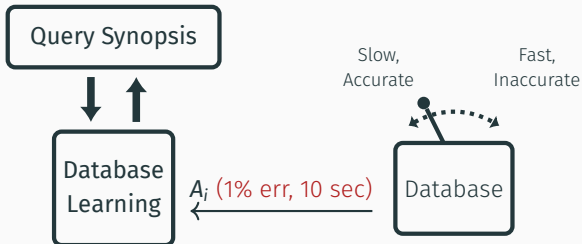
Users



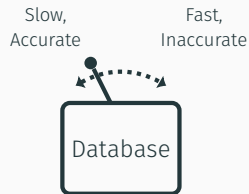
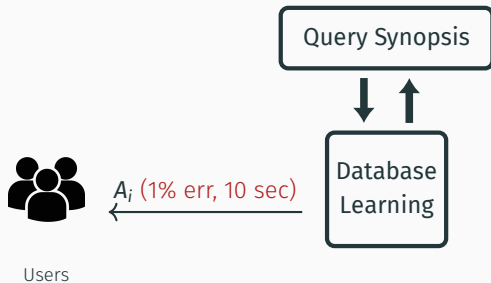
A New Paradigm in AQP Setting



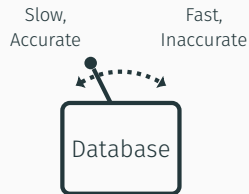
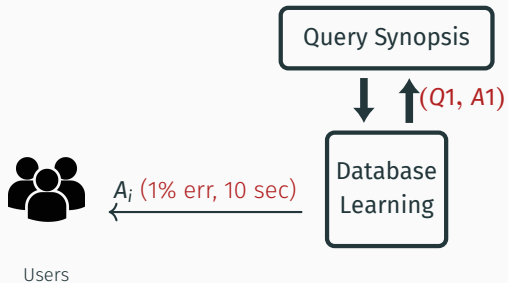
Users



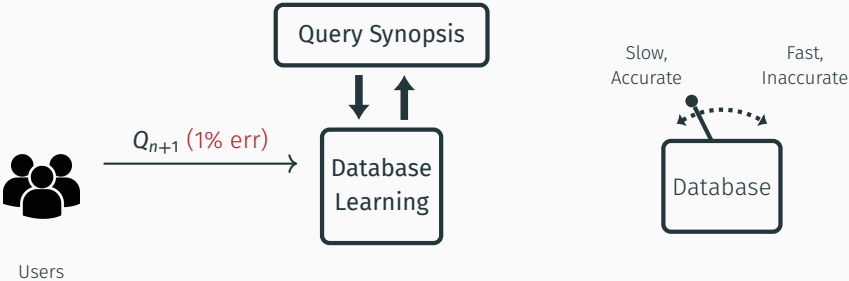
A New Paradigm in AQP Setting



A New Paradigm in AQP Setting



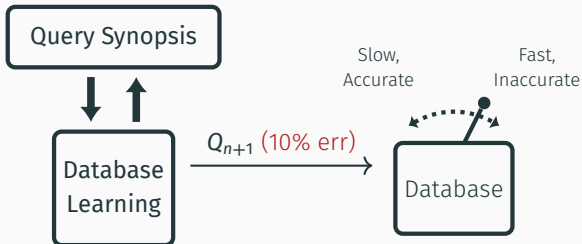
A New Paradigm in AQP Setting



A New Paradigm in AQP Setting



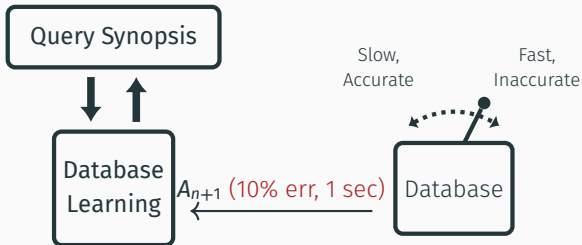
Users



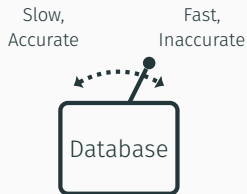
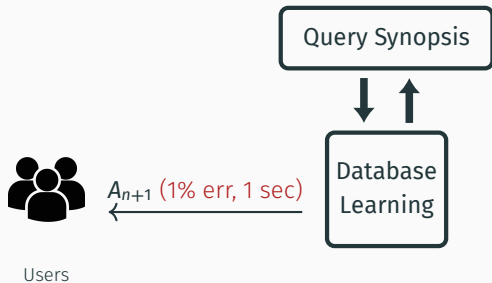
A New Paradigm in AQP Setting



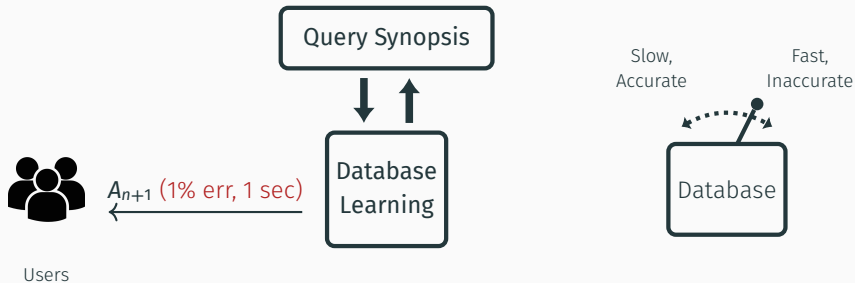
Users



A New Paradigm in AQP Setting

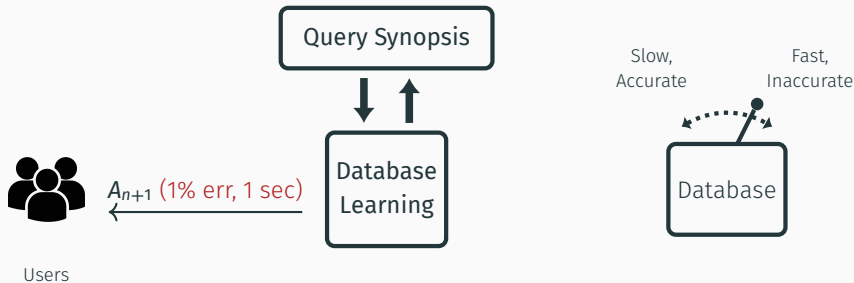


A New Paradigm in AQP Setting



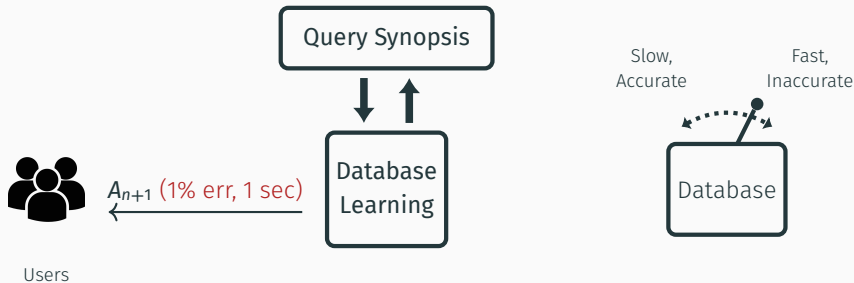
1. User: enjoys 1% error bound in 1 second!

A New Paradigm in AQP Setting



1. User: enjoys **1% error bound in 1 second!**
2. Formally, always more accurate

A New Paradigm in AQP Setting



Users

1. User: enjoys **1% error bound in 1 second!**
2. Formally, always more accurate
3. Popularity of analytic workloads \Rightarrow **Approximate solutions**
 - BlinkDB, SnappyData, Yahoo Druid, Facebook Presto, Infobright, etc.

From Machine Learning To Database Learning

Machine Learning: Past Observations \Rightarrow Future Predictions

From Machine Learning To Database Learning

Machine Learning: Past Observations \Rightarrow Future Predictions

Database Learning: **Past Answers** \Rightarrow **Future Answers**

From Machine Learning To Database Learning

Machine Learning: Past Observations \Rightarrow Future Predictions

Database Learning: Past Answers \Rightarrow Future Answers

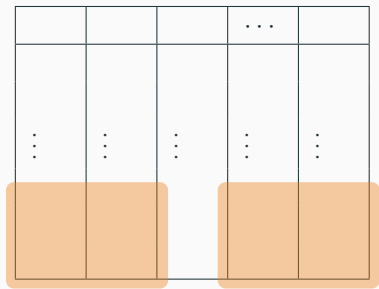


The more past queries, the more Accurate and Faster

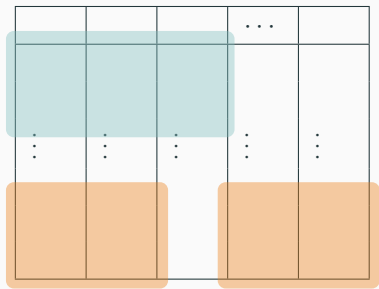
Primary Benefit: Exploratory Workloads

			...	
⋮	⋮	⋮	⋮	⋮

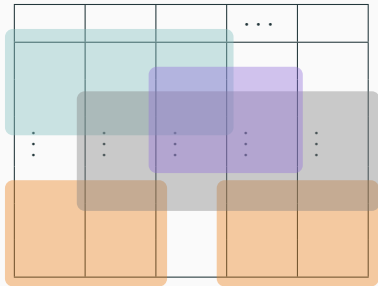
Primary Benefit: Exploratory Workloads



Primary Benefit: Exploratory Workloads

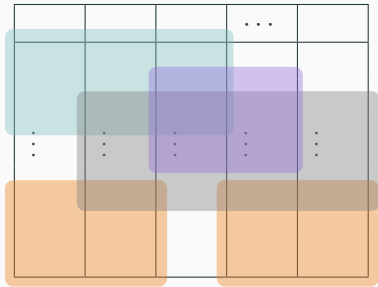


Primary Benefit: Exploratory Workloads



Queries use the data in different columns/rows.

Primary Benefit: Exploratory Workloads

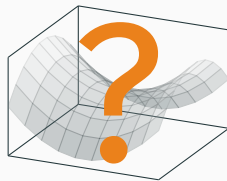


Queries use the data in different columns/rows.

How to leverage those queries for future queries?

Our Idea

			...	
⋮	⋮	⋮	⋮	⋮



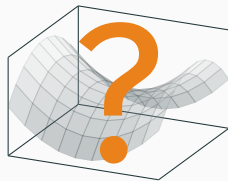
Our Idea

Q1



			...	
⋮	⋮	⋮	⋮	⋮

The table consists of five columns and three rows. The top row contains five cells, with the fourth cell containing an ellipsis (...). The second row contains five cells, each containing a vertical ellipsis (⋮). The bottom row contains five empty cells. The bottom two rows of the table are highlighted with a light orange background.

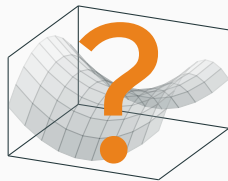


Our Idea

$(Q1, A1)$



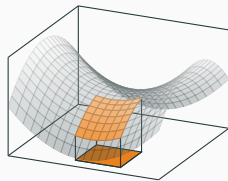
			...	
⋮	⋮	⋮	⋮	⋮



Our Idea

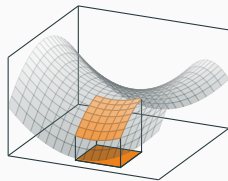
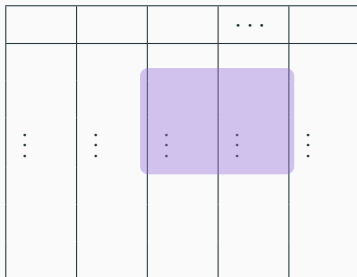
(Q_1, A_1)

			...	
⋮	⋮	⋮	⋮	⋮



Our Idea

Q2

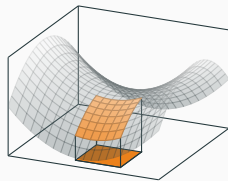


Our Idea

(Q_2, A_2)



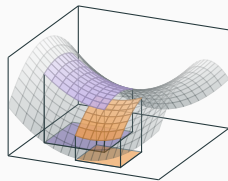
			...	
⋮	⋮	⋮	⋮	⋮



Our Idea

(Q_2, A_2)

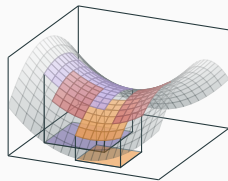
			...	
⋮	⋮	⋮	⋮	⋮



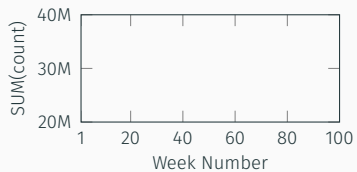
Our Idea

more queries
and answers

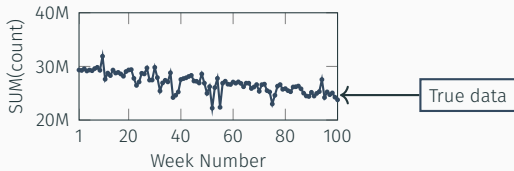
			...	
⋮	⋮	⋮	⋮	⋮



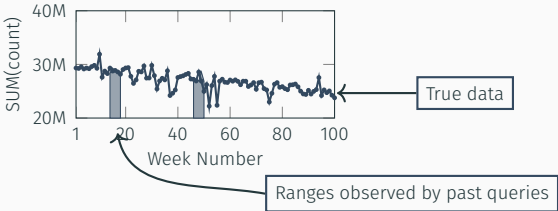
Concrete Example



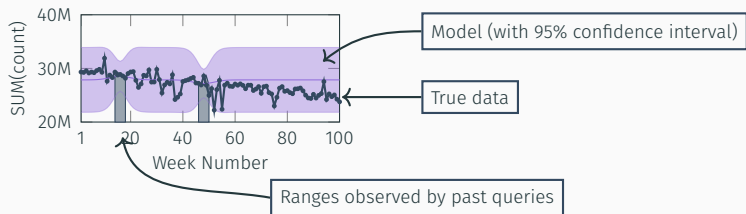
Concrete Example



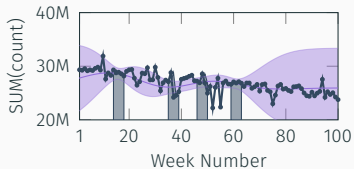
Concrete Example



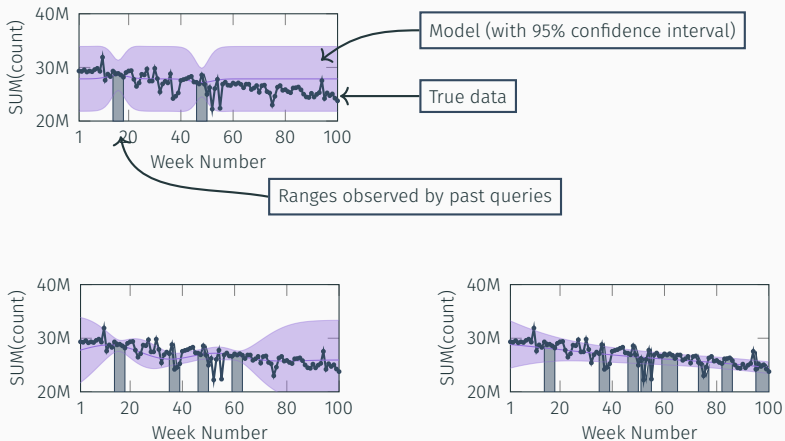
Concrete Example



Concrete Example



Concrete Example



Design Goals

```
select X3, avg(Y1)
from t
where 5 < Y1
```

```
select sum(Y2)
from t
where X2 between Apr and May
group by X3;
```

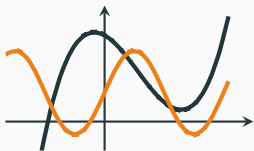
1. Support a **wide class of** SQL queries

Design Goals

```
select X3, avg(Y1)
from t
where 5 <
```

```
select sum(Y2)
from t
where X2 between Apr and May
group by X3;
```

1. Support a **wide class of** SQL queries



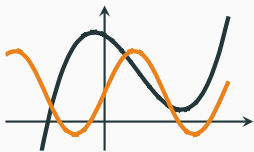
2. **No Assumptions** about Data

Design Goals

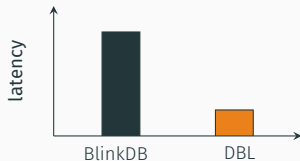
```
select X3, avg(Y1)
from t
where 5 <
```

```
select sum(Y2)
from t
where X2 between Apr and May
group by X3;
```

1. Support a **wide class of** SQL queries



2. **No Assumptions** about Data



3. **Lightweight**

Our Approach

Problem Statement

Problem Statement

Problem:

Given queries, $\{q_1, \dots, q_{n+1}\}$, and their approximate answers,

Find the **most likely** answer to q_{n+1} and **its estimated error**.

Problem Statement

Problem:

Given queries, $\{q_1, \dots, q_{n+1}\}$, and their approximate answers,

Find the **most likely** answer to q_{n+1} and **its estimated error**.

Our Result:

Under a *certain model assumption*,

our answer's error bound \leq **original answer's error bound** ,
(in practice, much more accurate)

if the error bounds provide the same probabilistic guarantees.

Our Model Assumption

Assumption:

Goodness of the principle of maximum entropy [?]

Our Model Assumption

Assumption:

Goodness of the principle of maximum entropy [?]

→ Not overly confident on the unobserved data.

Our Model Assumption

Assumption:

Goodness of the principle of maximum entropy [?]

→ Not overly confident on the unobserved data.

Justification:

“ agrees with everything that is known, but carefully avoids assuming anything that is not known — [?] ”

Our Model Assumption

Assumption:

Goodness of the **principle of maximum entropy** [?]
→ **Not** overly confident on the **unobserved data**.

Justification:

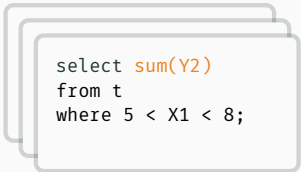
“ agrees with everything that is known, but carefully avoids assuming anything that is not known — [?] ”

We provide **empirical justifications** in our report [?].

Overview of Our Technique

```
select avg(Y2)
from t
where 6 < X1 < 8;
```

Overview of Our Technique



```
select sum(Y2)
from t
where 5 < X1 < 8;
```

Overview of Our Technique

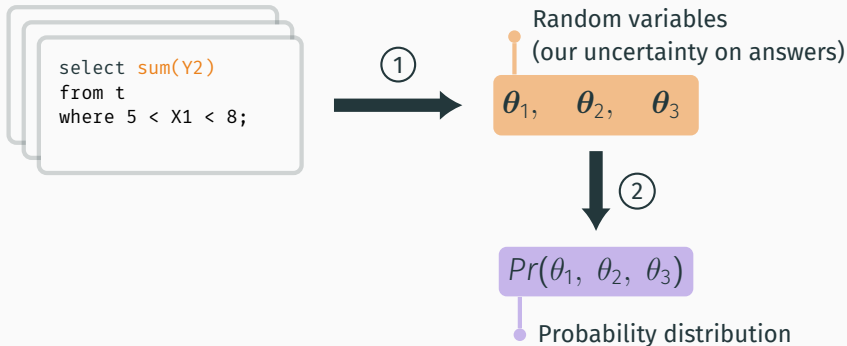
```
select sum(Y2)
from t
where 5 < X1 < 8;
```



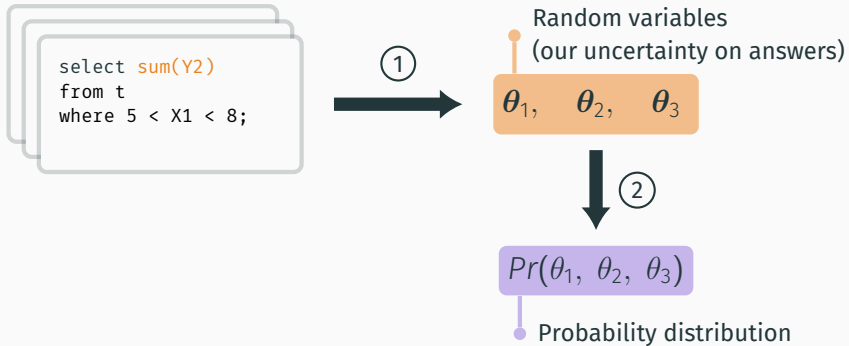
Random variables
(our uncertainty on answers)

$\theta_1, \theta_2, \theta_3$

Overview of Our Technique

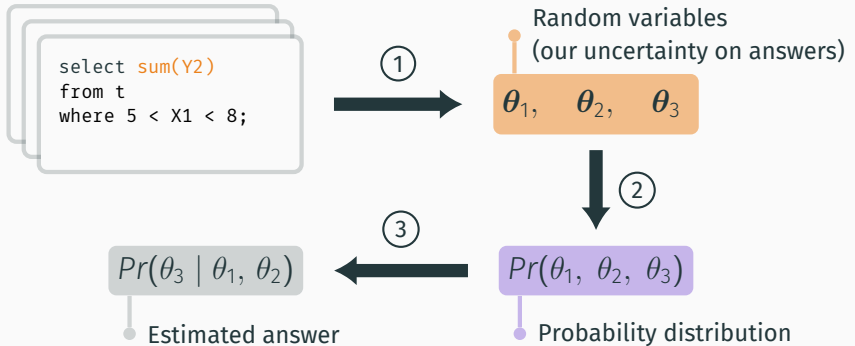


Overview of Our Technique



Two aggregations involve **common values**
→ **correlation** between answers

Overview of Our Technique



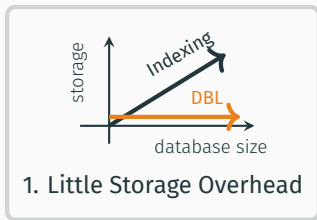
Two aggregations involve **common values**
→ **correlation** between answers

Benefits of Database Learning

Database Learning vs. Indexing

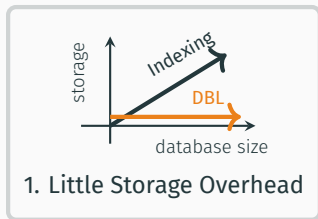
Benefits of Database Learning

Database Learning vs. Indexing



Benefits of Database Learning

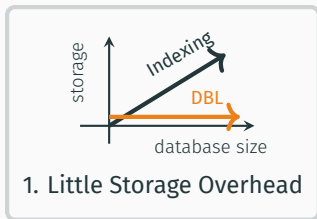
Database Learning vs. Indexing



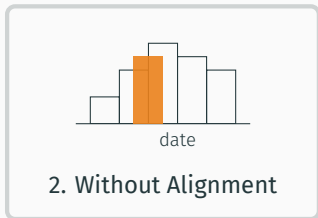
Database Learning vs. Materialized View Selection

Benefits of Database Learning

Database Learning vs. Indexing

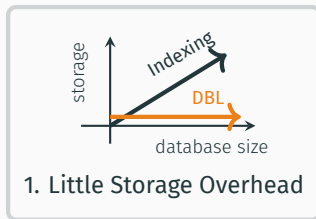


Database Learning vs. Materialized View Selection

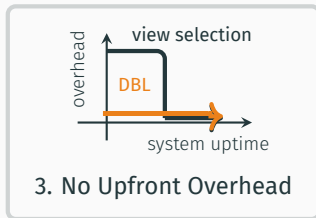
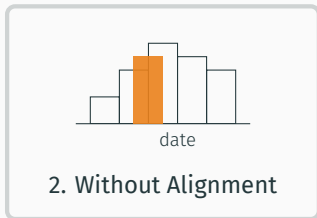


Benefits of Database Learning

Database Learning vs. Indexing



Database Learning vs. Materialized View Selection



Experiment



1. Using **Spark SQL** as a backend

- NOLEARN: Sampling-based AQP engine
- VERDICT: Our database learning system

Implementation and Experiments



1. Using **Spark SQL** as a backend

- NOLEARN: Sampling-based AQP engine
- VERDICT: Our database learning system

2. Datasets:

- **Customer1**: Query log from an analytic DB vendor
- TPC-H: 100G TPC-H dataset

Implementation and Experiments

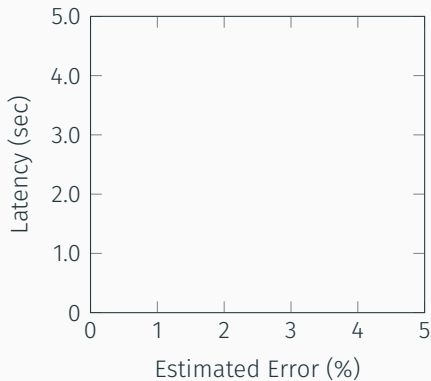


1. Using **Spark SQL** as a backend
 - NOLEARN: Sampling-based AQP engine
 - VERDICT: Our database learning system
2. Datasets:
 - **Customer1**: Query log from an analytic DB vendor
 - TPC-H: 100G TPC-H dataset
3. Environment:
 - 5 Amazon EC2 workers (**m4.2xlarge**)
 - SSD-backed HDFS for Spark's data loading

Generality of VERDICT

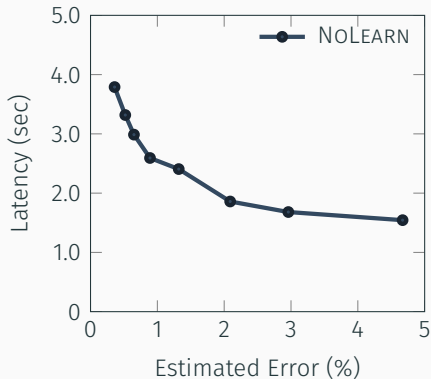
Dataset	# Analyzed	# Supported	Percentage
Customer1	3,342	2,463	73.7%
TPC-H	21	14	63.6%

Latency-Error Trade-off



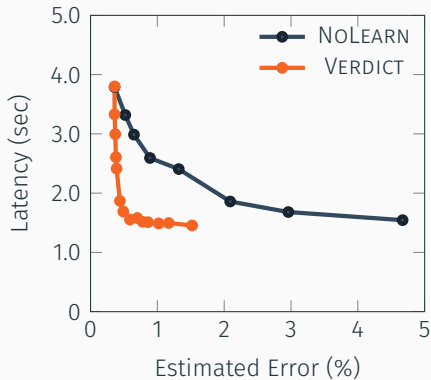
Customer1 dataset in memory

Latency-Error Trade-off



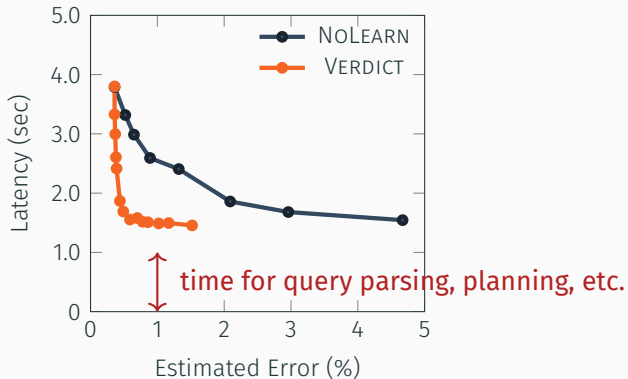
Customer1 dataset in memory

Latency-Error Trade-off



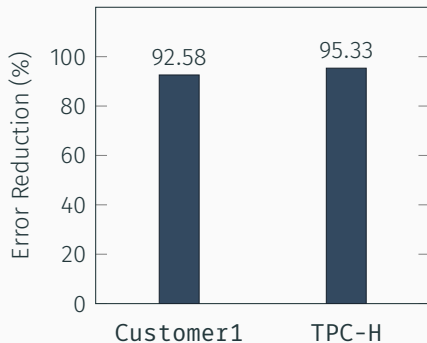
Customer1 dataset in memory

Latency-Error Trade-off



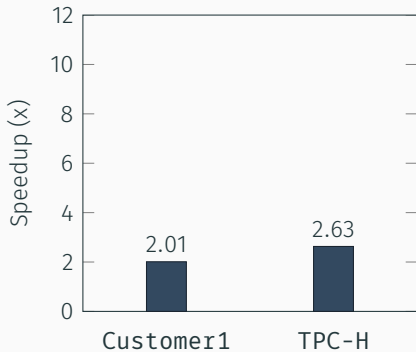
Customer1 dataset in memory

Error Reduction



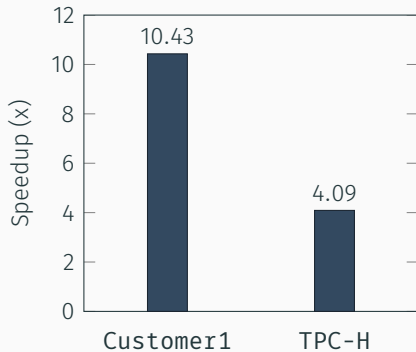
For the same time budget; Data on SSD

Speedup (Data in Memory)



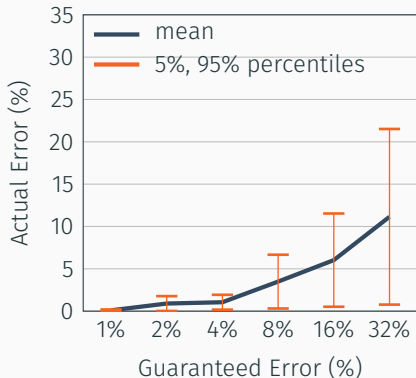
For the same target error of 2%.

Speedup (Data on SSD)



For the same target error of 2%.

Reliability of Estimated Error Guarantees



Guaranteed error = from 95% confidence interval

Memory and Computational Overhead

1. Memory Overhead:

- 150KB per query for the Customer1 dataset
- 8KB per query for the TPC-H dataset

Memory and Computational Overhead

1. Memory Overhead:

- 150KB per query for the Customer1 dataset
- 8KB per query for the TPC-H dataset

2. Computational Overhead:

Latency	Cached	No-Cache
NOLEARN	2.083 sec	52.50 sec
VERDICT	2.093 sec	52.51 sec
Overhead	0.010 sec (0.48%)	0.010 sec (0.02%)

1. Database Learning:

Answers to past queries → boost your AQP!

Conclusion and Future Work

1. Database Learning:

Answers to past queries → boost your AQP!

2. Our prototype, **VERDICT**, demonstrated:

Conclusion and Future Work

1. Database Learning:

Answers to past queries → boost your AQP!

2. Our prototype, **VERDICT**, demonstrated:

- Support **73.7%** real-world analytical queries

1. Database Learning:

Answers to past queries → boost your AQP!

2. Our prototype, **VERDICT**, demonstrated:

- Support **73.7%** real-world analytical queries
- Error reduction up to **95%** compared to existing AQP engines

Conclusion and Future Work

1. Database Learning:

Answers to past queries → boost your AQP!

2. Our prototype, **VERDICT**, demonstrated:

- Support **73.7%** real-world analytical queries
- Error reduction up to **95%** compared to existing AQP engines

Our next goal: Active Database Learning

Aims to build a probabilistic model of data
even before any queries submitted

Thank You!

