# Learning-based Cost Management for Cloud Databases

**Olga Papaemmanouil**

**Brandeis University**

# Outline

Motivation

Offline Learning

Online Learning

Conclusions

# Outline
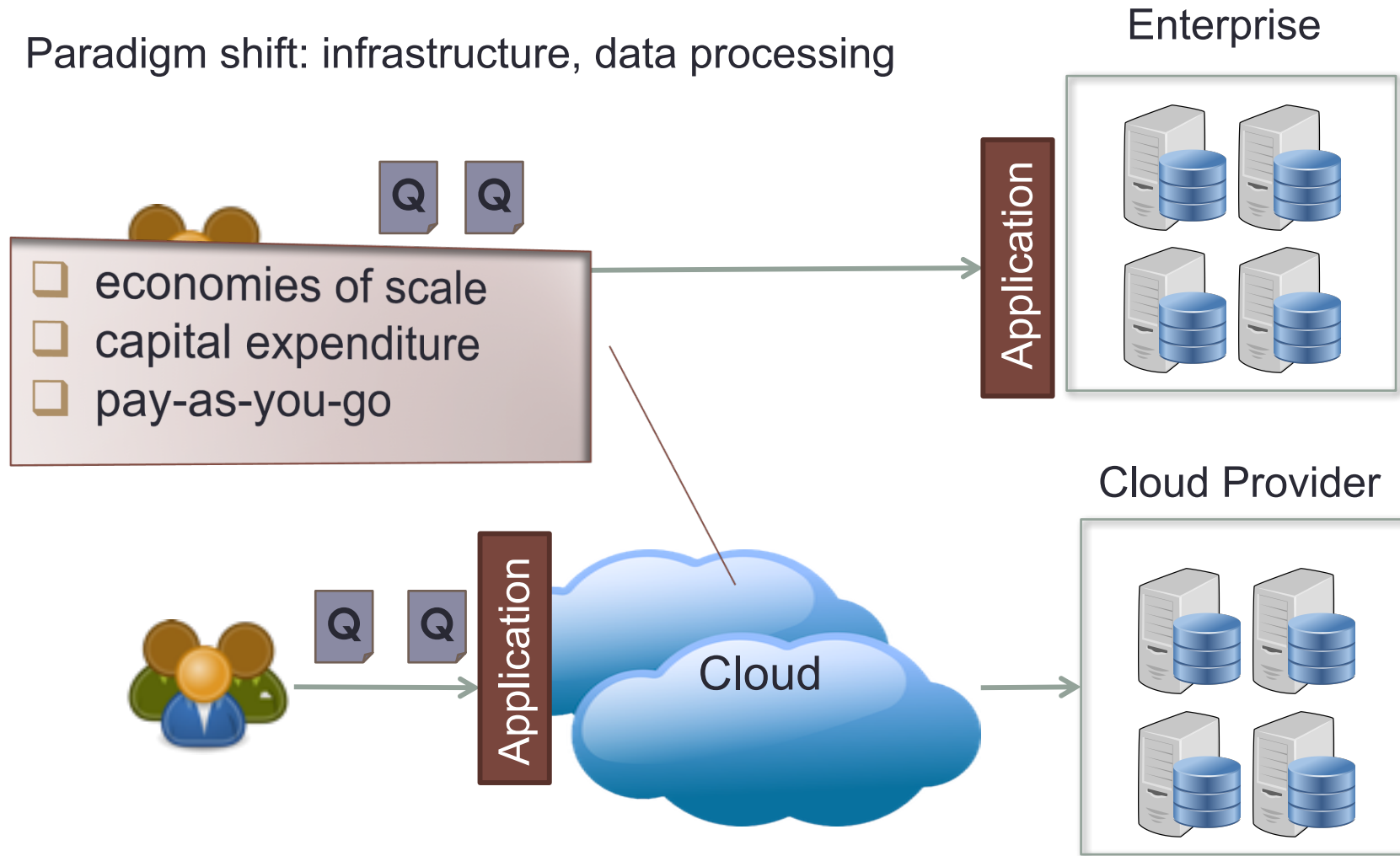
**Motivation**

Offline Learning

Online Learning

Conclusions

- ☐ Cloud Databases
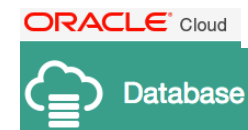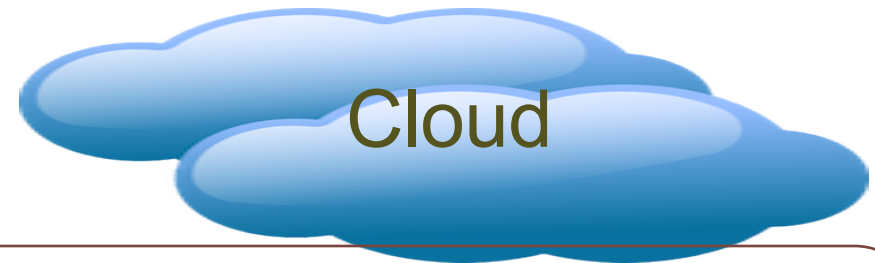- ☐ Challenges
- ☐ Why Machine Learning ?

# Cloud Computing

Paradigm shift: infrastructure, data processing

Enterprise

Q Q

- ❑ economies of scale
- ❑ capital expenditure
- ❑ pay-as-you-go

Application

Q Q

Application

Cloud

Cloud Provider

# Cloud Databases Landscape

Cloud

## Database-as-a-Service
- ❏ Managed DBMS
- ❏ Relational & NoSQL DBs

Microsoft® SQL Azure™

ORACLE Cloud
Database

Amazon RDS

**Infrastructure as a Service (IaaS)**

rackspace®
the open cloud company

amazon web services™

Azure

Google Compute Engine

## IaaS-based DB Instances
- ❏ Non managed DBMS
- ❏ Do It Yourself model

# IaaS-deployed Databases

**App Management Tools**

- Monitoring resources, performance, cost

- Event-driven scaling

- NO cost vs performance optimization

OpsWorks

Trusted Advisor
AWS Cloud Optimization Expert

StackDriver Monitoring

amazon web services CloudWatch

Q Q Q Q

Data Management Application

Microsoft SQL Server

MySQL

PostgreSQ

ORACLE

IaaS Provider

VM VM VM VM

# Deployment Challenges

Q Q Q Q

## Data Management Application

Custom-built application management tools

Microsoft SQL Server  MySQL  PostgreSQL  ORACLE

IaaS Provider

VM VM VM VM

# Deployment Challenges

**Meet SLOs**
(Service Level Objective)

- ❑ Query-level: response time
- ❑ Workload level: average, total, max, percentile

**Offer SLAs**
(Service Level Agreement)

- ❑ SLO+ Violation penalties

**Pay-as-you-go Model**

Q Q Q Q

## Data Management Application

**Cost Management**

**Performance Management**

Microsoft® SQL Server®

MySQL®

PostgreSQL

ORACLE®

IaaS Provider

$$  VM

$$  VM

VM  $$

# Deployment Challenges

**NP-hard problem**

## Beyond monitoring & alerts

- ❑ Automatic scale up & down
- ❑ Query routing & scheduling
- ❑ Cost-driven decisions
- ❑ SLA-awareness

Q Q Q Q

## Data Management Application

| Cost Management | Performance Management |
| Resource Provisioning | Workload Scheduling |

Microsoft SQL Server   MySQL   PostgreSQL   ORACLE

**IaaS Provider**

VM  VM  VM  VM

# State-of-the-art

| Placement | Provisioning | | Scheduling |
|---|---|---|---|
| **PMAX** (Liu et al.) | **Auto** (Rogers et al.) | **Dolly** (Cecchet et all) | **Shepherd** (Chi et al.) |
| **SLATree** (Chi et al.) | | | |
| **Multi-tenant SLOs** (Lang et al.) | | | **iCBS** (Chi et al.) |
| **Delphi / Pythia** (Elmore et al.) | **Hypergraph** (Çatalyürek et al.) | | |
| **SCOPE** (Chaiken et al.) | **Bazaar** (Jalaparti et al.) | | many traditional methods ... |

# State-of-the-art

| Query deadline | Workload deadline |
|---|---|
| Average latency | Percentile deadline |
| | Piecewise linear |

| Placement | Provisioning | | Scheduling |
|---|---|---|---|
| PMAX (Liu et al.) | Auto (Rogers et al.) | Dolly (Cecchet et all) | Shepherd (Chi et al.) |
| SLATree (Chi et al.) | | | |
| Multi-tenant SLOs (Lang et al.) | | | iCBS (Chi et al.) |
| Delphi / Pythia (Elmore et al.) | Hypergraph (Çatalyürek et al.) | | |
| SCOPE (Chaiken et al.) | Bazaar (Jalaparti et al.) | | many traditional methods ... |

# Wish List

## Challenges

| End-to-end cost-aware service<br>(resource provisioning, workload scheduling) | complex interactions |

| Application-defined performance goals<br>(per query deadline, percentile, average latency, max latency ) | arbitrary goals |

| Agnostic to workload semantics | arbitrary workloads |

machine learning: auto modeling and insight

# WiSeDB Advisor

## Offline Learning

- ☐ batch scheduling

## Online Learning

- ☐ online scheduling

- ☐ performance model free

## Data Management Application

| Cost Management | SLA Management |
| Resource Provisioning | Workload Scheduling |

Microsoft SQL Server    MySQL    PostgreSQL    ORACLE

IaaS Provider    VM    VM    VM    VM

# Outline

Motivation

**Offline Learning**

Online Learning

Conclusions

❑System Overview

❑Supervised Learning

❑Adaptive Learning

*WiSeDB: A Learning-based Workload Management Advisor for Cloud Databases, Ryan Marcus, Olga Papaemmanouil, **VLDB 2016***

# WiSeDB – Batch Processing

**Workload & SLO Spec**

2min
SLO: 3min

0.5min
SLO: 1min

**Penalty Function**
$$/sec past deadline

**Data Management Application**

**(Offline) Training**

Model
Generator

# WiSeDB – Batch Processing

**Workload & SLO Spec**

**2min**
SLO: 3min

**0.5min**
SLO: 1min

**SLA Spec**
$$/sec past deadline

- ☐ OLAP on full replicas (no updates)
- ☐ Known queries
- ☐ Performance model

## Data Management Application

**(Offline) Training**

Model Generator

Microsoft® SQL Server®

MySQL®

PostgreSQL

ORACLE®

**IaaS Provider**

VM  VM  VM  VM

# WiSeDB – Batch Processing

## Original SLO

$0.12
**3min**
SLO: 3min

$0.20
**1min**
SLO: 1min

## Stricter SLO

$0.15
**2.5min**
SLO: 3min

$0.13
**0.15min**
SLO: 1min

## Data Management Application

**(Offline) Training**

Model Generator

Strategy Recommendations

# Batch Execution

Runtime Query Batch

Q, …, Q

Q, …, Q

## Resources to rent

- ❑ # VMs/ type

## Query scheduling

- ❑ Query execution order for each recommended VM

## Data Management Application

**(Offline) Training**

Model Generator

Strategy Recommendations

**(Online) Resource & Workload Management**

Strategy Generator

IaaS Provider

Q Q Q Q Q

VM VM VM VM

**ASSUMPTIONS**

- ❑ OLAP on full replicas (no updates)
- ❑ Known query types
- ❑ Performance prediction model

# Supervised Learning

Model Generator

| identify classes | classes == actions → | ❑ dispatch a query to a VM<br>❑ provision new VM |

| create training data | context of actions → | ❑ identify best decisions<br>❑ extract cost-related features |

| generate classifier | decision tree → | ❑ describe (context, action)<br>❑ interpretable: offers insight |

*"To be the best, learn from the best"* ( *D. LaCroix*)

Model Generator

## Offline Learning

**identify best decisions**

1. Generate small workload
2. Build decision graph
   - query assignment
   - VM provisioning
3. Find optimal (minimum cost) solution (path)
4. Extract context of optimal decisions

**generate model**

1. Repeat for many sample workloads
2. Build a training set of (feature, action)
3. Train a classifier

## Runtime Scheduling

**apply model**

- Use classifier for
  - batch scheduling
  - online scheduling
  - performance vs cost exploration

# Decision Graph



Model Generator

A

9    6

C    B

Monetary Cost
- ❑ Resource usage ($$/time)
    - ❑ time = VM start up + query execution

- ❑ Violation fees
    - ❑ Penalty function (black box)

# Search for Optimal

A* search (best-first)
for optimal

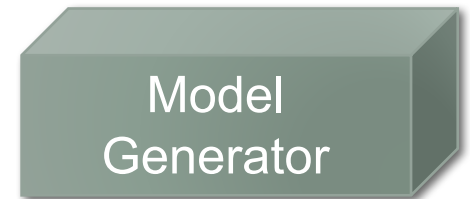Model Generator

# Search for Optimal



A* search (best-first)
   for optimal

**Model Generator**

9    6    6    50    9    10

**Graph-based Approach Pros**
- ☐ Step-by-step decisions
- ☐ Graph reduction techniques
- ☐ Fast search for optimal

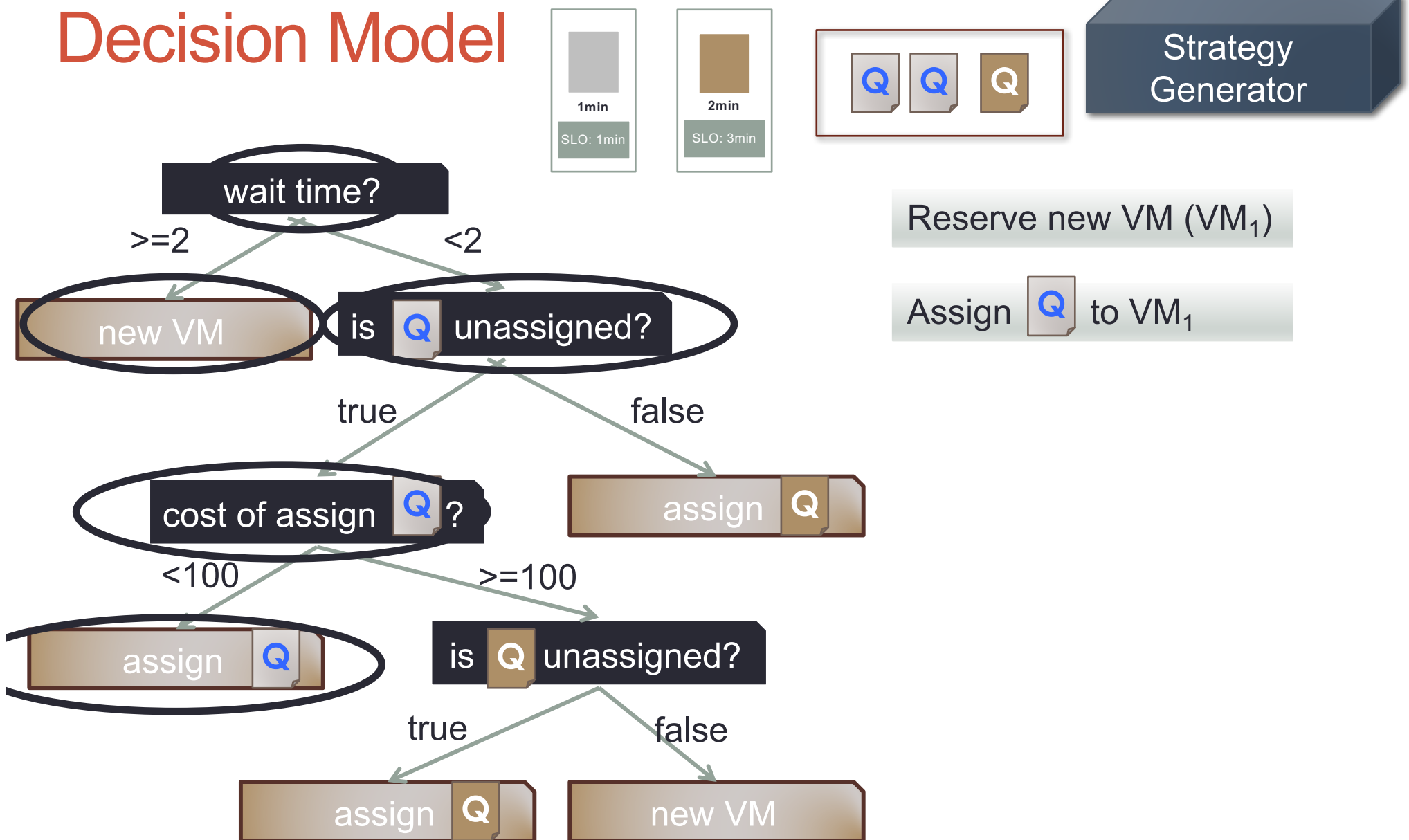# Feature Extraction
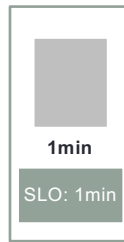


Model Generator

**Decision:** Assign Q to VM

**Features:**
- unassigned Q : true
- unassigned Q : false
- cost of assigning Q : $2
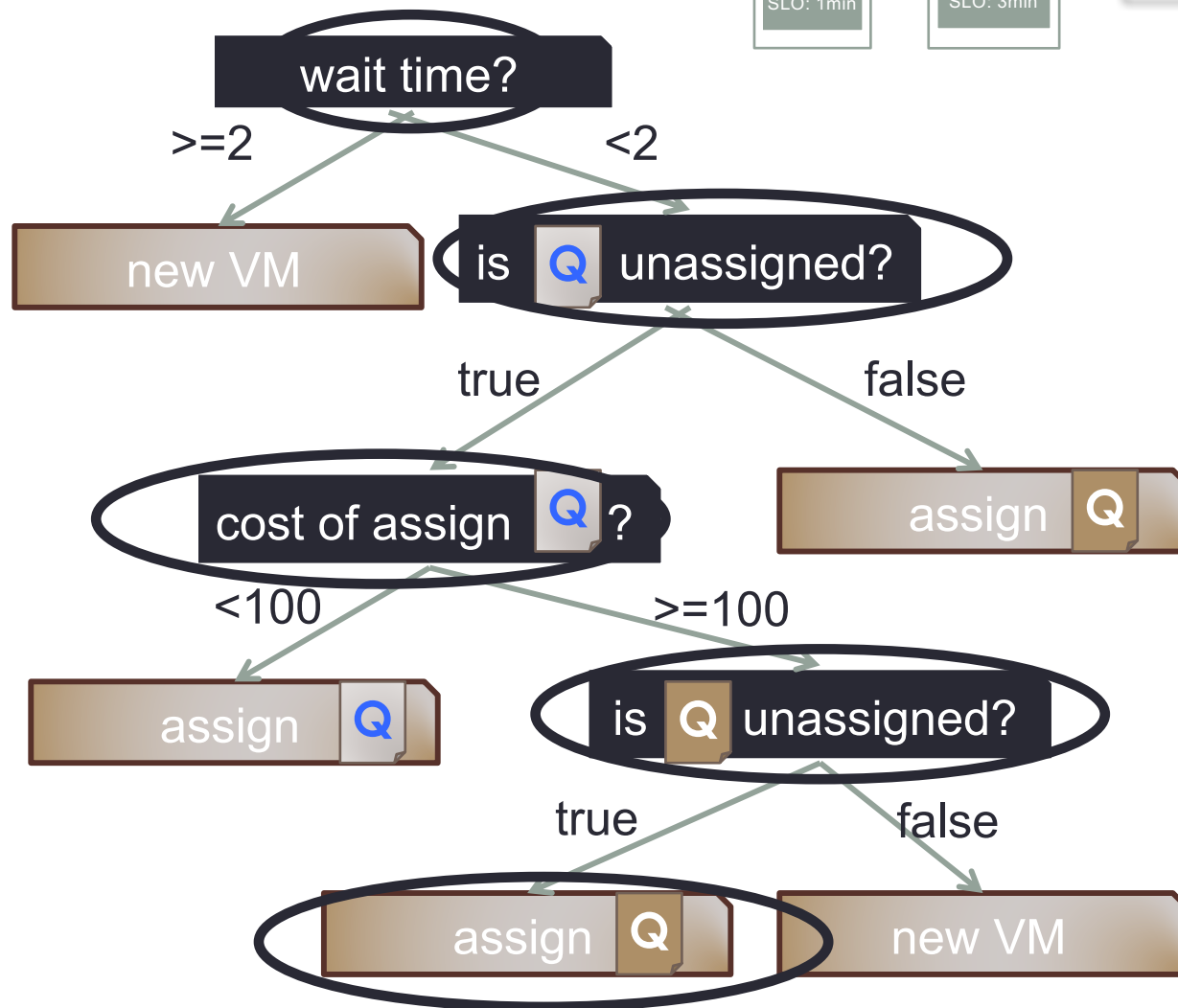- wait time on VM: 1min
- % of Q in VM: 50%
- % of Q in VM: 50%

9

6

**Agnostic to**
- Query semantics
- Performance goal (SLO)
- Workload size

# Decision Model

# Decision Model

# Decision Model



Strategy Generator

- Reserve new VM (VM$_1$)
- Assign Q to VM$_1$
- Assign Q to VM$_1$
- Reserve new VM (VM$_2$)
- Assign Q to VM$_2$

wait time?

>=2 → new VM

<2 → is Q unassigned?

true → cost of assign Q ?

false → assign Q

<100 → assign Q

>=100 → is Q unassigned?

true → assign Q

false → new VM
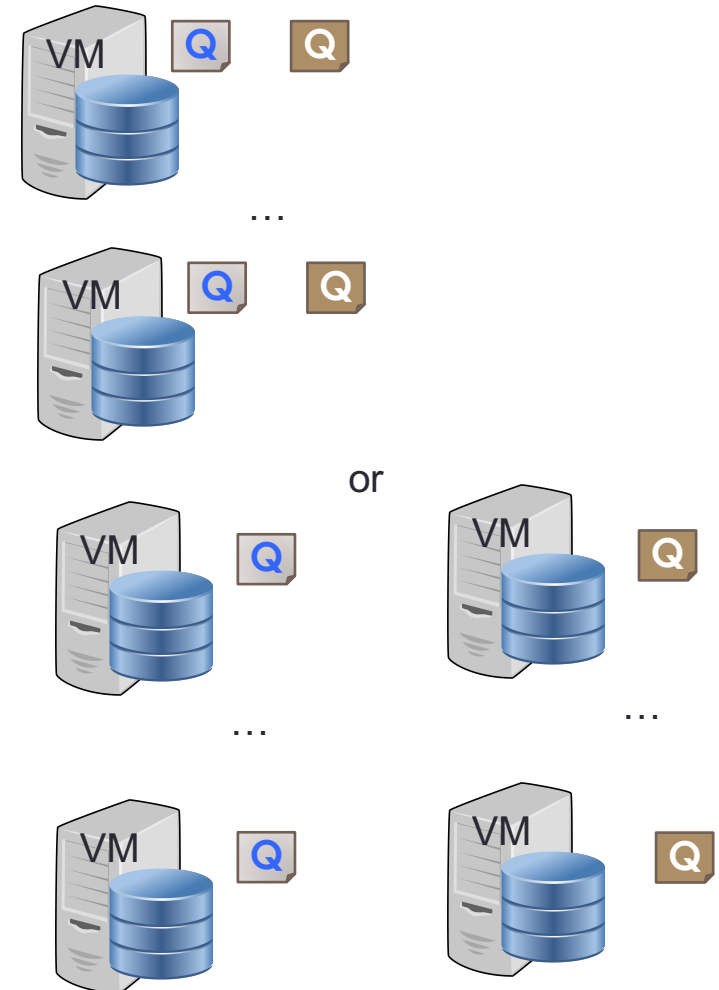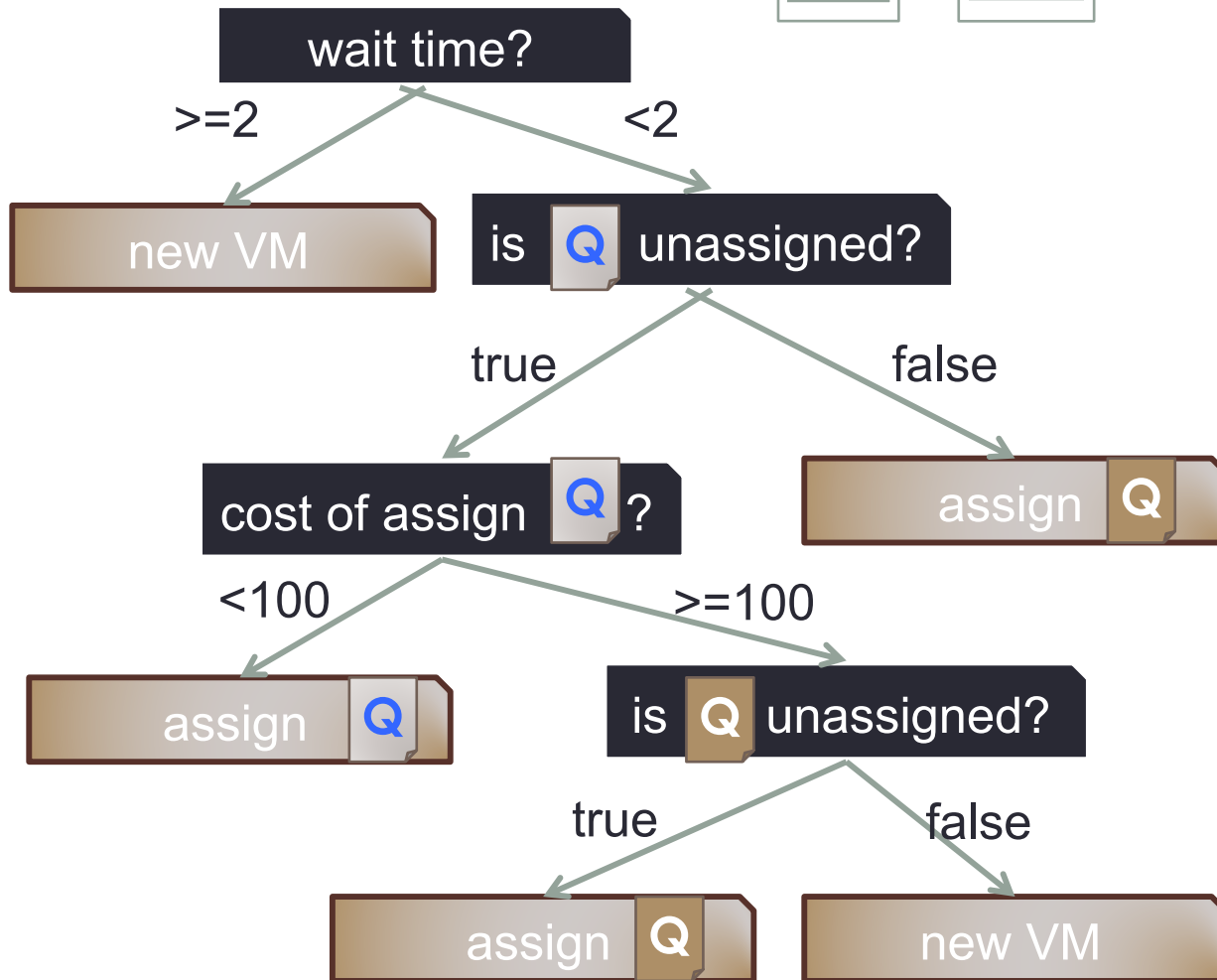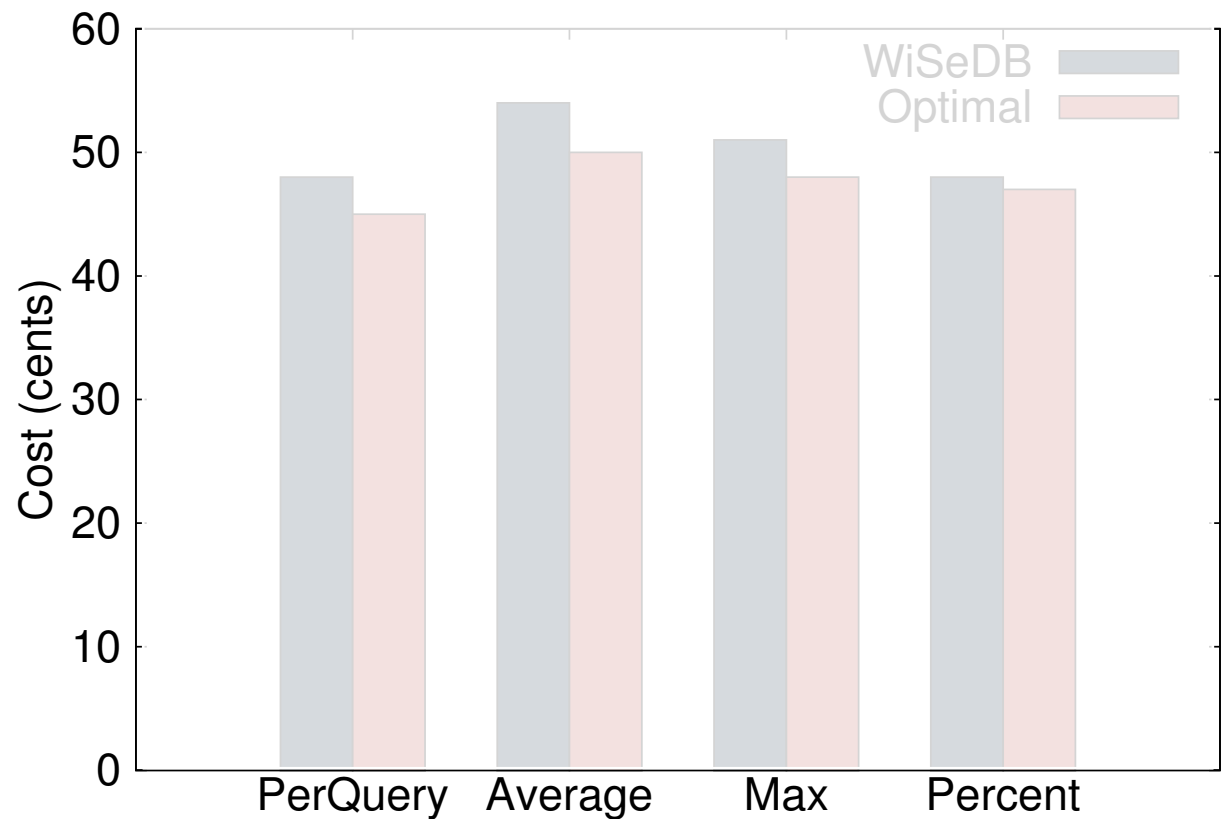
# Strategy

# Experimental Setup

**Training Data**

3000 samples

10 TPC-H templates
18 queries/sample



query execution time <=x secs
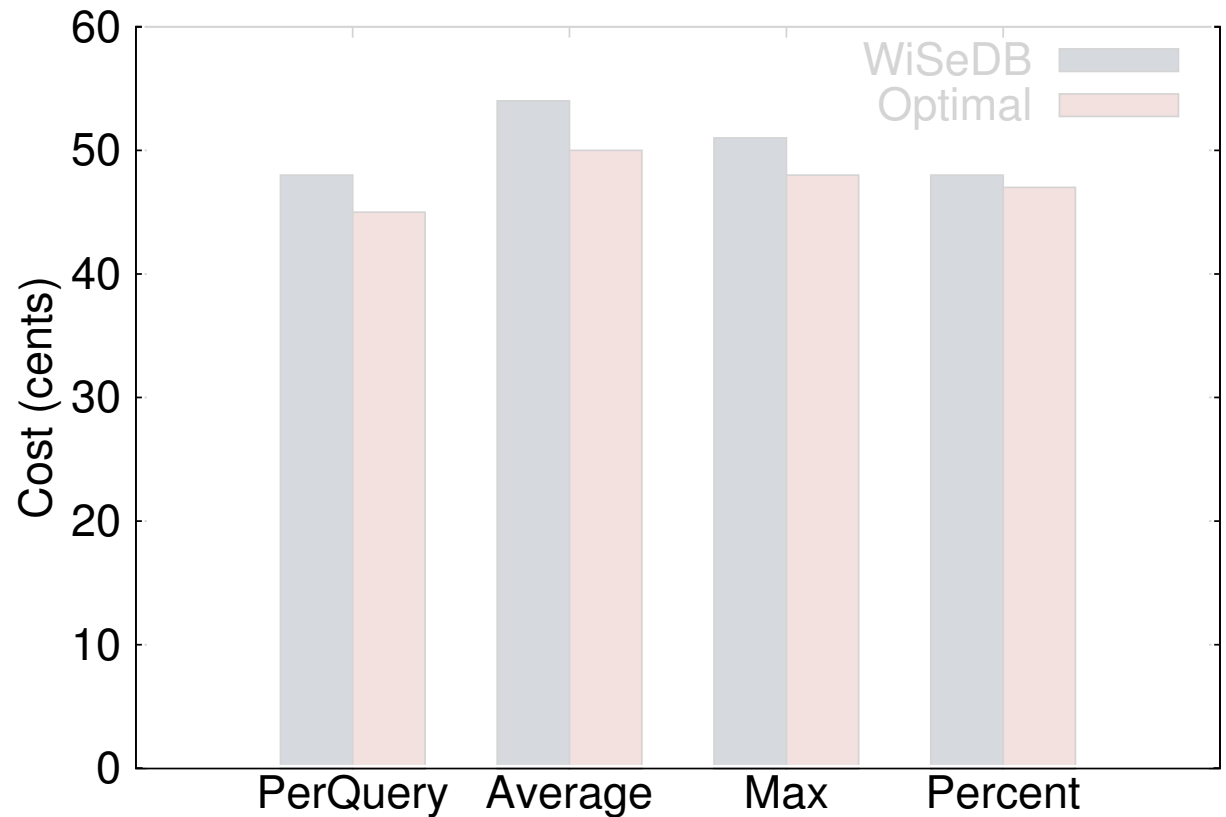(same deadline per template)

# Experimental Setup

**Training Data**

3000 samples

10 TPC-H templates
18 queries/sample



average latency of the
workload <= x secs

# Experimental Setup

**Training Data**

3000 samples

10 TPC-H templates
18 queries/sample



max latency <=x secs
(longest query in the workload )

# Experimental Setup

**Training Data**

3000 samples

10 TPC-H templates
18 queries/sample



execution time of 90% of queries
in the workload <= x secs

# Experimental Setup

**Training Data**

3000 samples

10 TPC-H templates
18 queries/sample

**Testing Data**

10 TPC-H templates

varied queries/workload

# Experimental Setup

**Training Data**

3000 samples

10 TPC-H templates
18 queries/sample

**Testing Data**

10 TPC-H templates

varied queries/workload

*cost: resource utilization+ penalties*

**AWS Cloud**

fees penalty $0.01/sec of violation

# Effectiveness (small workloads)

**Training Data**

3000 samples
10 TPC-H templates
18 queries/sample

**Testing Data**

10 TPC-H templates
**30 queries/workload**
*Optimal*: **Brute force**



**WiSeDB models are within 8% of the minimum cost solution**

# Effectiveness (large workloads)

**Training Data**

3000 samples

10 TPC-H templates

18 queries/sample

**Testing Data**

10 TPC-H templates

**5000 queries/workload**

**One heuristic cannot fit all**

**WiSeDB learns the right heuristic**



Best: shortest query first

Best: longest query first

Best: top-90% shortest then 10% longest queries

# Training Overhead

**Training Data**

3000 samples

10 TPC-H templates

18 queries/sample





**Offline learning overhead
20sec – 120 sec**

# Beyond Batch Scheduling

- Efficient performance vs cost trade off exploration
  - Recommend strategies with stricter/looser performance goals
  - Reuse original training set to generate **quickly** alternative models
    - Best-first heuristic reduces search time (dominant training factor)
  - Training overhead improvement by **96-98%**

- Online scheduling (query at a time)
  - Challenge: arrival times are unknown and hence not modeled
  - Generate a new model upon arrival of new query: too expensive
  - Optimization 1: Adapt previous model to reduce training overhead
  - Optimization 2: Reuse past models, when feasible

# Offline Learning

## Advantages

❑ Provides insight on complex decisions

❑ Learns custom strategies per application

❑ Explores performance vs cost trade-offs

### Data Management Application

**(Offline) Training**

Model Generator

Strategy Recommendations

**(Online) Resource & Workload Management**

Strategy Generator

IaaS Provider

VM  VM  VM  VM

# Offline Learning



**Limitations**

- ❏ Static decision models
- ❏ Batch scheduling
- ❏ Performance model

**Data Management Application**

**(Offline) Training**

Model Generator

Strategy Recommendations

**(Online) Resource & Workload Management**

Strategy Generator

IaaS Provider

VM  VM  VM  VM

# Outline

Motivation

Offline Learning

Online Learning

Conclusions

☐ Explicit vs Implicit Modeling

☐ Reinforcement Learning

*Releasing Cloud Databases from the Chains of Predictions Models.*
*Ryan Marcus, Olga Papaemmanouil,* **CIDR 2017**

# (Explicit) Performance Prediction

- ❑ DBMS-related challenges
    - ❑ isolated vs. concurrent query execution
    - ❑ low accuracy for new query types ("templates")
    - ❑ extensive off-line training
    - ❑ <span style="color:red">state-of-the-art: 15-20% prediction error*</span>

- ❑ Cloud-related challenges
    - ❑ "noisy neighbors"
    - ❑ numerous resource configurations
    - ❑ predictions errors accumulation

*Contender: A Resource Modeling Approach for Concurrent Query Performance Prediction*,
Jenny Duggan, Olga Papaemmanouil, Ugur Cetintemel, Eli Upfal, **EDBT 2015**

*Performance Prediction for Concurrent Database Workloads*,
Jennie Rogers, Ugur Cetintemel, Olga Papaemmanouil, Eli Upfal, **SIGMOD 2011**

# WiSeDB: Implicit Performance Modeling
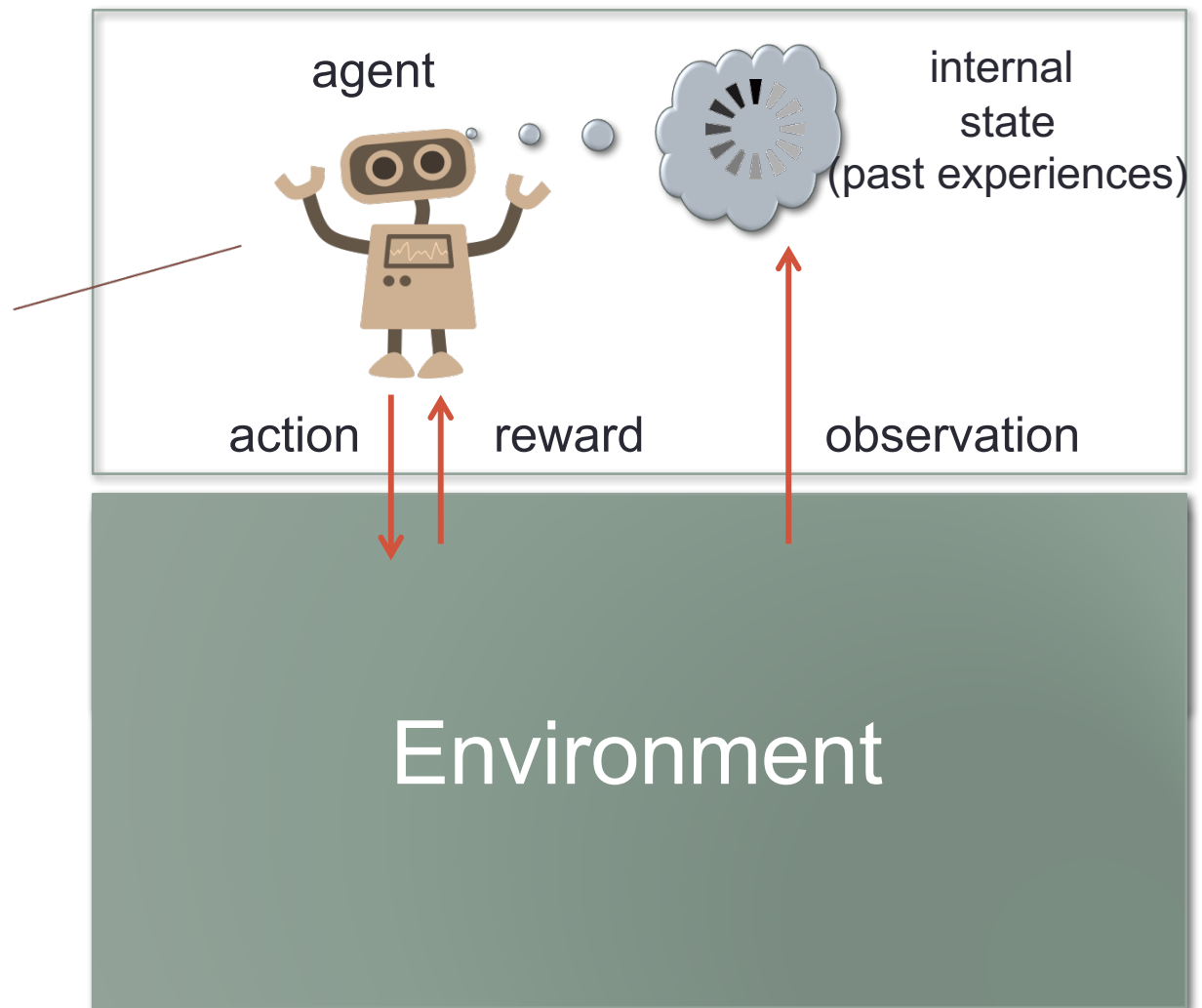
- ❑ Explicit performance models are NOT necessary for:
  - ❑ monetary cost management
  - ❑ resource & workload management
  - ❑ offer performance SLA and **keep penalties low**

- ❑ Implicitly model query latency
  - ❑ predict *monetary cost* ( *& violation penalties)*
- ❑ Online training for dynamic environments
  - ❑ Automatic scaling & workload distribution

**Wish List #2**

# Reinforcement Learning

- ❑ Continuous learning

- ❑ Explicit reward modeling

- ❑ Action selection
  - ❑ maximize reward

agent

internal state (past experiences)
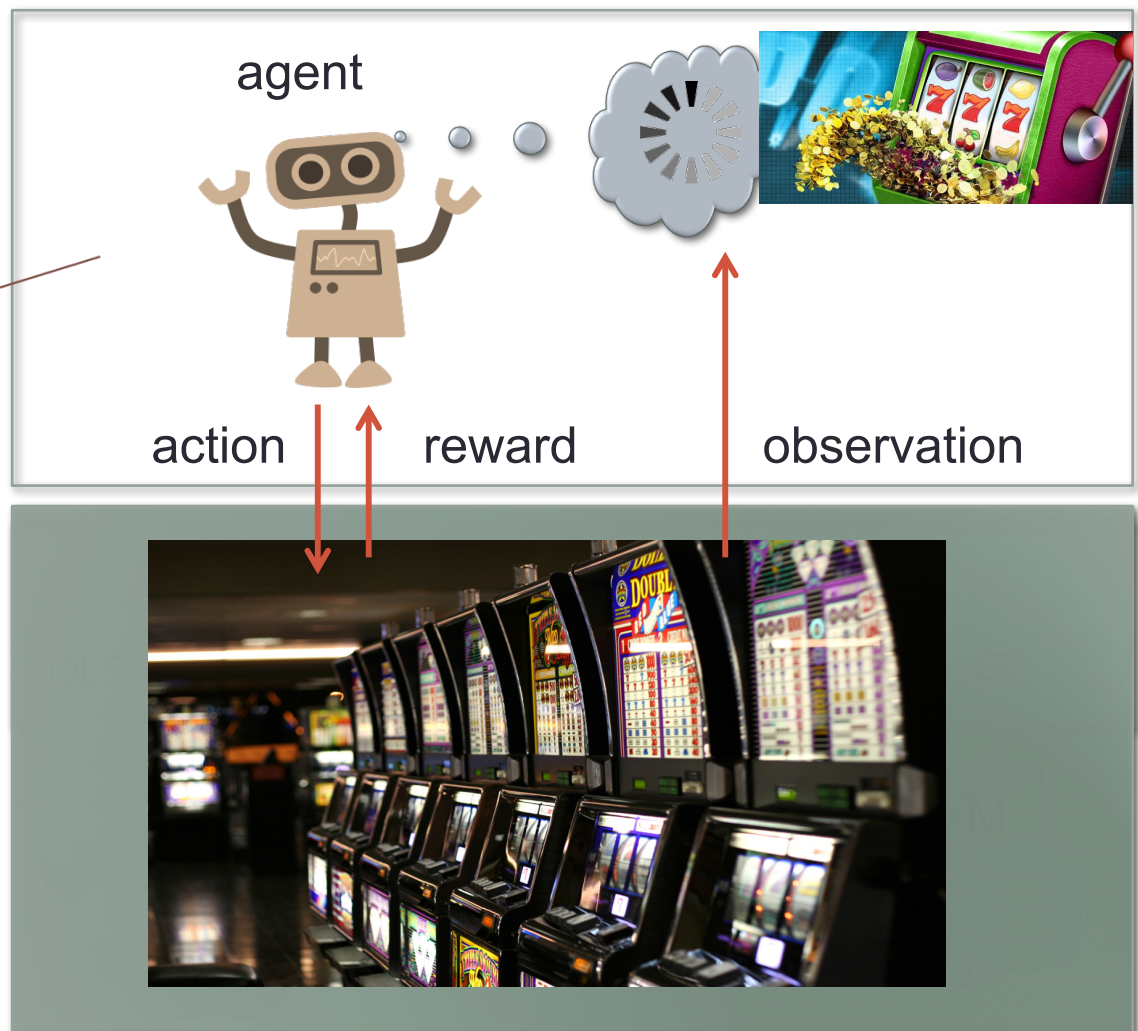
action    reward    observation

Environment

# CMABs
## (Contextual Multi-Armed Bandits)

**Contextual Multi-Armed Bandit Problem**

Armed Bandit = Slot Machine

*Which slot machine to play (**action**) so that you walk out with the most $$$ (**reward**)?*

agent

action     reward     observation

# CMABs in WiSeDB
(Contextual Multi-Armed Bandits)
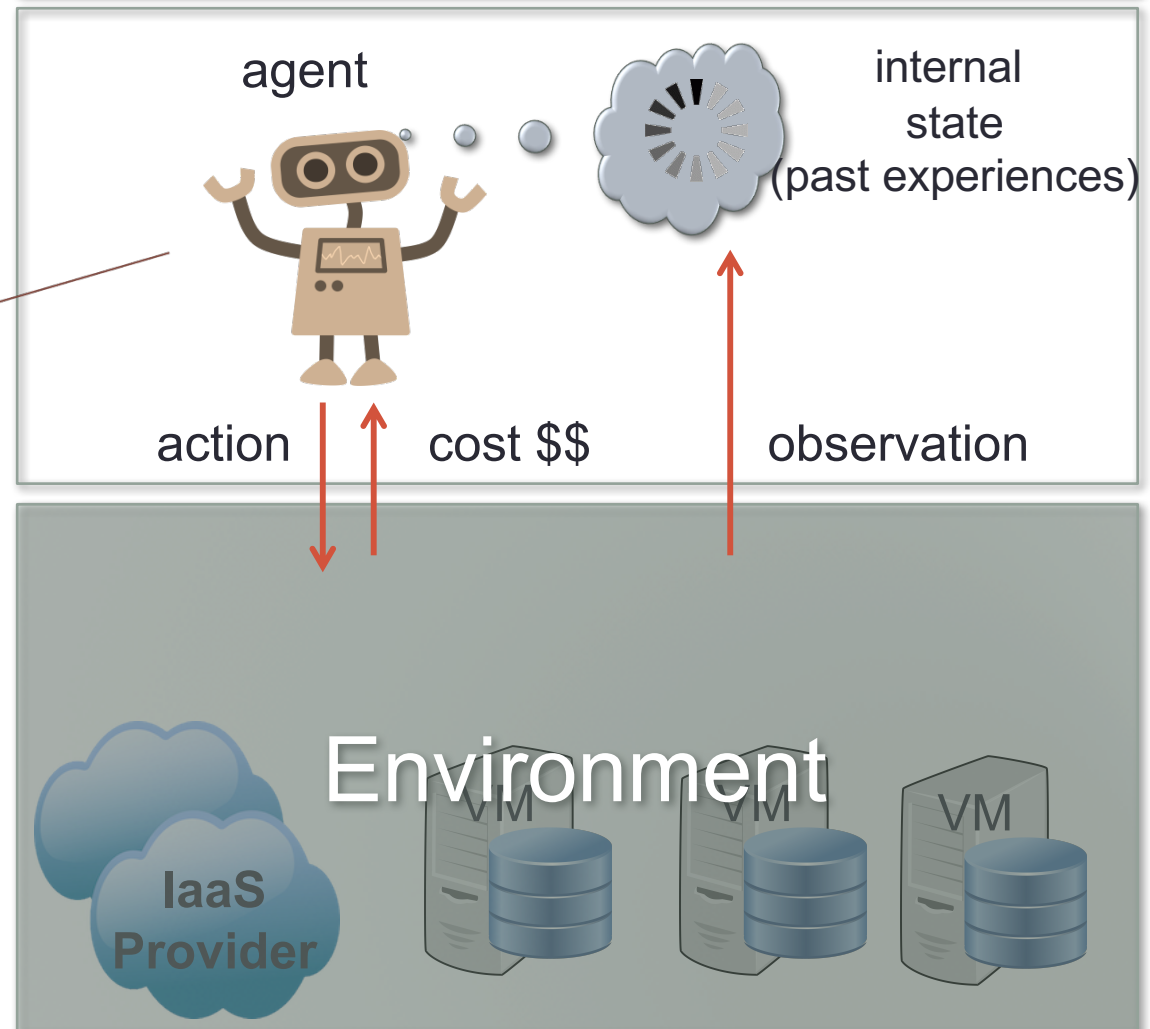


**Contextual Multi-Armed Bandit Problem**

Slot Machine = Virtual Machine

*Which machine to use (new/old) (**action**) so that you execute the incoming query with minimum cost $$ (**cost**)?*

Data Management Application

agent

internal state (past experiences)

action      cost $$      observation

Environment

IaaS Provider

VM      VM      VM

# CMABs in WiSeDB
(Contextual Multi-Armed Bandits)

Q Q Q Q

## Action (per VM)
- ❑ Accept
- ❑ Pass to next /new VM
- ❑ Down one VM tier

## Reward
- ❑ $$ cost: processing & SLA violation penalties

## Observation
- ❑ environment context (query, VM)
- ❑ action
- ❑ $$ cost

## Data Management Application

SLA

internal state (past experiences)

action        cost $$        observation

VM Tier 1

IaaS Provider

VM Tier 2

# CMABs in WiSeDB
(Contextual Multi-Armed Bandits)

Q Q Q

## Action (per VM)
- ❑ Accept
- ❑ Pass to next /new VM
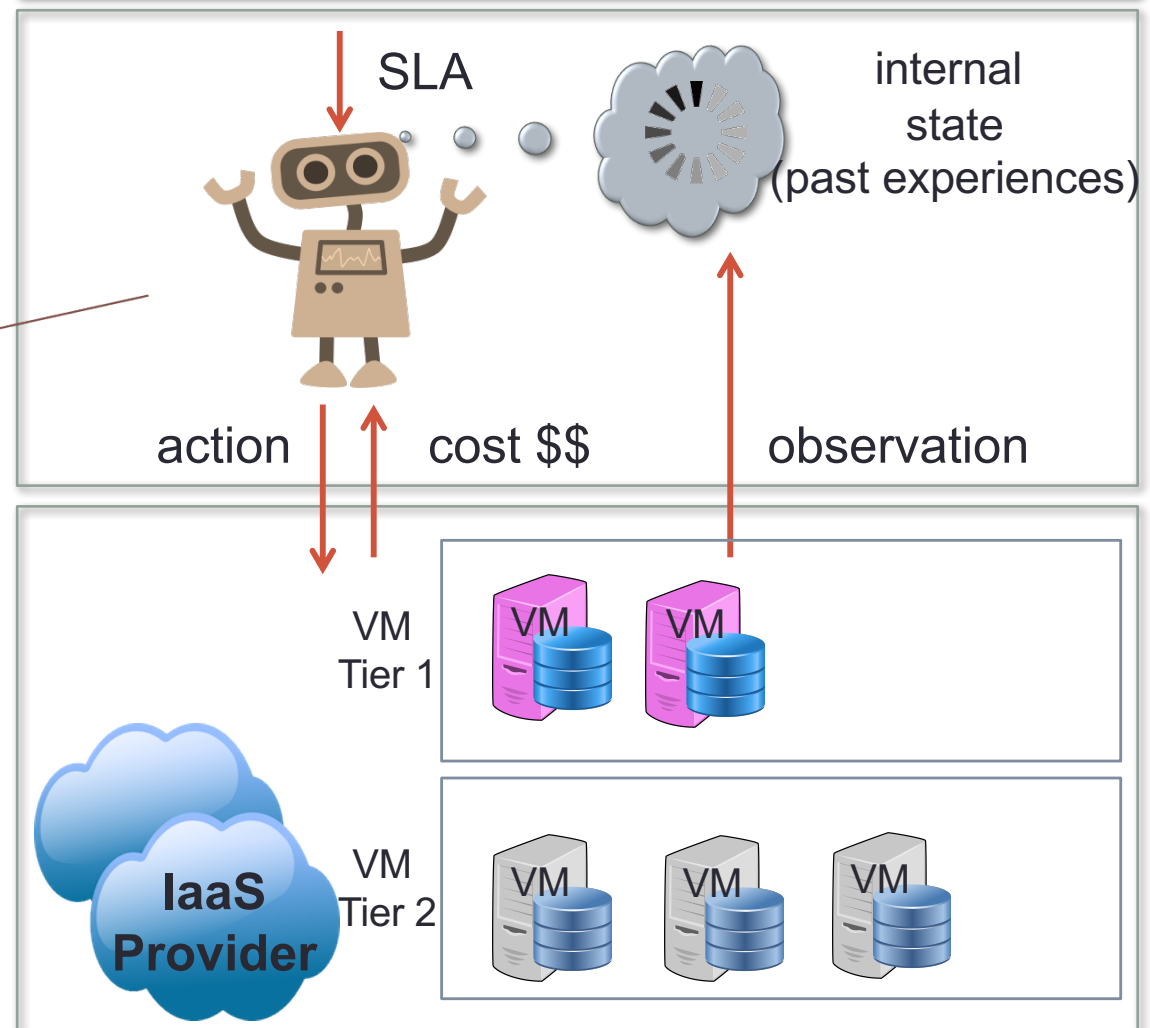- ❑ Down one VM type

## Reward
- ❑ $$ cost: processing & SLA violation penalties

## Observation
- ❑ environment context (query, VM)
- ❑ action
- ❑ $$ cost

## Data Management Application

Q → SLA

internal state (past experiences)

action      cost $$      observation

VM Tier 1      **pass**   **down**

IaaS Provider      VM Tier 2   **accept**   VM   VM

# CMABs in WiSeDB
(Contextual Multi-Armed Bandits)

Q Q Q

## Action (per VM)
- ☐ Accept
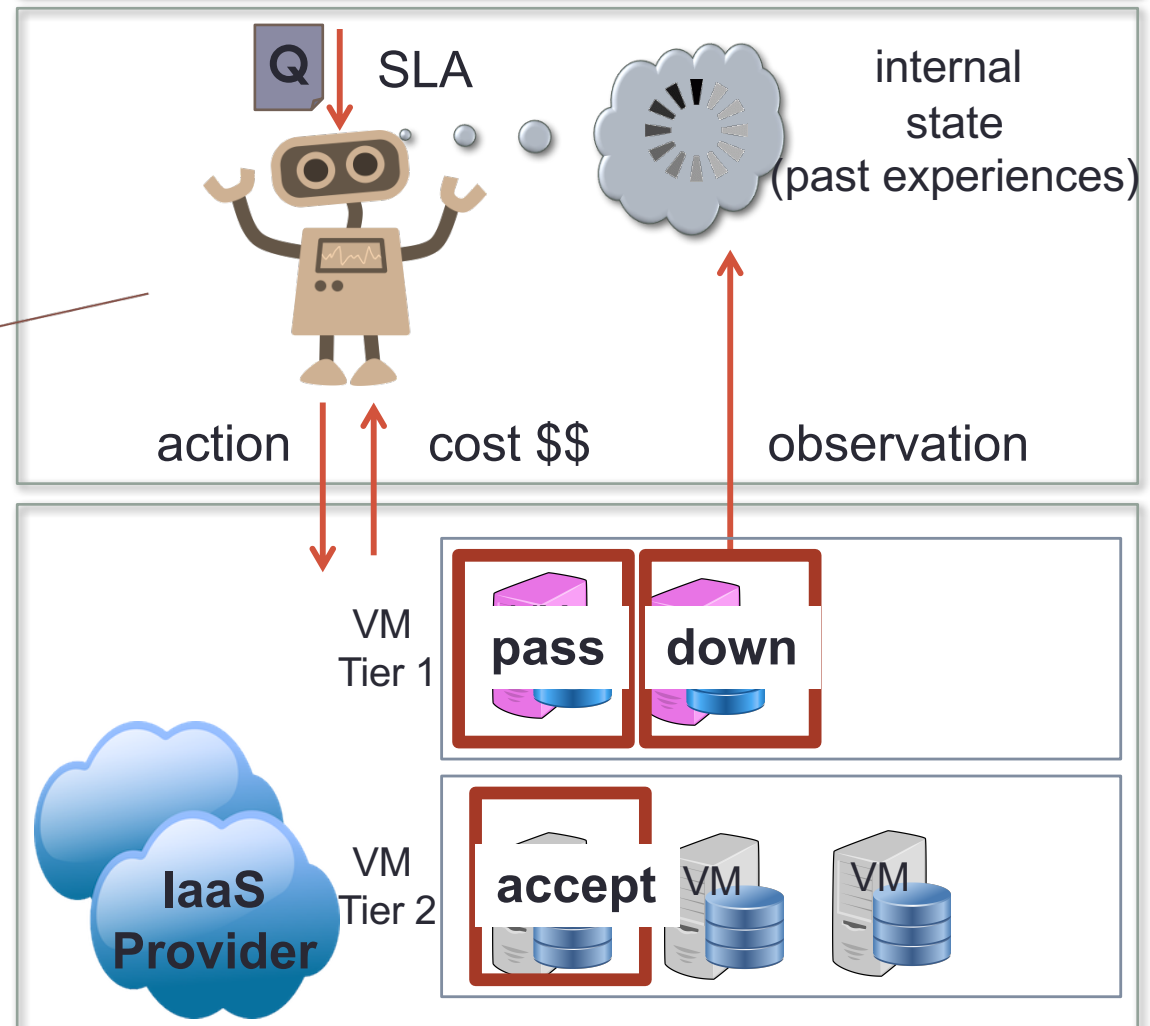- ☐ Pass to next /new VM
- ☐ Down one VM type

## Reward
- ☐ $$ cost: processing & SLA violation penalties

## Observation
- ☐ environment context (query, VM)
- ☐ action
- ☐ $$ cost

## Data Management Application

Q   SLA

internal state (past experiences)

action        cost $$        observation

VM Tier 1

VM VM

IaaS Provider

VM Tier 2

VM VM VM

# CMABs in WiSeDB
(Contextual Multi-Armed Bandits)

Q  Q  Q

## Action (per VM)
- ❑ Accept
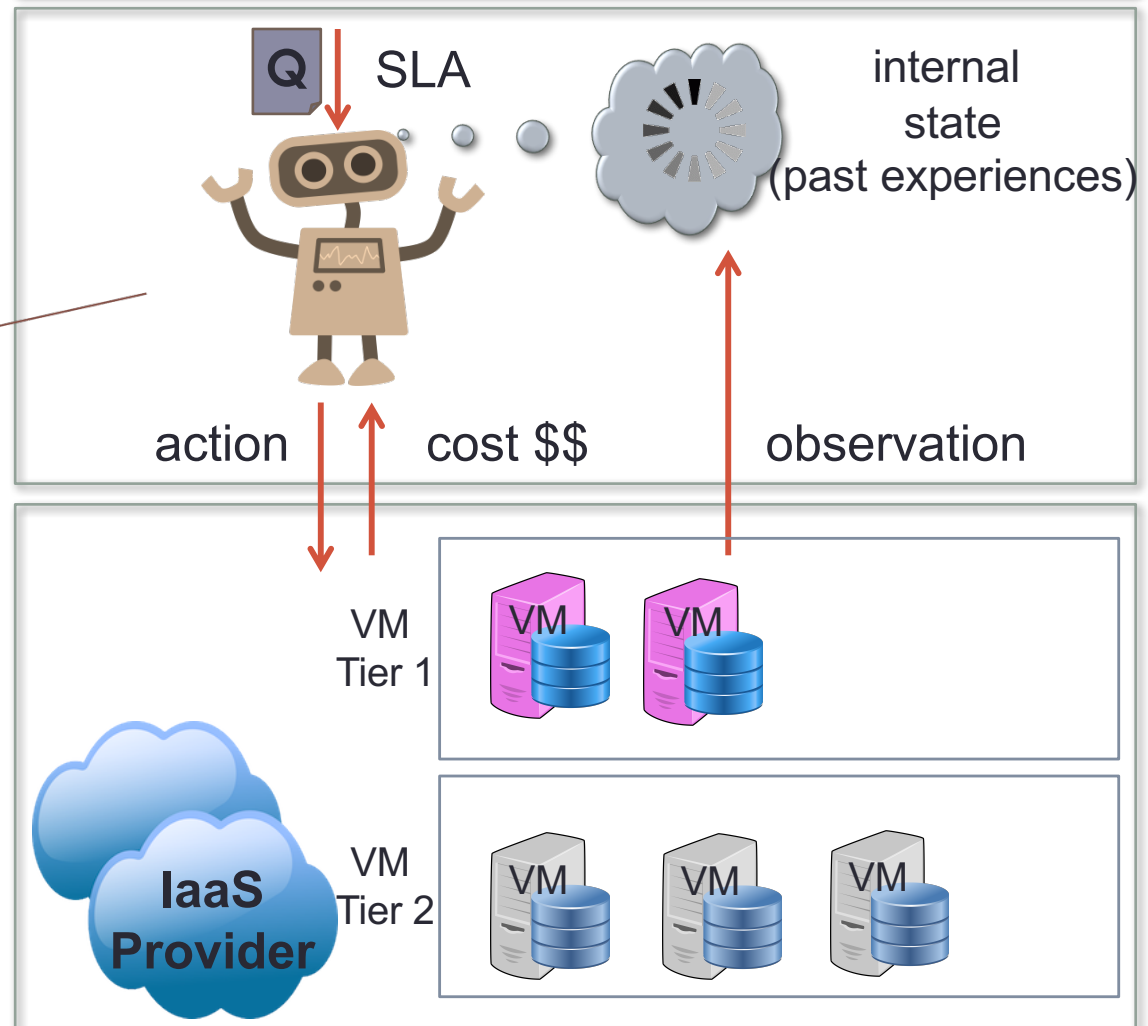- ❑ Pass to next /new VM
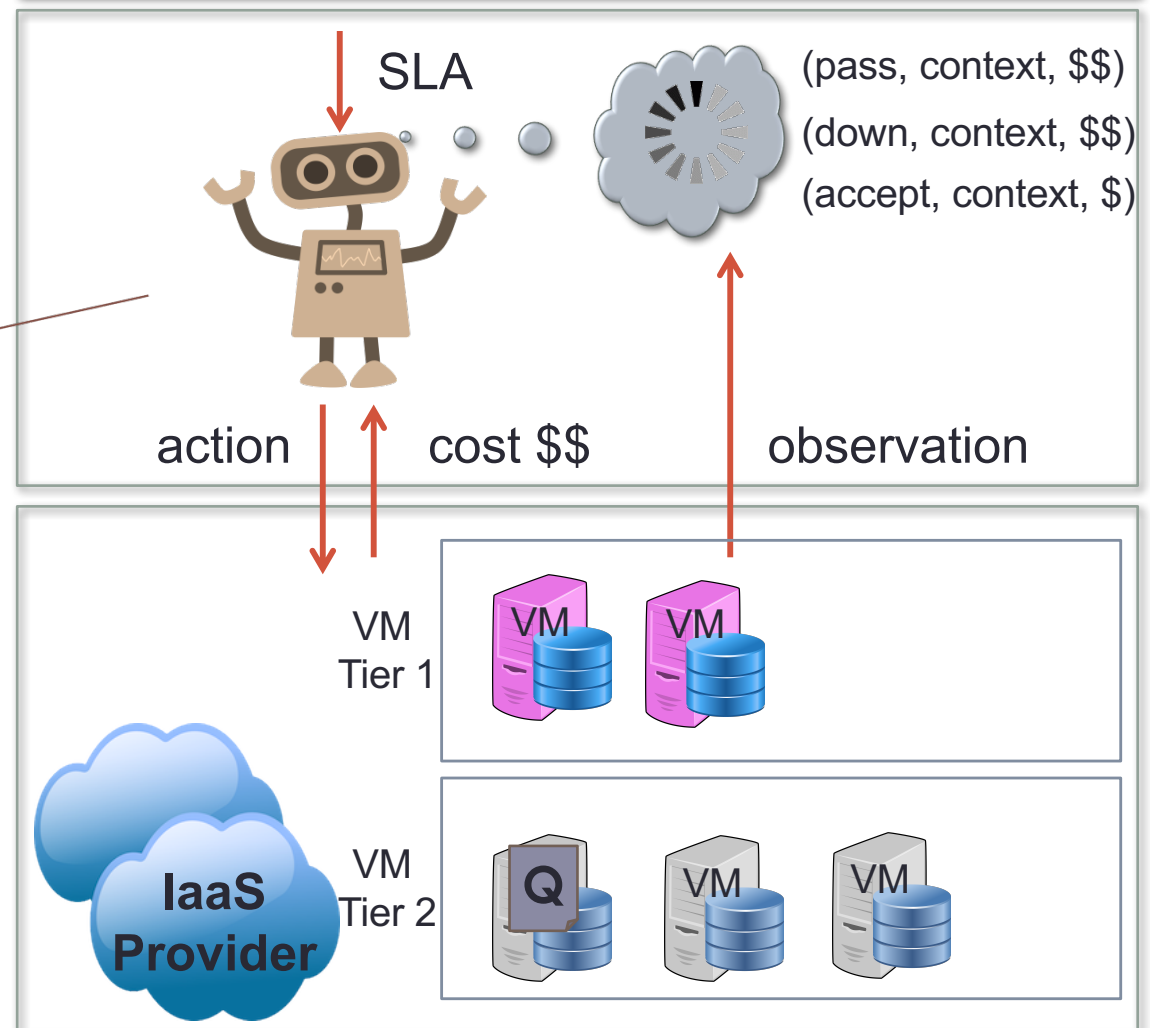- ❑ Down one VM type

## Reward
- ❑ $$ cost: processing & SLA violation penalties

## Observation
- ❑ environment context (query, VM)
- ❑ action
- ❑ $$ cost

## Data Management Application

SLA

(pass, context, $$)

(down, context, $$)

(accept, context, $)

action     cost $$     observation

VM Tier 1

VM Tier 2

IaaS Provider

VM   VM

Q   VM   VM

# CMABs in WiSeDB
(Contextual Multi-Armed Bandits)

Q  Q

## Action (per VM)
- ❑ Accept
- ❑ Pass to next /new VM
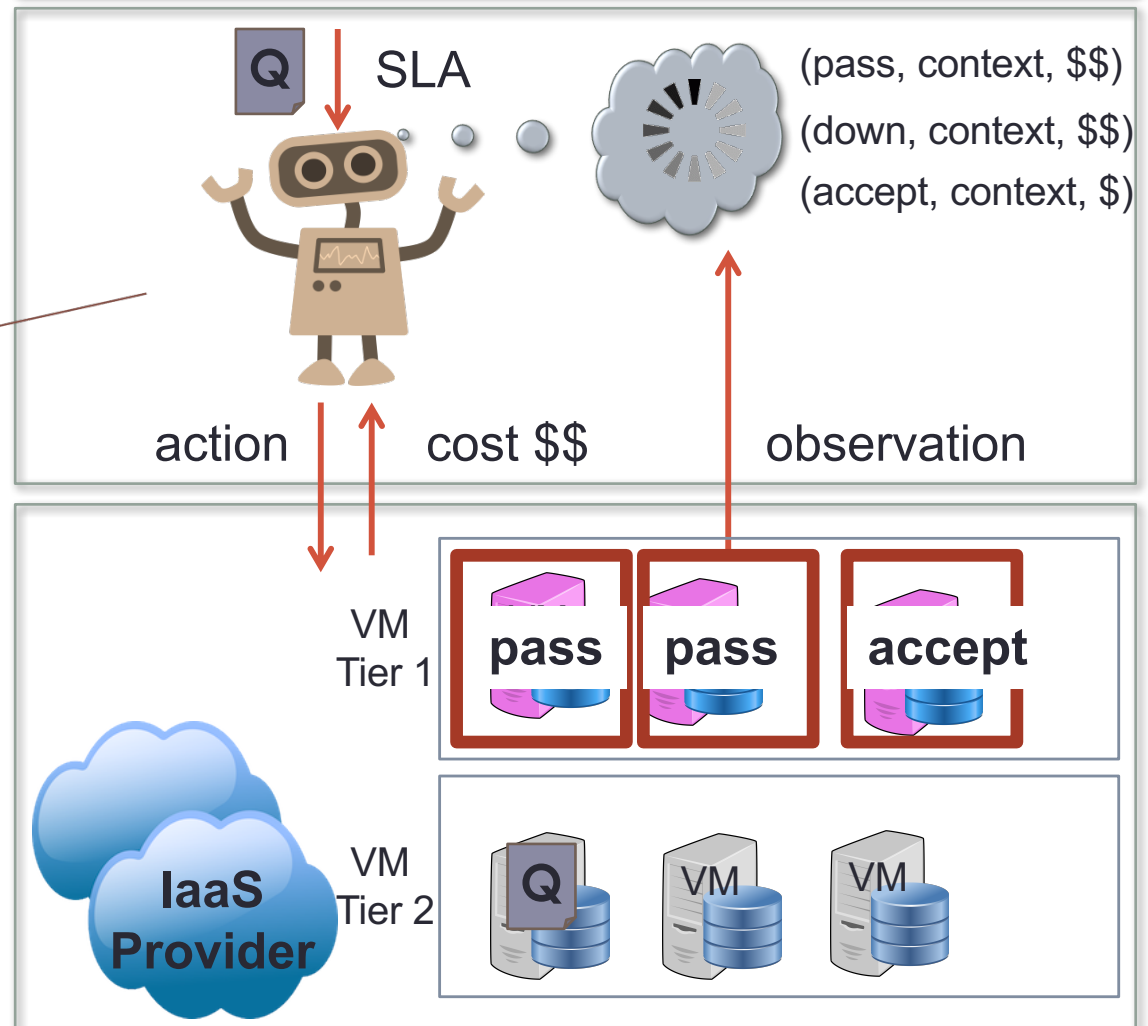- ❑ Down one VM type

## Reward
- ❑ $$ cost: processing & SLA violation penalties

## Observation
- ❑ environment context (query, VM)
- ❑ action
- ❑ $$ cost

## Data Management Application

Q   SLA

(pass, context, $$)
(down, context, $$)
(accept, context, $)

action        cost $$        observation

VM Tier 1     **pass**   **pass**   **accept**

IaaS Provider

VM Tier 2     Q   VM   VM

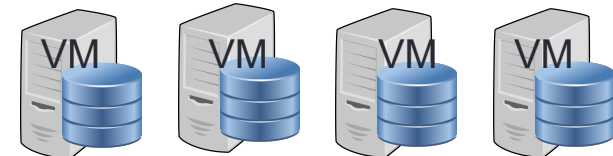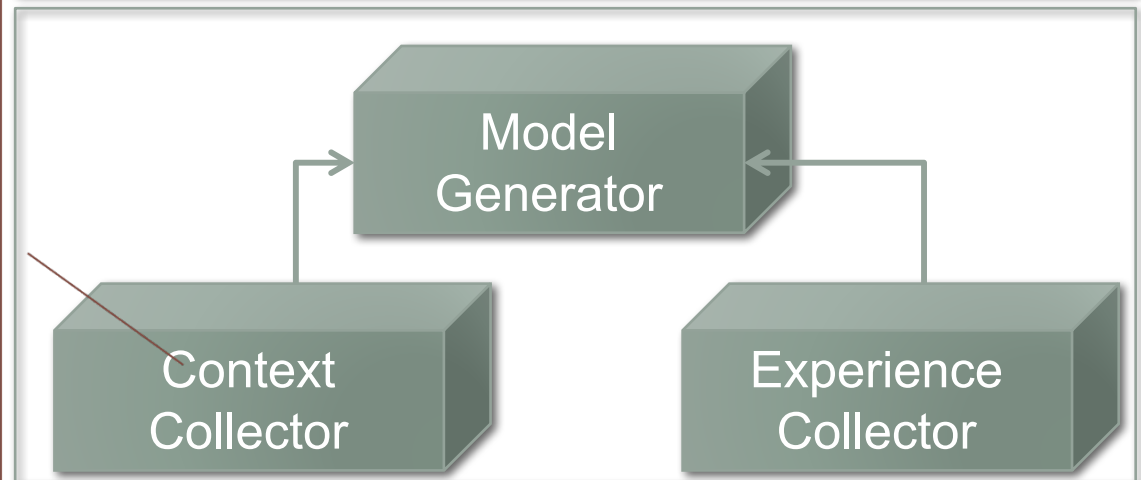# Online Learning



## Context Features

- **VM context**
  - memory, I/O rate
  - #queries in queue

- **Query context**
  - tables used by current query
  - tables used by old query
  - # table scans
  - # joins
  - # spill joins
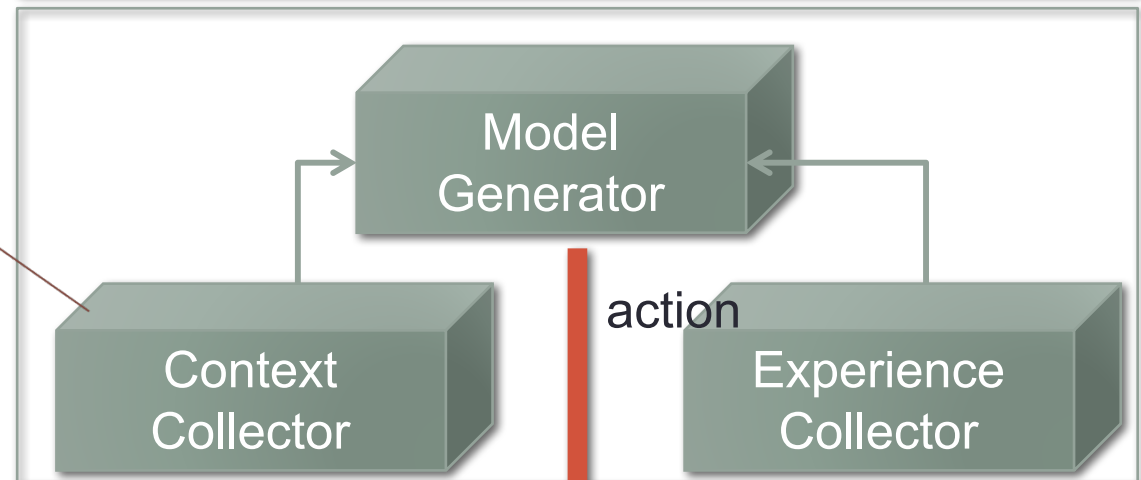  - cache reads in the plan

**Data Management Application**

Model Generator

Context Collector

Experience Collector

IaaS Provider

VM    VM    VM    VM

# Online Learning

## Action Selection

- **Explore** opportunities
  - gather information

- **Exploit** "safe" actions
  - make best decision given current information

**Data Management Application**

Model Generator

Context Collector

action

Experience Collector

IaaS Provider

VM  VM  VM  VM

# Probabilistic Action Selection

- Select action according to probability of being the best
- Past observations (action, context, cost) $D = \{(x_i, a_i, c_i)\}$
  - modeled by likelihood function over cost $c$ : $P(c \mid \alpha, x, \theta)$
  - $\theta$: **parameters of likelihood function: splits of a regression tree**
    - *if (# joins in the query =1) and (queries in the queue =3 ) => cost = $$*

- Posterior distribution of $\theta$ (Bayes rule)        **perfect decision tree is unknown**

$$P(\theta \mid D) \propto \prod P(c_i \mid a_i, x_i, \theta) P(\theta)$$

  - $P(\theta)$: prior distribution of parameters $\theta$

- Choose action $\alpha'$ to minimize cost for perfect model $\theta*$

$$\min_{a'} E(c \mid a', x, \theta^*)]$$

# Probabilistic Action Selection

❑ Exploitation: pick action based on mean of posterior $P(\theta|D)$

$$\min_{a'} \mathrm{E}(c \mid a', x) = \int \mathrm{E}(c \mid a', x, \theta) P(\theta \mid D) \, d\theta$$

❑ Exploration: pick a random action

❑ Thompson Sampling: balance exploration/exploitation

**Select <u>random</u> action according to probability that it is the best**

# WiSeDB Action Selection

context $x_i$

$$\underset{\alpha_i}{\text{argmin}}\, \text{E}(c \mid x_i, a_i, \theta_i)]$$

$$D = D \cup (x_i, a_i, c_i)$$

**Sample** random parameter $\theta_i$ according to $P(\theta \mid D)$

Select best action $\alpha_i$ according to $\theta^t$

Observe cost $c_i$ update model

**Select a <u>random training set,</u> <u>generate the regression tree</u> and pick <u>best action</u> according to it**

**Update the experience set**

**Create new model**
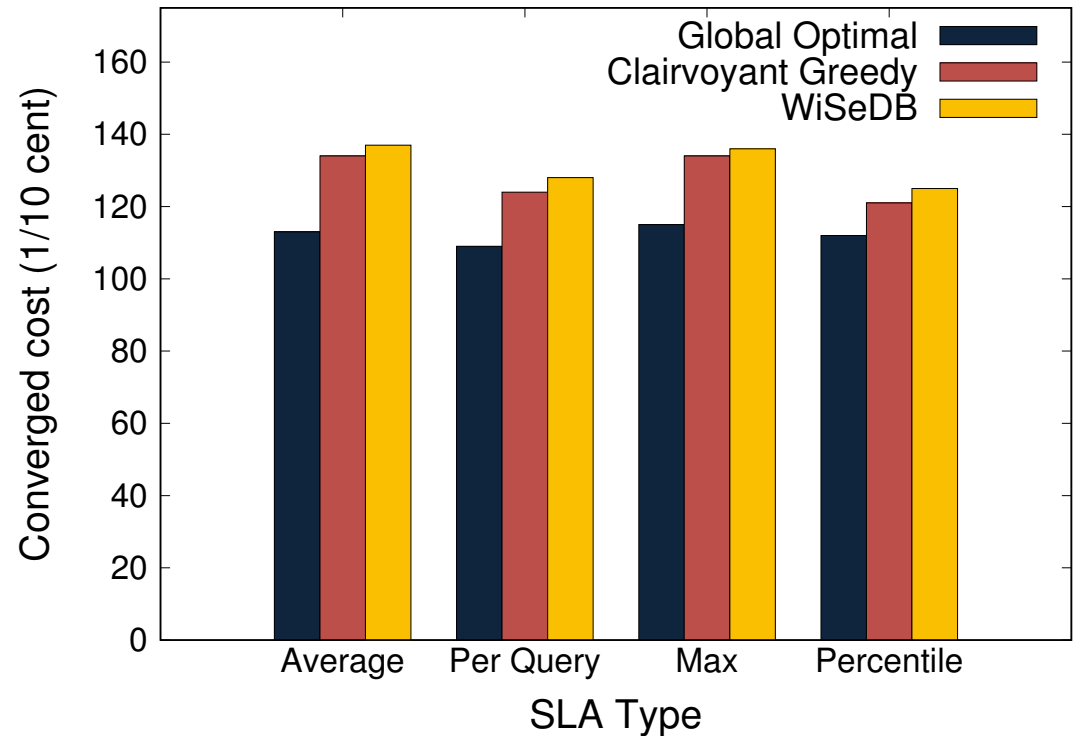
# Effectiveness

**Training Data**

30 query sequence

22 TPC-H templates
repeat until convergence

*Optimal*: **brute force (NP-hard)**

*Clairvoyant*: **perfect cost model**

**Amazon AWS**

t2.large, t2.medium, t2.small



**WiSeDB models can perform at the same cost as a perfect cost model**
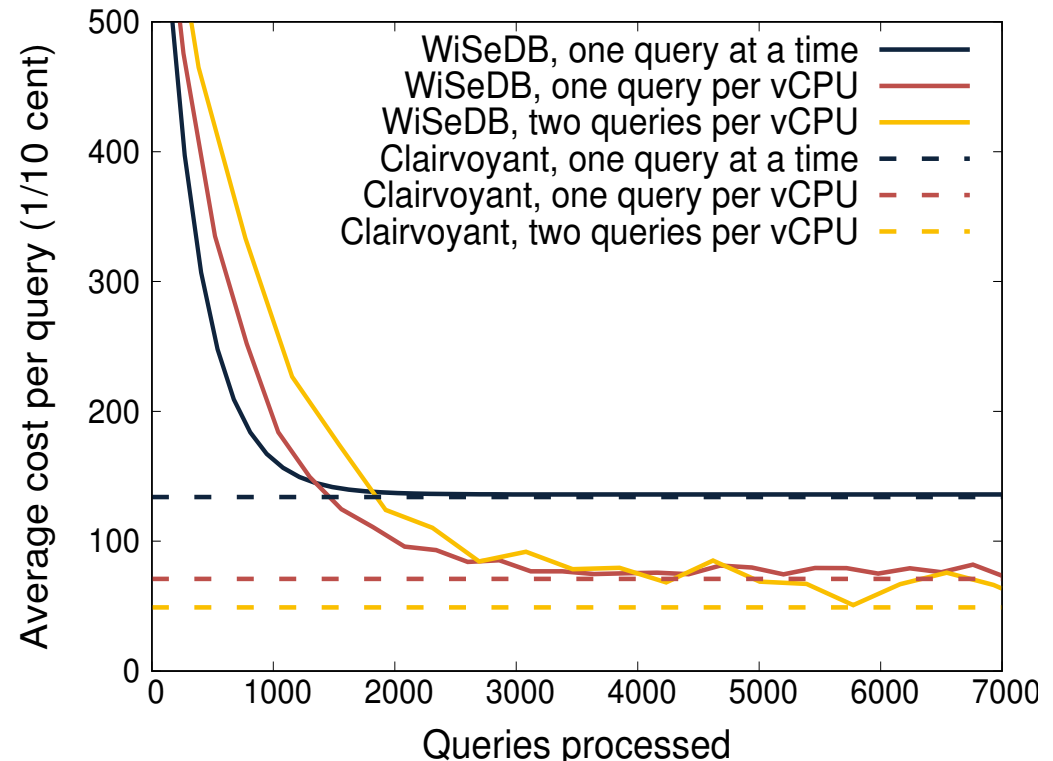
# Effectiveness (concurrency)

**Training Data**

22 TPC-H templates
900 queries/hour
Poison distribution

*Clairvoyant*: perfect cost model

*One query/vCPU*: 1-2 queries

*Two queries/vCPU*: 2-4 queries



**WiSeDB models handles concurrency
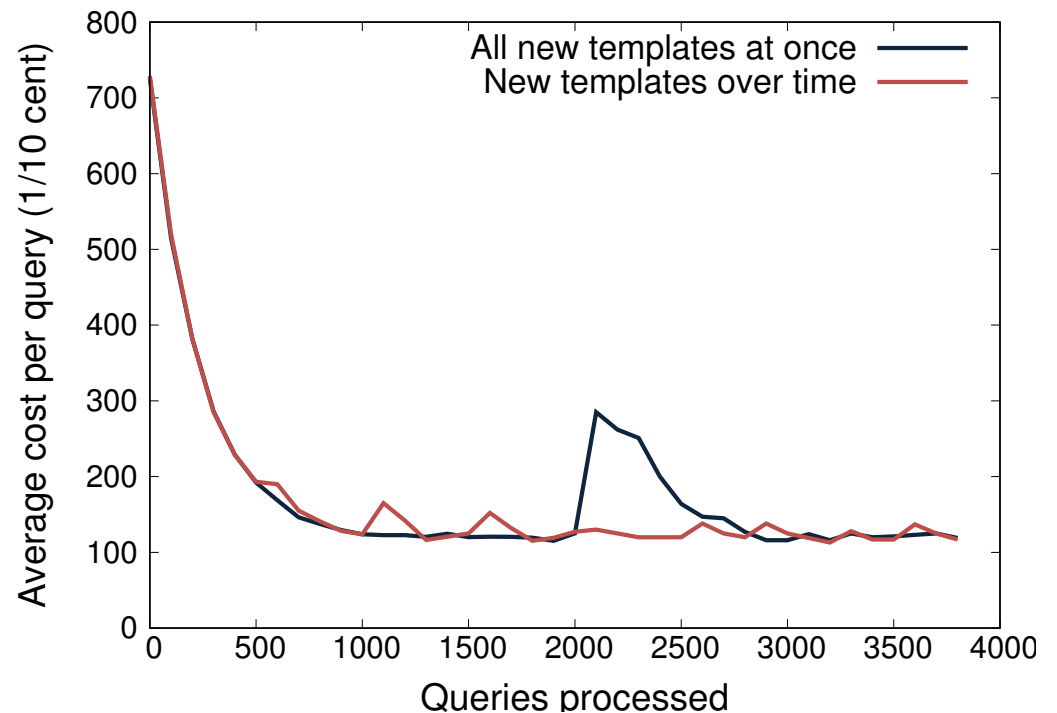levels with no pre-training or tuning**

# Adaptivity

**Training Data**

13 TPC-H templates
900 queries/hour
Poison distribution
Max SLO

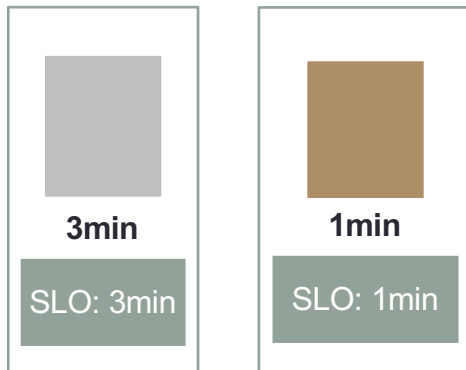*all new at once*: 7 new templates every 2000 queries (after convergence)

*new over time*: 1 new template every 500 queries



**WiSeDB models quickly adapt to new unseen before templates**

# Next Steps: Prediction-free Batch Scheduling

- ❑ Train once, use "**forever**"?
  - ❑ obsolescence detection and correction
- ❑ SVMs: Support Vector Machines
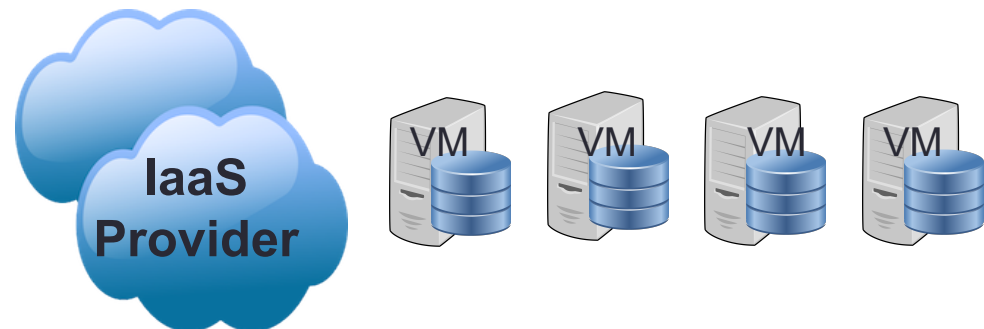  - ❑ detect decision boundaries based on cost, SLO slack, SLA violation risk

**3min**

SLO: 3min

**1min**

SLO: 1min

**Data Management Application**

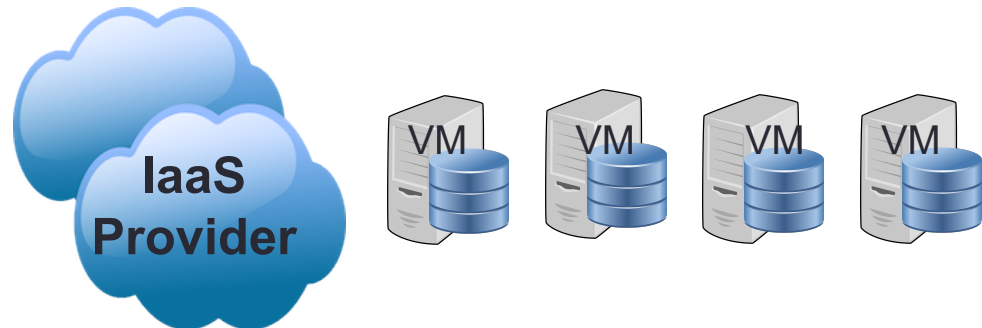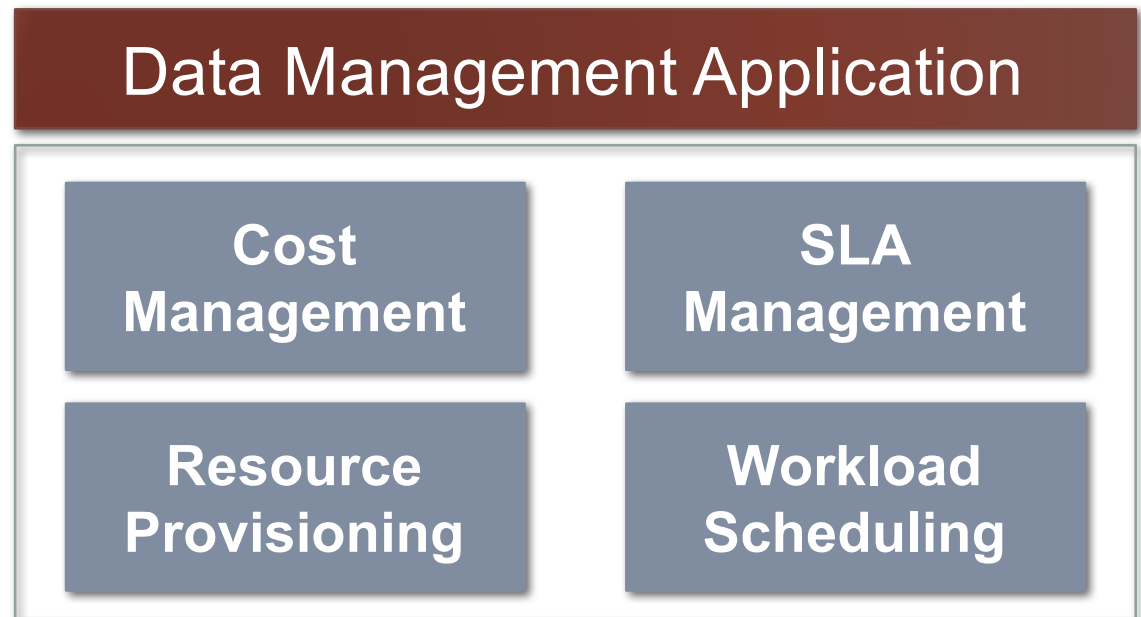| Cost Management | SLA Management |
|---|---|
| Resource Provisioning | Workload Scheduling |

**IaaS Provider**

VM  VM  VM  VM

# Next Steps: End-to-End Online Learning

- ❑ Query Scheduling
  - ❑ query ordering actions

- ❑ Shut-down strategy
  - ❑ hill-climbing learning

- ❑ Training overhead
  - ❑ search space reduction
  - ❑ warm bootstrapping

## Data Management Application

| Cost Management | SLA Management |
| --- | --- |
| Resource Provisioning | Workload Scheduling |

IaaS Provider

VM  VM  VM  VM

# Next Steps: Learning-based Pricing

- Resource consumption & SLA pricing

- Predicted cost == minimum price
  - no SLA violation fees

- System & economics interplay
  - fairness & competition affects system design
  - "learn" the pricing scheme & system decisions that offers pricing fairness

# Conclusions

❑ Cost and performance management for IaaS-deployed data managements apps are becoming more complex

  ❑ human ability to derive insight remains the same

❑ WiSeDB demonstrates how **ML techniques can**

  ❑ **offer insight** on the complex interplay of cost vs performance

  ❑ **discover** customized solutions for app-specific SLAs

  ❑ **automate** complex application management decisions

  ❑ **adapt** to workload and resource configurations

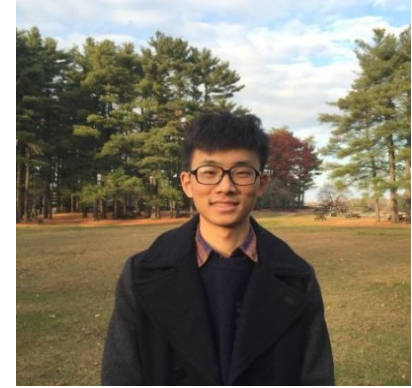  ❑ **build** systems that perform beyond unaided human heuristics

# Our Database Group



**Ryan Marcus**

Cloud Databases
Machine Learning



**Kyriaki Dimitriadou**

Interactive Data Exploration
Machine Learning



**Zhan Li**

Benchmarking Optimizers
Statistical Analysis

# THANK YOU

Questions?