# Peregrine: workload optimization for cloud query engines\*

Alekh Jindal

Gray Systems Lab

Microsoft

\* Peregrine: Workload Optimization for Cloud Query Engines. Alekh Jindal, Hiren Patel, Abhishek Roy, Shi Qiao, Jarod Yin, Rathijit Sen, Subru Krishnan. SOCC 2019.





# On-Premise

KI

DBA



ТАКСИ

Need to reach by 10, can we drive faster?

DAXI

K



Sure!

Announces Announces

### Cloud Query Engines

- Setup, installation, maintenance taken care of
- On-demand provisioning, pay as you go

#### .. ahhh!

Need to reach by 10, can we drive faster?

Sorry, we don't have a **DBA** 

Reality Check for providers:

- System developers == virtual DBAs!
- Too many cloud users, compared to system developers
- Too many support requests; often redundant
- Less time for feature development

Reality Check for customers:

- Lots of services to choose from (even within Azure, GCP, AWS)
- Lot of knobs to tune for **good perf** and **low cost**
- Lack of control; and lack of expertise
- And, the DBA is gone!

### Cosmos: big data infra at Microsoft

- 100s of thousands of machines
- Exabytes of data at rest; Petabytes ingress/egress daily
- 500k+ batch jobs / day
- 3B+ tasks executed / day
- 10s of millions interactive queries / day
- 10s of thousands of SCOPE developers
- 1000s of teams



### The missing DBA and the growing pain in Cosmos

- Large number of knobs/hints at script, data, plan level
  - Only few expert users
  - Rest need guidance
  - Survey: better tooling for improving SCOPE queries
- Support challenge
  - 10s of thousands incidents / years
  - 10 incidents per system developer on call
  - 100x users compared to system developers
  - ~10% growth in SCOPE workload in 2019



# On-premise pain -> Cloud pain









### The cloud opportunity

Fragmented on-premise workloads

## The Cosmos opportunity



Job metadata name, user, account, submit/start/end times

Query plans logical, physical, stage graph, estimates

Runtime statistics Operator-wise observables

Task level logs start/end events

Machine counters CPU, IO, etc.

### The case for a workload optimization platform

- DBA-as-a-Service
  - Another service in the cloud (easier integration)
  - Based on cloud workloads at hand (instance optimization)
- Engine agnostic
  - Not specific to different query engines, e.g., SCOPE, Spark, SQL DW, or etc.
  - E.g., view selection is still the same problem
- Global optimizations
  - Cloud workloads are organized into data pipelines
  - People often care about end-to-end aggregate costs in the cloud

### First cless lounge @ + Station reception 5 # Left loggage IE a implicates 8 7. Platforms 8 to 18 C at Relycant B X Toles El 2 Platforms 7 fb 7 Padoms1to7 0.3 Hickorne to Eustan Station Step 1: workload representation

**e a may may** 

Instrument, log, and collect workload characteristics

### Engine-agnostic workload representation



### Step 2: optimize for patterns









### Typical workload patterns

• Consider a simplified 2D space of data and queries



### Recurring pattern



### Recurring

Query templates appear over newer datasets

- Majority of production workloads
  - There is a regular ETL needed before other things can happen
- Opportunity to learn from the past
- Examples
  - Learned cardinality
  - Learned cost models
  - Learned resources
  - Learned etc.

### Recap from NWDS'19

### SCOPE Cardinality Estimation



#### Recap from NWDS'19

### SCOPE Cardinality Estimation



#### Towards a Learning Optimizer for Shared Clouds.

Chenggang Wu, Alekh Jindal, Saeed Amizadeh, Hiren Patel, Wangchao Le, Shi Qiao, Sriram Rao. VLDB 2019.

### **SCOPE** Cost Estimation

• Costs models are orders of magnitude off!



### SCOPE Cost Estimation

• Costs models are orders of magnitude off!



Manually tuned cost model

Feeding perfect cardinalities

- Pervasive use of user defined functions
- Complexity of big data systems
- Variance in the cloud environments

### Why cardinality is not enough?

- Incrementally add features
- Error drops from 110% to 40%
- Additional transformations needed
- Hard to come up with such heuristics





- Different feature weights
- Hard to instance optimize manually



### Ensemble of Models over Recurring Patterns



### **SCOPE** Cost Estimation

• Can learn pretty accurate cost models!



**Cost Models for Big Data Query Processing: Learning, Retrofitting, and Our Findings**. Tarique Siddiqui, Alekh Jindal, Shi Qiao, Hiren Patel, Wangchao Le. *SIGMOD 2020 (to appear)*.

### Similarity pattern



Queries

Similarity Queries over same

datasets have similarities

- Very typical in multi-user shared cloud environments
  - Cosmos, HDI, Ant Financial, ML workflows, etc.
- Opportunity for multi-query optimization

### • SCOPE compute reuse



#### Computation Reuse in Analytics Job Service at Microsoft.

Alekh Jindal, Shi Qiao, Hiren Patel, Jarod Yin, Jieming Di, Malay Bag, Marc Friedman, Yifung Lin, Konstantinos Karanasos, Sriram Rao. SIGMOD 2018.

### Spark Compute Reuse

- Instrument application log
- Analyze common subexpressions over Spark SQL plans
- Optimizer rules to automatically materialize/reuse in future queries
- Almost 30% improvement in total time on TPC-DS



#### SparkCruise: Handsfree Computation Reuse in Spark.

Abhishek Roy, Alekh Jindal, Hiren Patel, Ashit Gosalia, Subru Krishnan, Carlo Curino. VLDB 2019 (Demo).

### Step 3: feeding it back

- Actions
  - Insights
  - Recommendations
  - Self-tuning

Austrian 🗡

Reserved for your comfort

Austrian 🗡

Reserved for your comfort TOUS CONTRACTS

Austrian 🗡

EXIT AU

### Illustration: Scope and Spark query engines



### Peregrine Summary

- Easier to add newer features
- Easier to add newer engines
- Easier for people to participate
  - Researchers, developers, interns
  - Abstracts the painful steps
  - Build on top of each other
  - Focus on workload optimizations
- Gray Systems Lab: <u>aka.ms/gsl</u>



