# Towards an Adaptive, Fully Automated Performance Modeling Methodology for Cloud Applications

Ioannis Giannakopoulos [1], Dimitrios Tsoumakos [2] and Nectarios Koziris [1]

[1]:Computing Systems Laboratory, School of ECE, National Technical University of Athens, Greece
{ggian, nkoziris}@cslab.ece.ntua.gr
[2]: Department of Informatics, Ionian University, Greece
dtsouma@ionio.gr

# Introduction

# Performance Modeling (or profiling)

*"The ability to **predict** an application's behavior, given the parameters that affect it as input"*
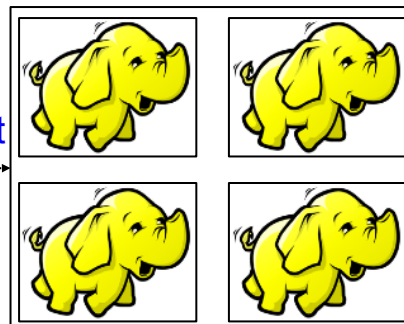
Valuable for:

- Optimal Resource Configuration
- Bottleneck Identification

- Application Scaling
- Multi-Engine Analytics

# Approaches

Existing approaches:

- White-box
  - Analytical Modeling
  - Interpretable Models
  - Good Accuracy for simple cases
  - Require Experience
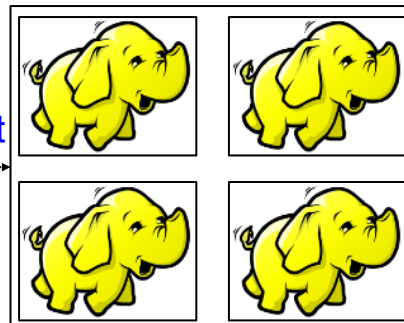  - Poor Results for complex cases

WordCount

# Approaches

Existing approaches:

- White-box
  - Analytical Modeling
  - Interpretable Models
  - Good Accuracy for simple cases
  - Require Experience
  - Poor Results for complex cases
- Black-box
  - Benchmarking & ML
  - Require no experience
  - Accuracy
  - Require time/computation

WordCount

WordCount
Nodes
Cores
Data Size

# "Black-box" Profiling Challenges

Traditional Challenges:

1. Which dimensions?
   a. data (e.g., dataset size, # of columns, distribution, etc.)
   b. resources (e.g., # nodes, # cores, RAM, etc.)
   c. workload (e.g., k parameter of k-means, clustering implementation, etc.)
2. Which model?
   a. Linear Models
   b. Neural Networks
   c. Support Vector Machines

New Challenge: **Massive Configuration Space**

# The challenges - Distributed Design

WordCount



Nodes = [2]

WordCount



Nodes = [2,4]

WordCount



Nodes = [2,4,6]

# The challenges - Cloud Computing

WordCount

Nodes = [2,4,6]
Cores = [1]                    3 combinations
RAM= [1]

WordCount

Nodes = [2,4,6]
Cores = [1,2]                  6 combinations
RAM= [1]

WordCount

Nodes = [2,4,6]
Cores = [1,2]                  12 combinations
RAM= [1,2]

**The configuration space grows exponentially with the
#dimensions and the #values per dimension!**

# Problem statement

Given:

1. Application A
2. Configuration space $D = d_1 \times d_2 \times \ldots \times d_n$
3. A positive number $B$

*"Estimate an application profile $F_A:D{\rightarrow}R$ using $D_s$ of size B,*

*with the highest possible accuracy!"*

or:

*"Find the top-B **representative** configurations to train a ML classifier"*

# Optimality & Observations

# Optimality & Observations

Optimal solution is **NP-Hard:**

- For each possible set of size B:
  - Deploy configuration
  - Train ML classifier
  - Evaluate it (test set)

Even deploying the entire D is not feasible:
- Prohibitive time + cost

Two observations regarding performance functions:
- Locality
- (piecewise) Linearity

Set 1 and Set 2

# Locality

**"Neighboring" configurations (tend to) give similar performance**

Hint:

*Divide-and-Conquer* to solve more (but simpler) problems

# (piecewise) Linearity

**"Neighborhoods" (tend to) present linear behavior (especially for resource-related dimensions)**

Hint:

Start with linear models

# Methodology

# Overview



**Algorithm 1** DT-based Adaptive Profiling Algorithm

1: **procedure** DTADAPTIVE($D$, $B$, $b$)
2:    $tree \leftarrow$ TREEINIT($\emptyset$), $samples \leftarrow \emptyset$
3:    **while** $|samples| \leq B$ **do**
4:      $tree \leftarrow$ PARTITION($tree$, $samples$)
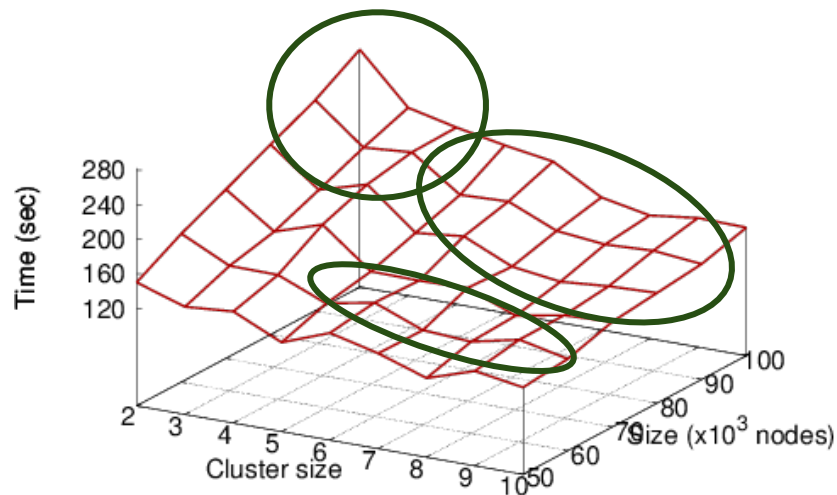5:      $s \leftarrow$ SAMPLE($D$, $tree$, $samples$, $b$)
6:      $d \leftarrow$ DEPLOY($s$)
7:      $samples \leftarrow samples \cup d$
8:    $model \leftarrow$ CREATEMODEL($samples$)
9:    **return** $model$

# Space Partitioning with Decision Trees

Decision Trees:

- Tree Structures
- Space Partitioning
- Node types:
  - Test nodes: Space Boundaries
  - Leaves: Linear Models

# Space Partitioning with Decision Trees

Space Partitioning → Decision Tree construction
- Substitute Leaves with Test nodes
- Which line to use?
  - Maximize leaf homogeneity
  - Well-known criteria: GINI Impurity, Entropy, Variance minimization, etc.
- Optimization Problem
  - Objective Function:   $Score(l) = -\dfrac{|L_1|R_1^2 + |L_2|R_2^2}{|L_1| + |L_2|}$

**Intuitively: the split line will generate two regions which best fit to linear models**

PARTITION → SAMPLE → DEPLOY → MODEL

# Space Sampling

Objective: distribute a deployment budget **b** (b<B) to each leaf

# Space Sampling

Objective: distribute a deployment budget **b** (b<B) to each leaf

# Space Sampling

- **How can we decide on the budget of each leaf?**
  - Leaf error (exploitation)
  - Leaf size (exploration)

- Leaf's score: $Score(leaf) = w_{error} \cdot \dfrac{error(leaf)}{maxError} + w_{size} \cdot \dfrac{size(leaf)}{maxSize}$

- # of samples assigned to each leaf *l*: *Score(l) x b*
- Uniform sampling inside the leaf

PARTITION → SAMPLE → DEPLOY → MODEL

17

# Space Sampling

Exploration vs Exploitation

$W_{error} = 1.0$
$W_{size} = 0.0$

# Space Sampling

## Exploration vs Exploitation

$$W_{error} = 1.0$$
$$W_{size} = 0.5$$

# Space Sampling

## Exploration vs Exploitation

$W_{error} = 0.0$
$W_{size} = 1.0$

# Deployment

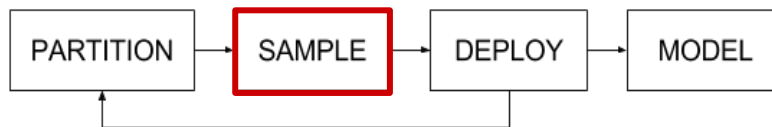- AURA: A Python-based deployment tool with error-recovery enhancements
- Application description: DAG of scripts for each module
- Transient failures (network glitches, poor coordination, etc.)
- AURA re-executes only scripts that failed
- Achieve idempotence via lightweight filesystem snapshotting

# Deployment DAG example



Hadoop Application Description:
- 1 Master node
- 2 Slave nodes

# Modeling

- Training a new Decision Tree from scratch
- When **B** comparable to #dimensions, DT degenerates into linear regression
  - Use multiple Regressors (e.g., ANNs, SVMs, etc.)
  - Keep the one with the lowest Cross Validation score
  - In practice, such budgets are **too low** to extract meaningful models
- Piecewise linearity
  - Approximates complex functions given higher deployment budgets

# Evaluation

# Methodology & Performance Functions

- Testbed: private Openstack Cluster (8 nodes, 200 cores, 600G RAM, 8TB disk)
- Competitors:
  - PANIC                                                                        exploitation
    - More samples between points samples with highest variance
  - Active Learning (Uncertainty Sampling)                   exploitation
    - More samples to areas of highest uncertainty
  - Uniform Sampling                                                   exploration
    - Uniform distribution of samples
- Modeling with WEKA regressors

# Methodology & Performance Functions

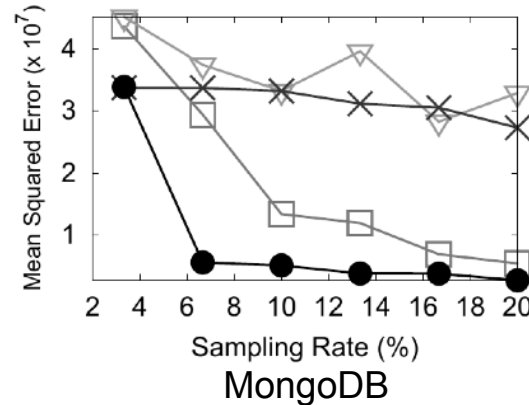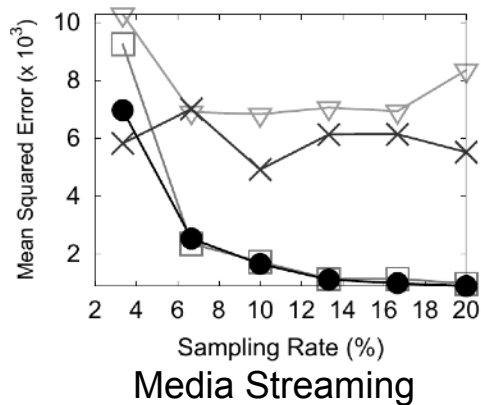| Application (perf. metric) | Dimensions | Values |
|---|---|---|
| Spark Bayes (execution time) | YARN nodes<br># cores per node<br>memory per node<br># of documents<br># of classes | 4–20<br>2–8<br>2–8 GB<br>0.5–2.5 ($\times 10^6$)<br>50–200 classes |
| Hadoop Wordcount (execution time) | YARN nodes<br># cores per node<br>memory per node<br>dataset size | 2–20<br>2–8<br>2–8 GB<br>5–50 GB |
| Media Streaming (throughput) | # of servers<br>video quality<br>request rate | 1–10<br>144p–1440p<br>50–500 req/s |
| MongoDB (throughput) | # of MongoD<br># of MongoS<br>request rate | 2–10<br>2–10<br>5–75 ($\times 10^3$) req/s |

SR: $|D_s|$ / $|D|$ x 100%

# Modeling Error



Bayes

Wordcount

Media Streaming

MongoDB

UNI
ACTL
PANIC
DTA

22

# Cost-aware profiling

$$Score(leaf) = w_{error} \frac{error(leaf)}{maxError} + w_{size} \frac{size(leaf)}{maxSize} - w_{cost} \frac{cost(leaf)}{maxCost}$$

| App/tions | SR | MSE | | | Cost | | |
|---|---|---|---|---|---|---|---|
| | | 0.2 | 0.5 | 1.0 | 0.2 | 0.5 | 1.0 |
| Bayes | 3% | +1% | -1% | -2% | -1% | -1% | -1% |
| | 20% | +4% | +10% | +5% | -7% | -9% | -12% |
| Wordcount | 3% | -1% | -5% | -1% | -4% | -6% | -1% |
| | 20% | 0% | +13% | +19% | -6% | -8% | -18% |
| Media Str. | 3% | -1% | -3% | +3% | -1% | -5% | -6% |
| | 20% | -6% | -2% | -8% | -7% | -11% | -26% |
| MongoDB | 3% | +6% | +12% | +13% | -2% | -3% | -4% |
| | 20% | -1% | -7% | -7% | -6% | -9% | -12% |

# Conclusions

- Divide-and-Conquer strategy
- "zoom-in" to *interesting* areas of the configuration space
- exploration vs exploitation
- Space Partitioning: important dimensions are partitioned more frequently
- Dimension Importance: linearity