# CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics

O. Alipourfard et al.

**Presented by Dmitry Kazhdan**

# Overview

- Background

- Prior work

- CherryPick

- Evaluation

- Criticism

- Recent work

- Conclusions

- Questions

# Background

# Background

Opportunities:

- Cloud computing

- Big data analytics

- Cost savings

# Background

Challenges:

- Complex performance model

- Cost model tradeoffs

- Heterogeneous applications

- Limited number of samples (from a large configuration space)

# Prior Work

- Ernest

- Coordinate descent

- Exhaustive search

- Random search

# CherryPick

# CherryPick

- Uses Bayesian Optimisation to build performance models

- Finds optimal/near-optimal configurations in only a few test runs
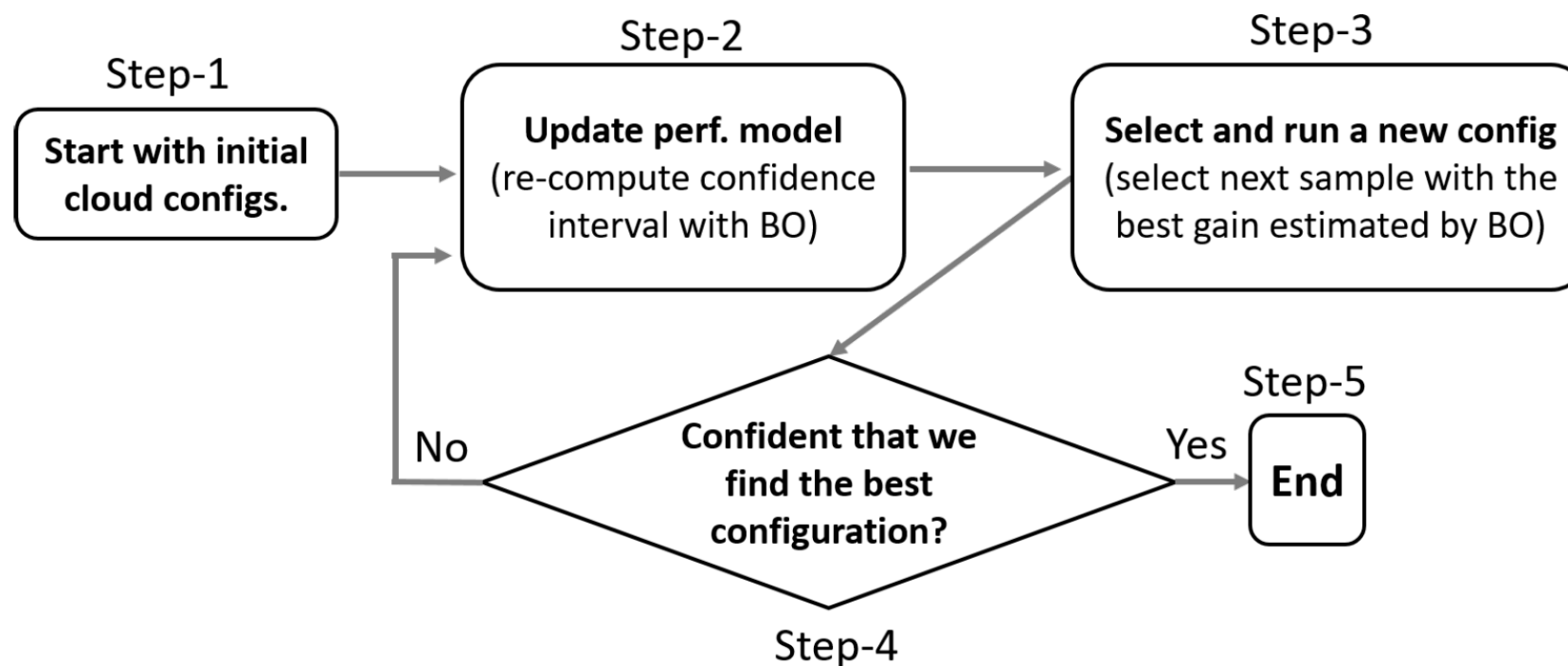
- Uses the acquisition function to draw samples

# CherryPick

Initial:

$$\underset{\vec{x}}{\text{minimize}} \quad C(\vec{x}) = P(\vec{x}) \times T(\vec{x})$$

$$\text{subject to} \quad T(\vec{x}) \le \mathscr{T}_{max}$$

Modified:

$$\log \tilde{C}(\vec{x}) = \log C(\vec{x}) + \log (1 + \varepsilon_c)$$

$$\text{subject to} \quad \log T(\vec{x}) \le \log \mathscr{T}_{max}$$

# CherryPick Workflow

# CherryPick Implementation

- Search Controller

- Cloud Monitor

- Bayesian Optimisation Engine

- Cloud Controller
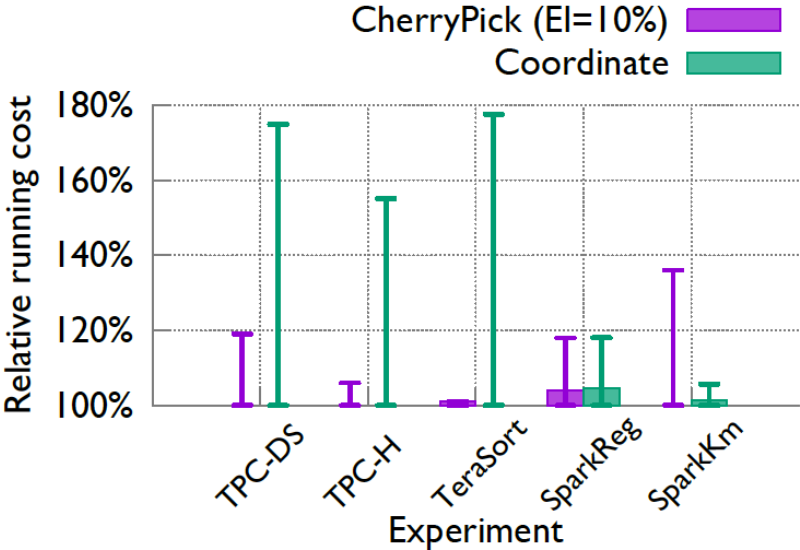
# Evaluation

# Evaluation

- Applications: TPC-DS, TPC-H, TeraSort, SparkReg, SparkKm

- 66 cloud configurations

- Objective: reduce cost of execution under runtime constraint

- Compared with:
  - Exhaustive search
  - Coordinate Descent
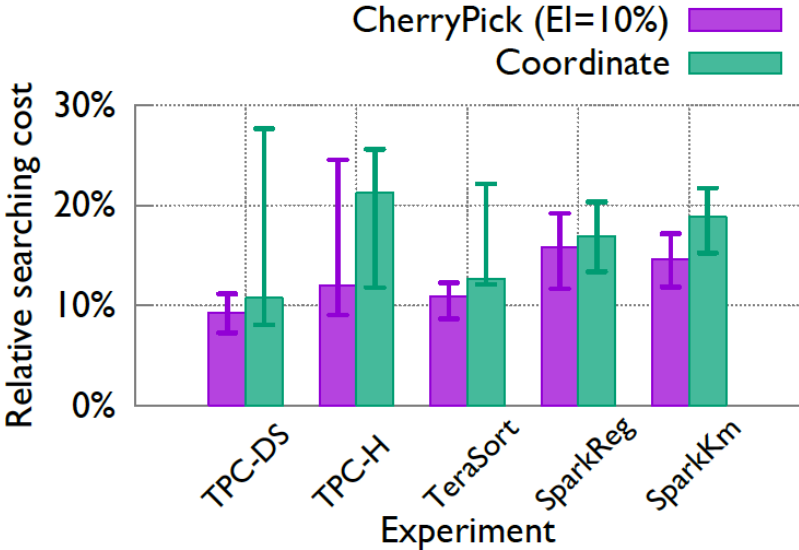  - Random Search (with a budget)
  - Ernest

# Evaluation

- Metric 1: the expense to run a job with the selected configuration

- Metric 2: the expense to run all sampled configurations
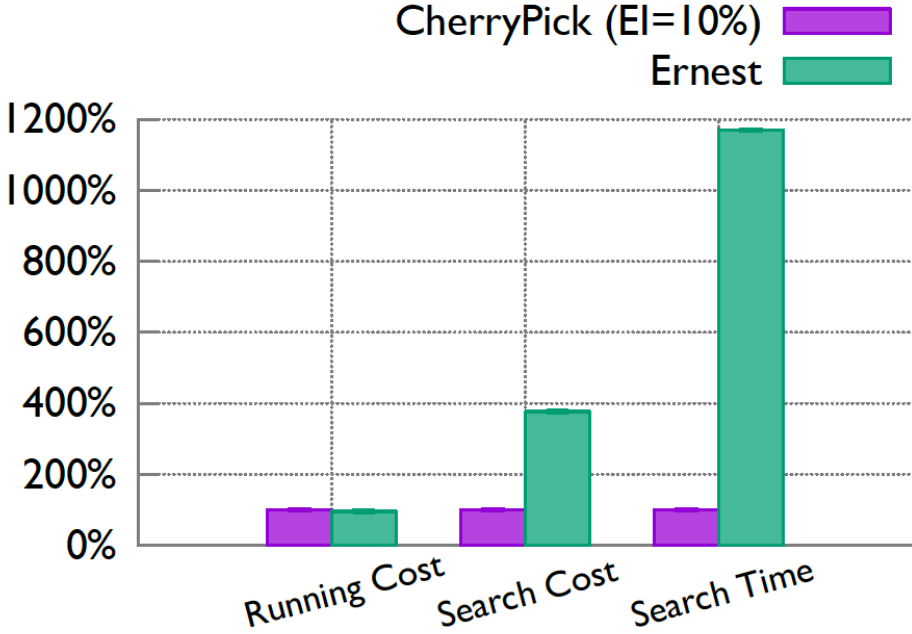
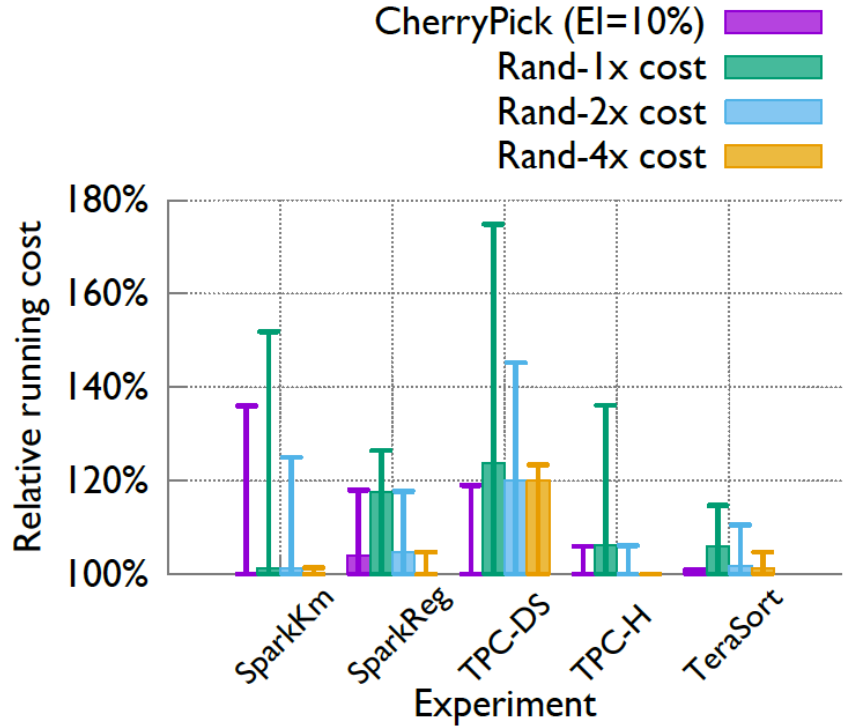- 20 independent runs

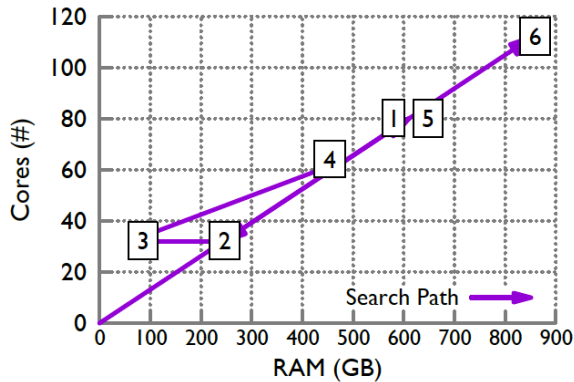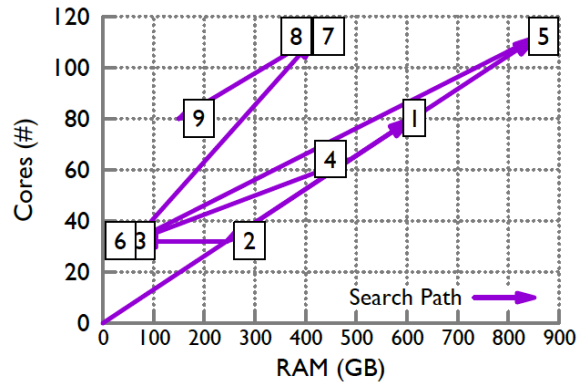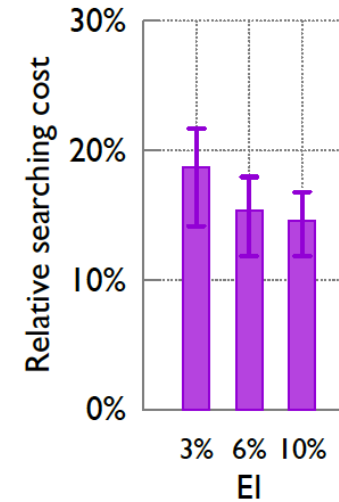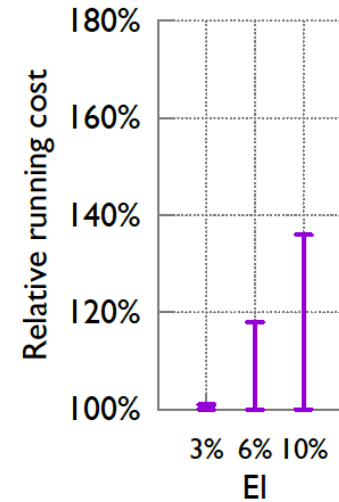- 10th, 50th and 90th percentiles computed

(a) Running cost

(b) Search cost

- Investigated parameter tuning
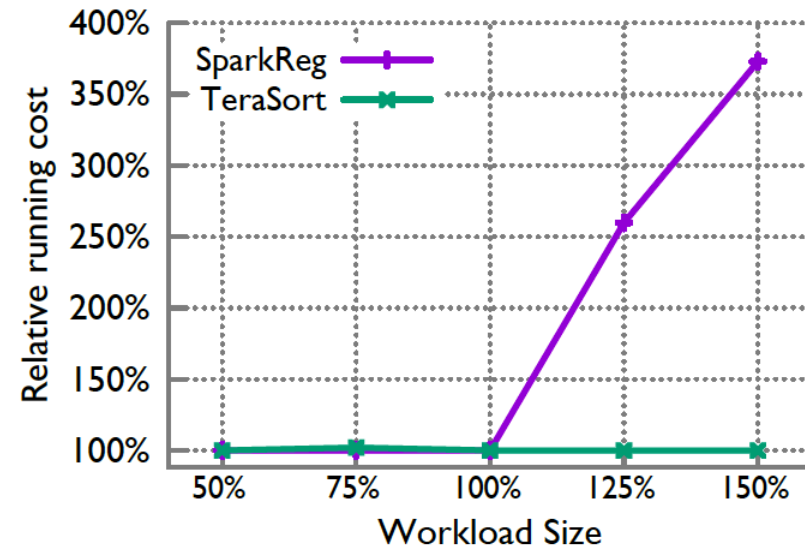
- Investigated performance behaviour



(a) SparkReg

(b) TPC-DS

- Handling workload variation

# Criticism

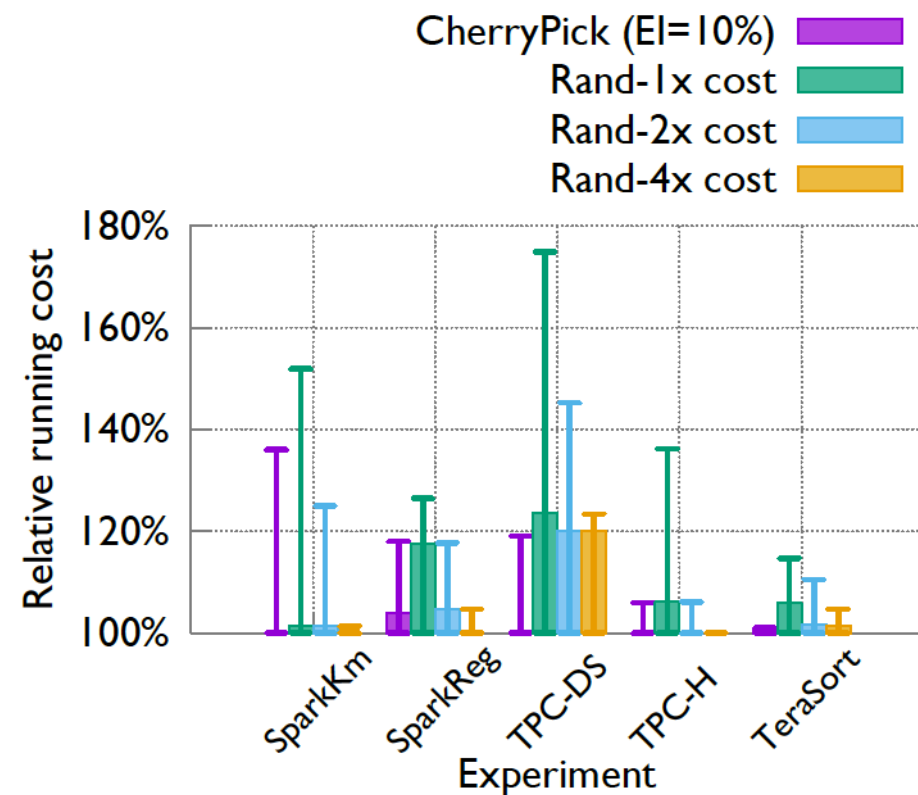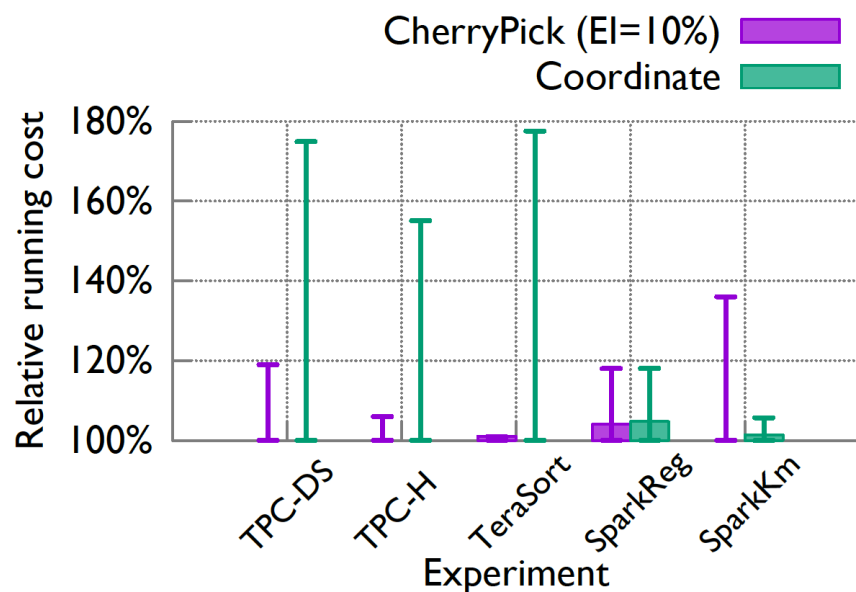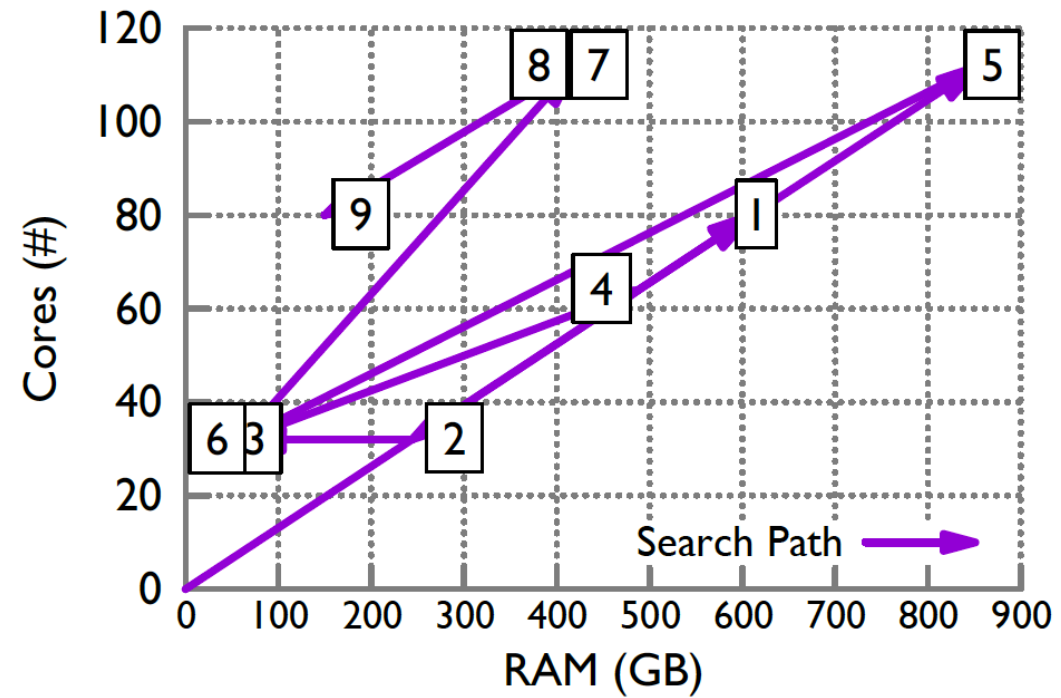# Criticism/Discussion

"With 4x cost, random search can find similar
configurations to CherryPick on the median"

- 3/4 comparison tasks are easy to beat (nothing to compare with)

- Not using available information efficiently

# Recent Work

# Recent Work

- PARIS

- Scout

- Arrow

- Micky

# Conclusions

# Conclusions

- Introduced CherryPick

- Compared to existing systems

- Presented evaluation results

- Criticism

# Questions?