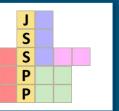
Performance Models to support HPC Co-Scheduling

Athanasios Tsoukleidis-Karydakis, Efstratios Karapanagiotis, Nikolaos Triantafyllis, Nectarios Koziris and Georgios Goumas







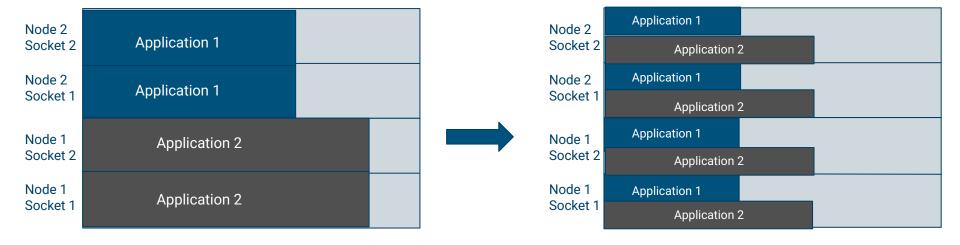
What is Co-Scheduling?

Problem: System Throughput



Proposed Solution: Co-Scheduling

An application typically spreads across twice the number of nodes compared to the classical scheduling case, but occupies half the number of cores on each node. Thus, two jobs can be allocated to the same node at the same time





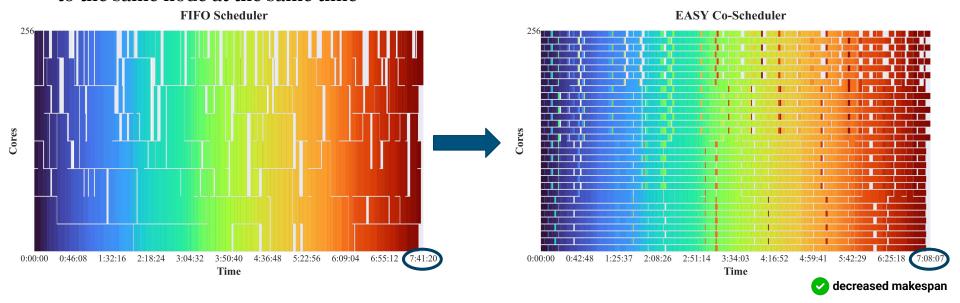
What is Co-Scheduling?

Problem: System Throughput



Proposed Solution: Co-Scheduling

• An application typically spreads across twice the number of nodes compared to the classical scheduling case, but occupies half the number of cores on each node. Thus, two jobs can be allocated to the same node at the same time

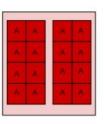


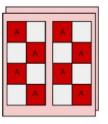


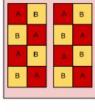
What is Co-Scheduling?

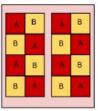
Compact mode:

- Traditional scheduling
- Exclusive node allocation







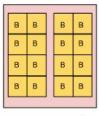


Co-location (striped) mode:

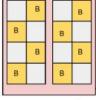
- Spread mode execution
- Idle halves occupied by other jobs

Spread mode:

- **Urgent computing**
- Double node spread
- Half-node occupancy
- Remaining half left idle



compact



spread

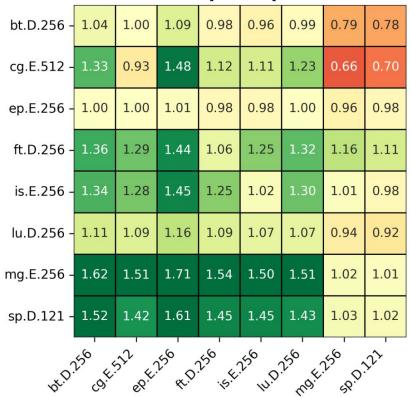


striped

resource contention in shared node resources can introduce performance degradation, leading to job slowdowns and counteracting these benefits



Speedup



$$JobSpeedup = \frac{Compact Execution Time}{Co-located Execution Time}$$

- pairwise execution of NAS benchmarks
- Classes : D, E

- 1.4

- 1.2

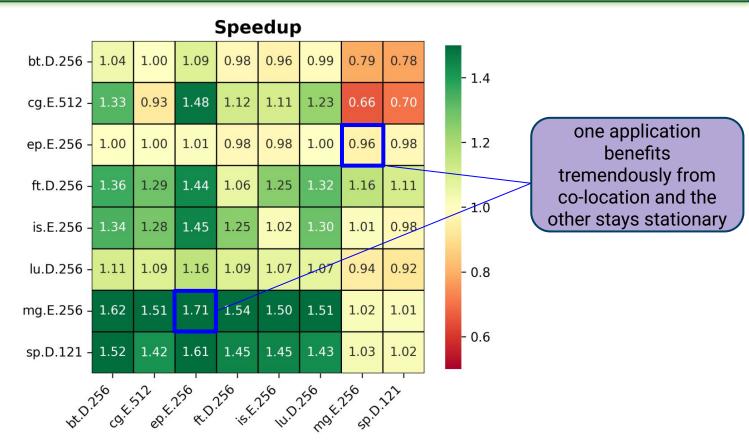
- 1.0

- 0.8

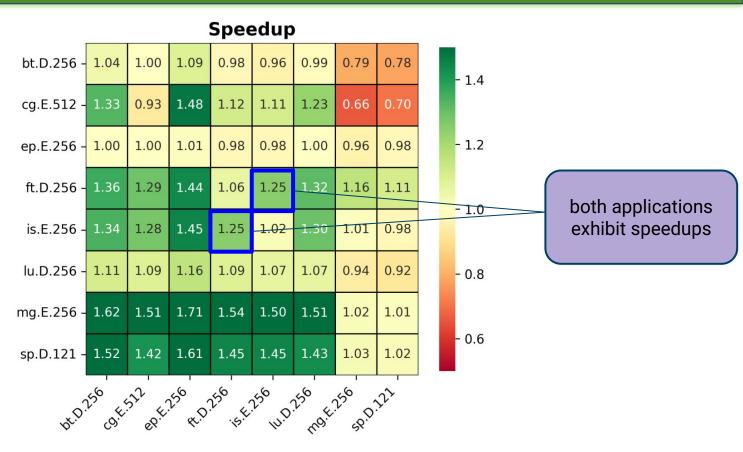
- 0.6

Various number of processes

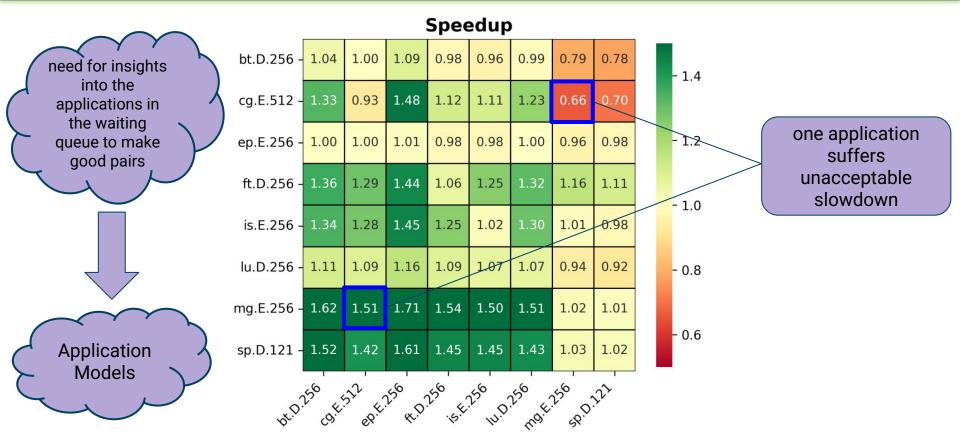












Goal and Infrastructure

- **Goal of this work :** Develop models that predict an application's performance improvement or degradation when co-located with other applications
 - **ARIS supercomputer** (operated by GRNET, Athens, Greece)
 - We utilized the **perf** Linux tool and **mpiP**
 - Collected Metrics:
 - o FLOPS
 - Memory Bandwidth
 - LLC hits/misses
 - MPI stats
 - MPI time
 - MPI time per directive
 - NAS Parallel Benchmarks (NPB) and SPEChpc 2021 benchmarks



Model Requirements



Need to have:

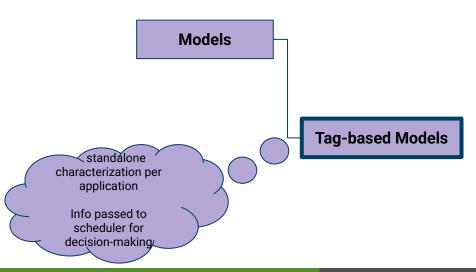
- High accuracy and precision
- Minimal data use
- Lightweight profiling
 - few hardware counters
 - few application runs
- Fast decision



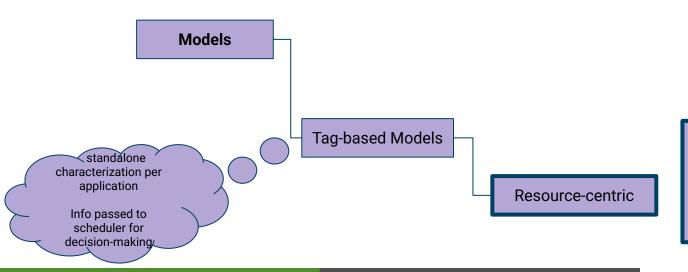
Nice to have:

- Interpretability
- Generalizable to other machines





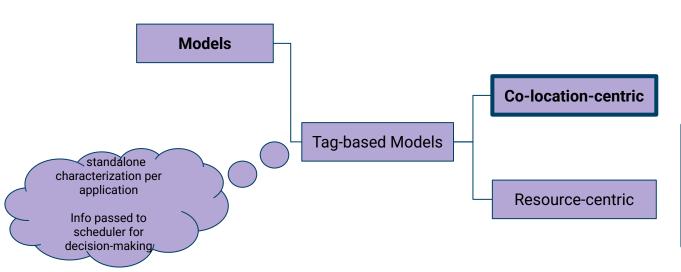




characterize applications according to their resource consumption patterns



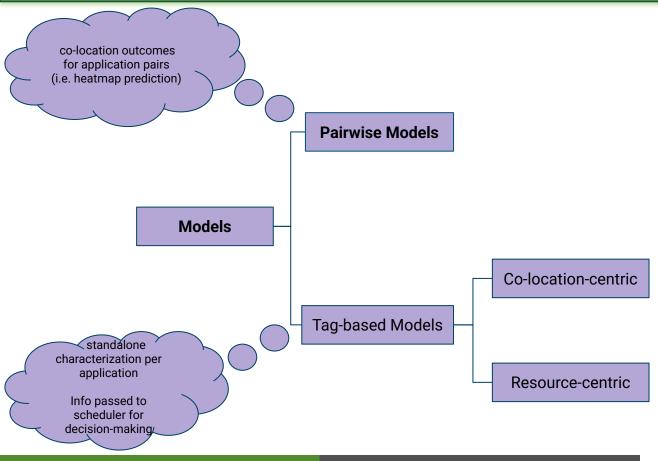




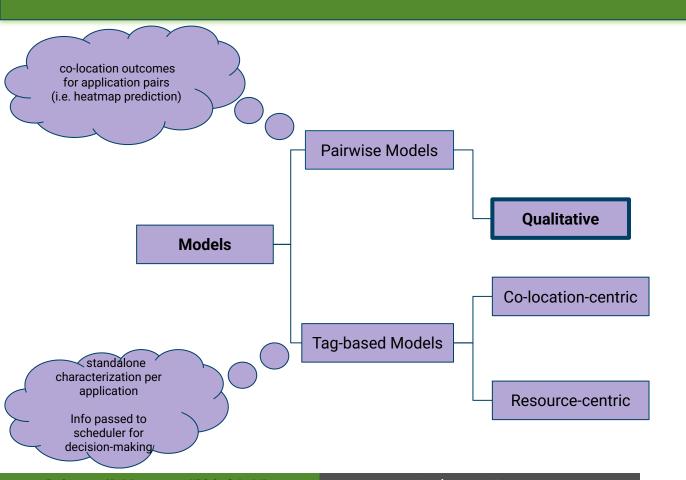
each application is labeled as either co-location friendly or co-location unfriendly







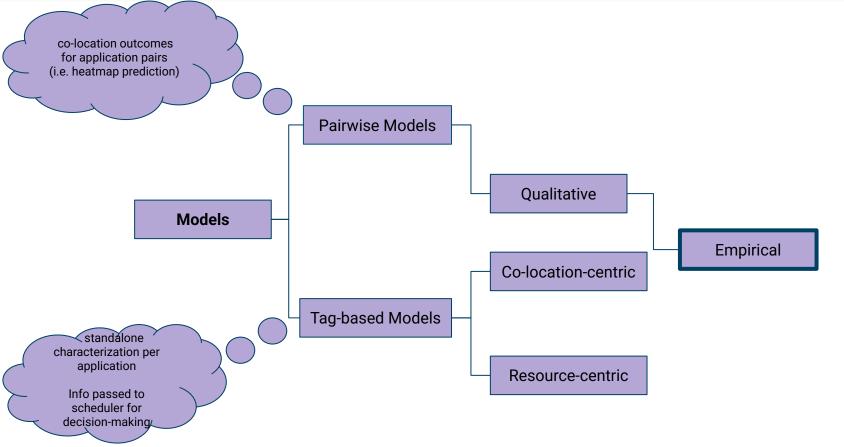




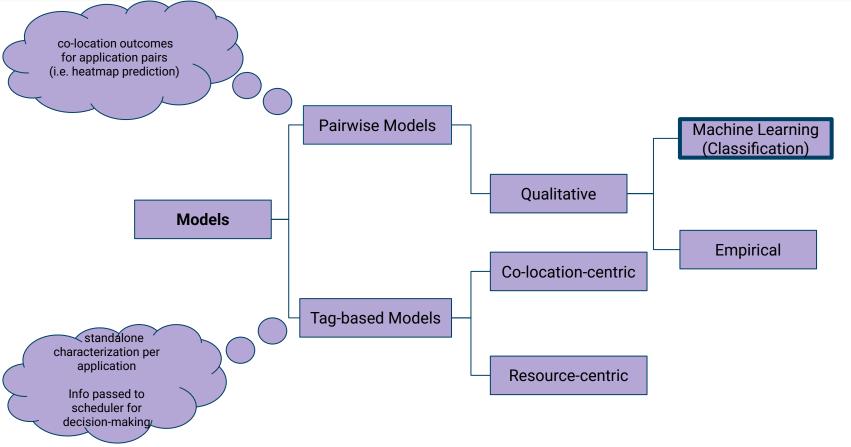
use of labels (e.g. good, bad, stationary)



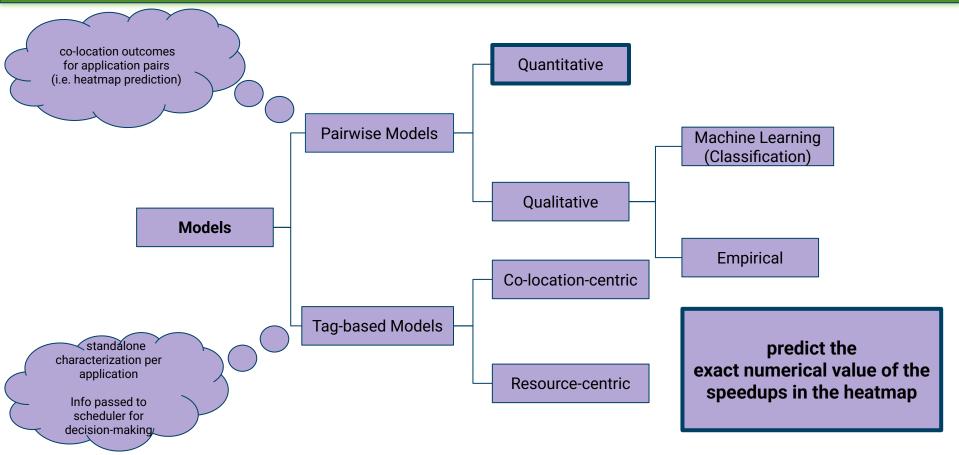


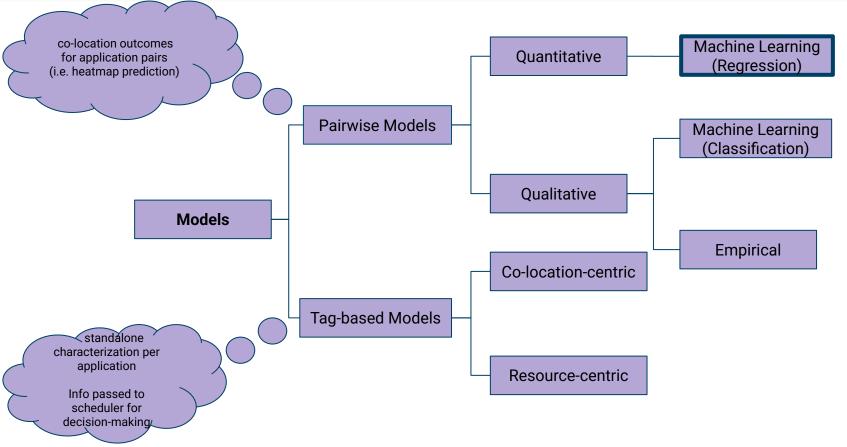




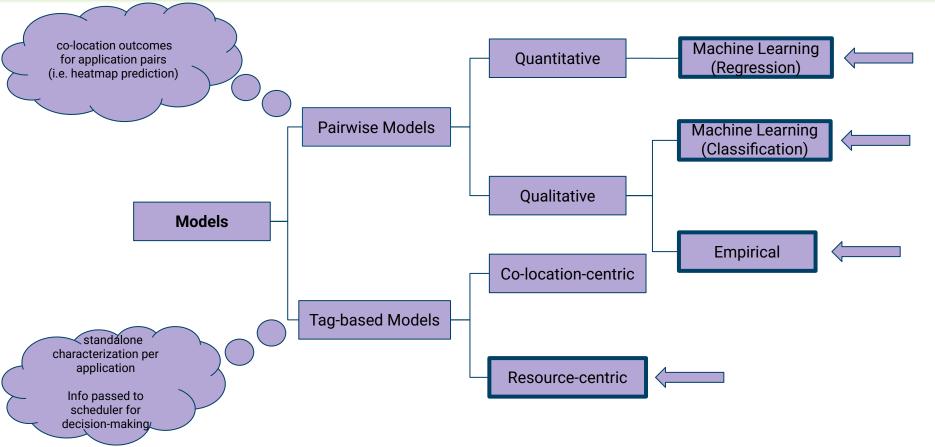




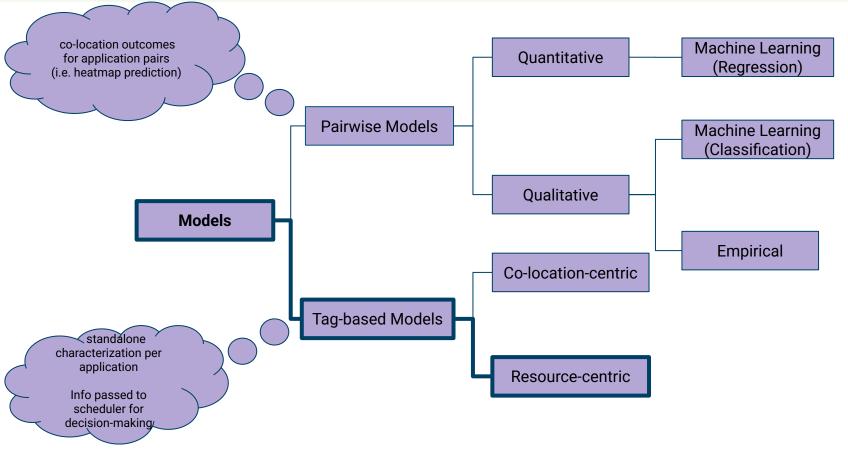














Resource-centric tag-based models

Goal: Retrieve a <u>standalone resource consumption-based characterization</u> for each application in the queue sp.D.121

aood

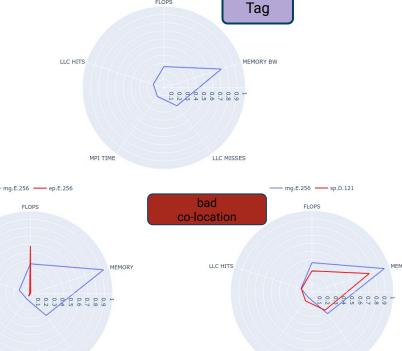
co-location

LLC HITS

FLOPS

LLC MISSES

Starting Point: «one can reduce the average slowdown experienced by co-located applications by simply preventing instances of the same program from being co-located together» (de Blanche & Lundqvist, 2016)



Extension:

Minimal overlap between spider plots of two applications



MPI TIME





LLC MISSES

MPI TIME

Resource-centric tag-based models

Static Experiment: (300 applications, 100 experiments)

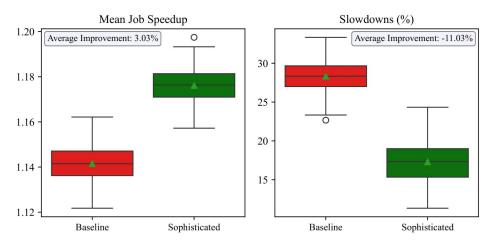
Baseline Case

random co-location

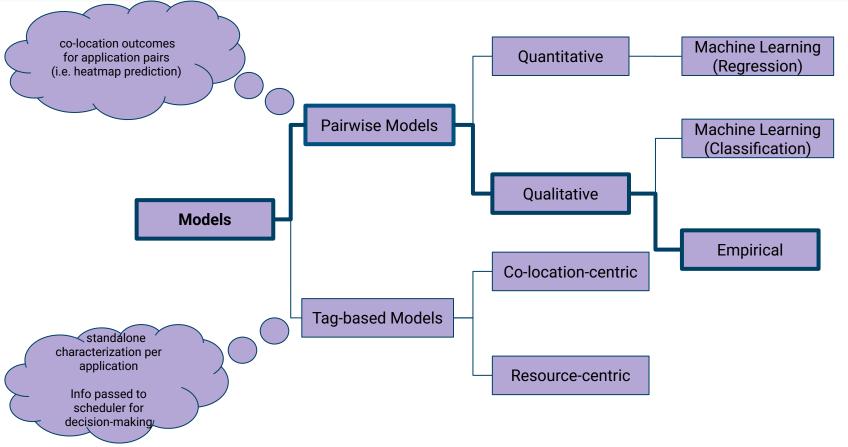
Sophisticated Case

minimize spider plot overlap greedily

Evaluation:



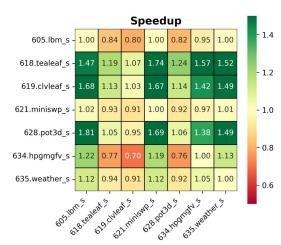




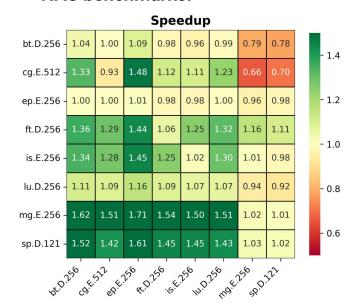


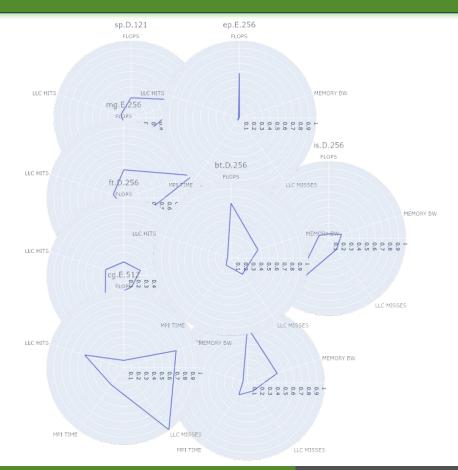
- **Goal:** devise an empirical way to qualitatively predict the heatmap of an application pool
- use of categorical labels : good, stationary, bad

SPEChpc 2021 benchmarks:



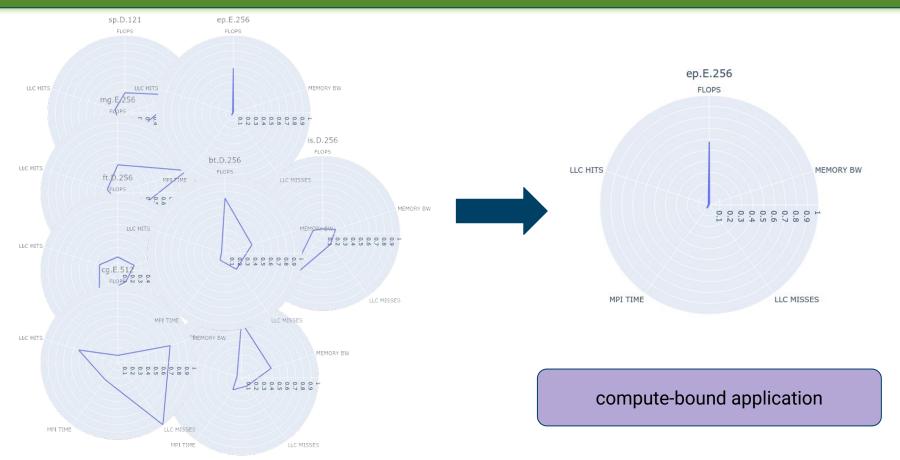
NAS benchmarks:





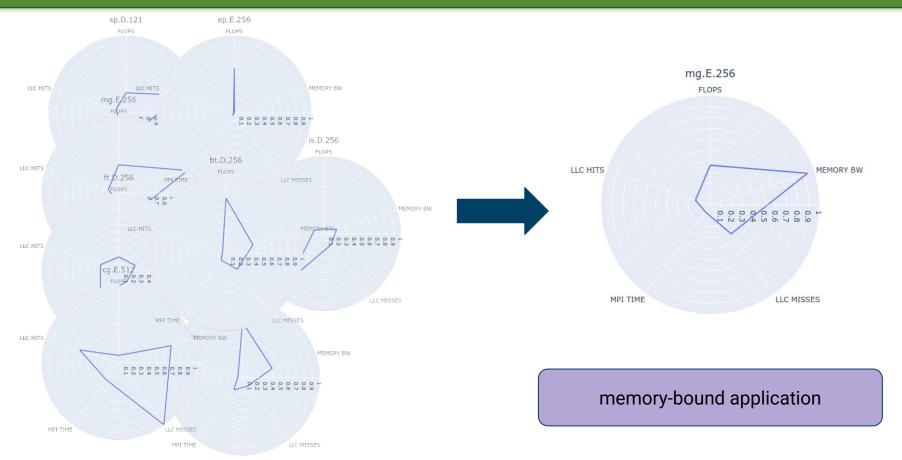




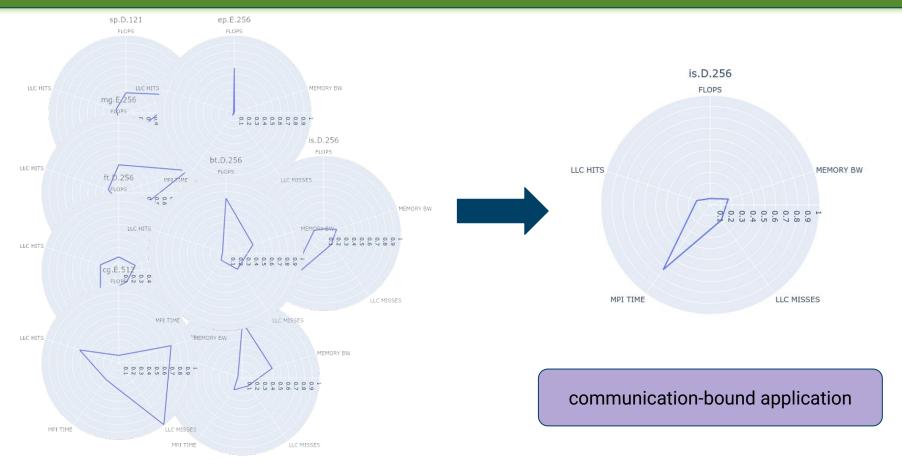






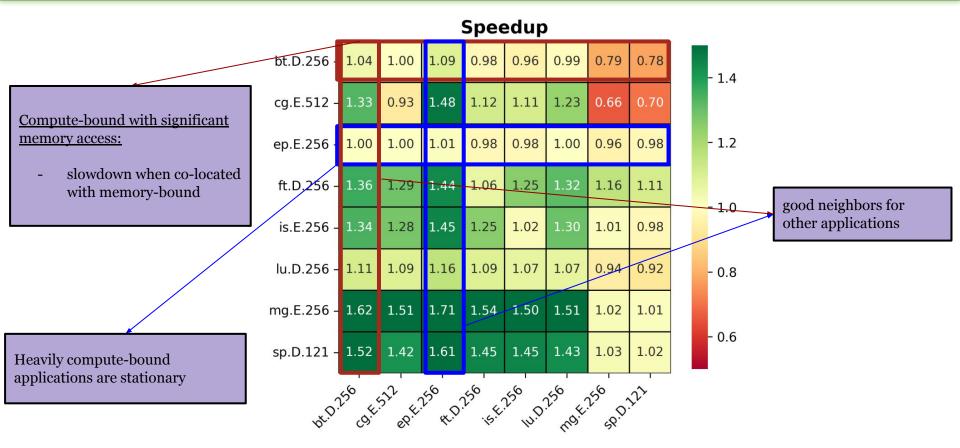








Compute-bound applications

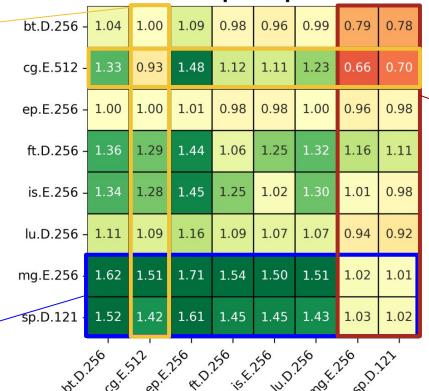


Memory-bound applications

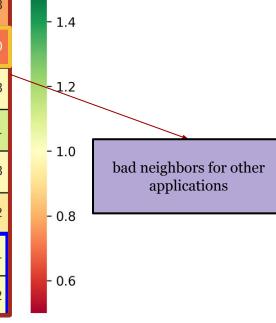
Memory-bound with irregular memory access:

- Low speedups
- slowdowns with memory-bound applications
- good neighbours

- remarkable speedups
- worsened performances with other memory-bound applications

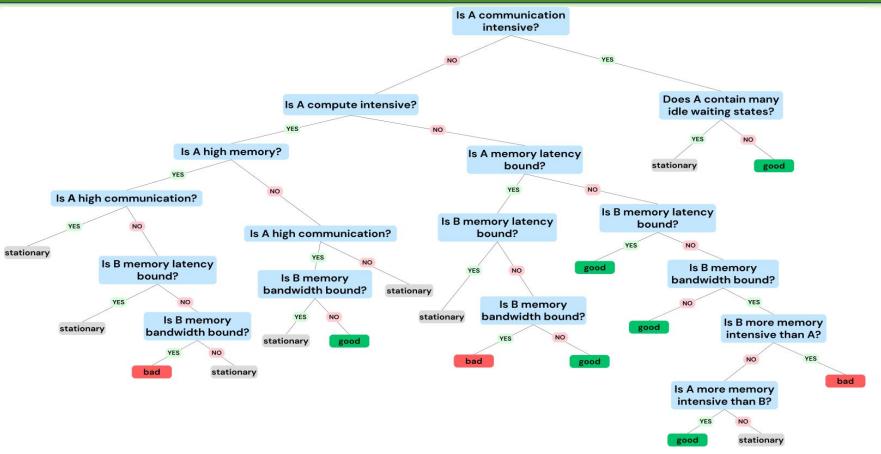


Speedup





Empirical Model

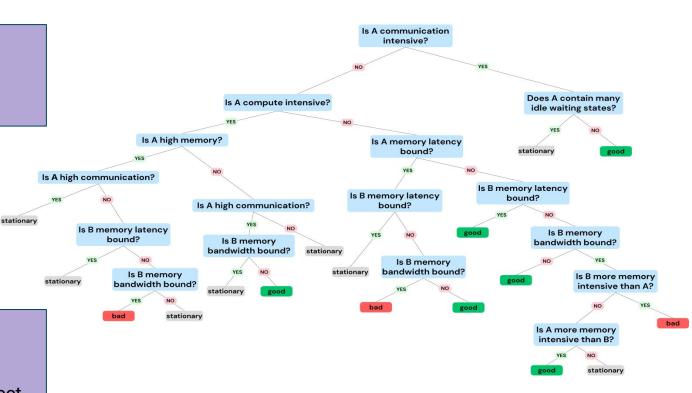


Empirical Model - Evaluation

NAS:

★ 83.43% fully correct

★ 7.1% sufficiently correct



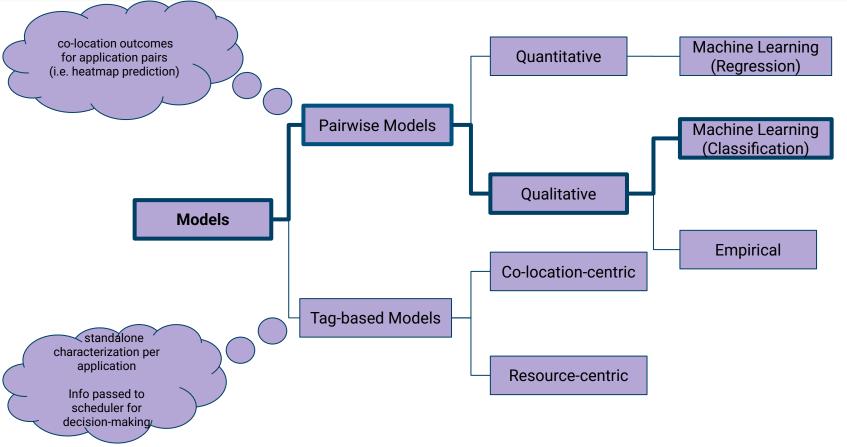
SPEChpc 2021:

- ★ 63.27% fully correct
- ★ 22.45% sufficiently correct









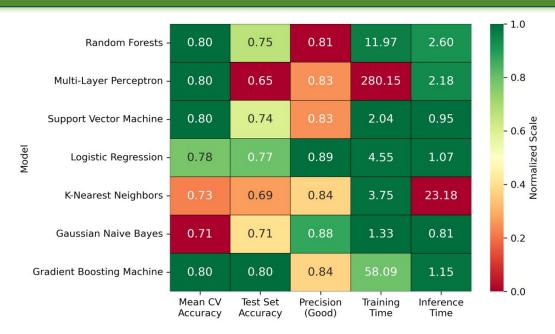


ML-based Classifiers

- **Features:** the five axes previously presented in the spider plots, normalized according to system bounds
- **Dataset:** 218 entries in total, with 70% allocated for training and the remaining 30% \rightarrow for testing
- **Labels**: good, stationary, bad
- **Methodology:** Train many models, with different hyperparameter combinations, from 7 different types of models, using 5-fold cross validation and keep the model from each type with the highest mean cross-validation accuracy score



ML-based Classifiers

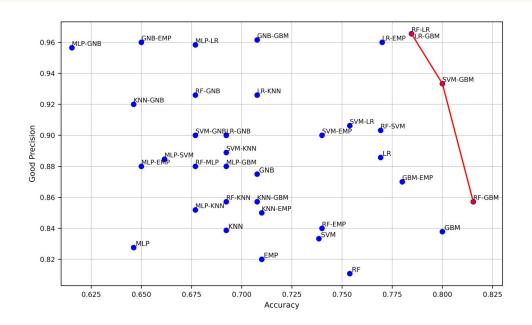


- **Ensemble boosting techniques** (e.g. RF, GBM) achieved the best overall results
- **Low-complexity models** (KNN, GNB) failed to grasp the complex, non-linear relations between the input features
- **Neural Networks** (MLP) need to be re-evaluated when bigger datasets are available
- Training and inference times appear to be reasonable

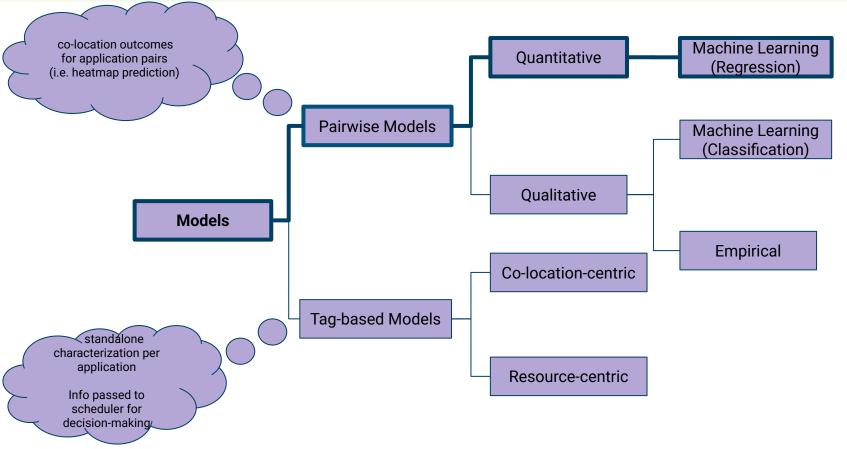




ML-based Classifiers: Combining two models



- Mispredicting many bad co-locations as good (sacrificing precision) can result in a slowdown in the overall makespan
- Achieving <u>high accuracy</u> and <u>precision for the "good" label</u> concurrently, is a two-criteria optimization problem that can be solved by utilizing combinations of models as shown in the Pareto plot



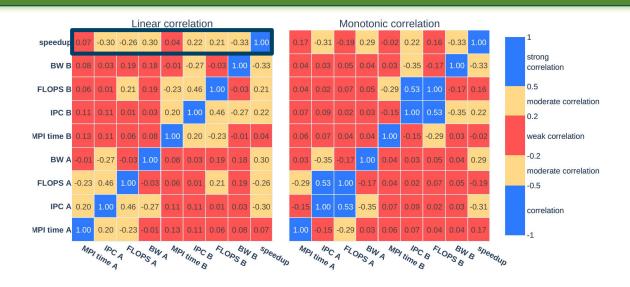


ML-based Regression

- → **Features:** MPI time, IPC, FLOPS, and Memory Bandwidth (BW)
- → **Dataset :** 788 entries in total, with 65% allocated for training and the remaining 35% for testing
- → **Labels**: Numerical value of the speedup for each heatmap cell
- → **Methodology:** Train many models, with different hyperparameter combinations, from 7 different types of models (plus a Dummy Regressor used for comparison purposes), using 5-fold cross validation and keep the model from each type with the highest mean cross-validation accuracy score



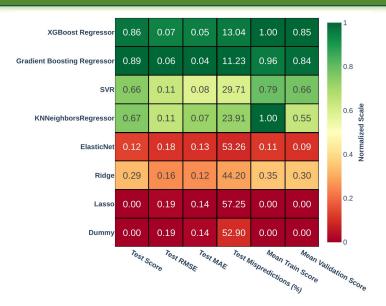
ML-based Regression



- moderate correlation of the input features with the dependent variable (speedup)
- Hint that <u>non-linear relations</u> are at play



ML-based Regression



- Linear models underperformed due to low feature-speedup correlation
- Ridge best among linear models; SVR and KNN better but KNN overfitted
- **Boosting models** highest R², lowest RMSE/MAE, consistent performance across trials



Conclusion

Review of the presented methods with regards to the requirements:

	Minimal Data use	Lightweight profiling	Fast Inference	High Accuracy and Precision	Interpretability	Generalizable to other machines
Tag-based				N/A		
Empirical						
Linear/Simple ML						
Boosting techniques						
Neural Networks				Needs further research		

Future Work

Re-evaluation of Neural Networks' accuracy and performance utilizing bigger datasets

Use of **new/composite features** with bigger correlation with co-location speedup

Integration of these application models to a real co-scheduler or a simulation tool to test them in a dynamic co-scheduling scenario

Thank you for your attention!

Q & A session

