

Auto-tuning Spark Big Data Workloads on POWER8: Prediction- Based Dynamic SMT Threading

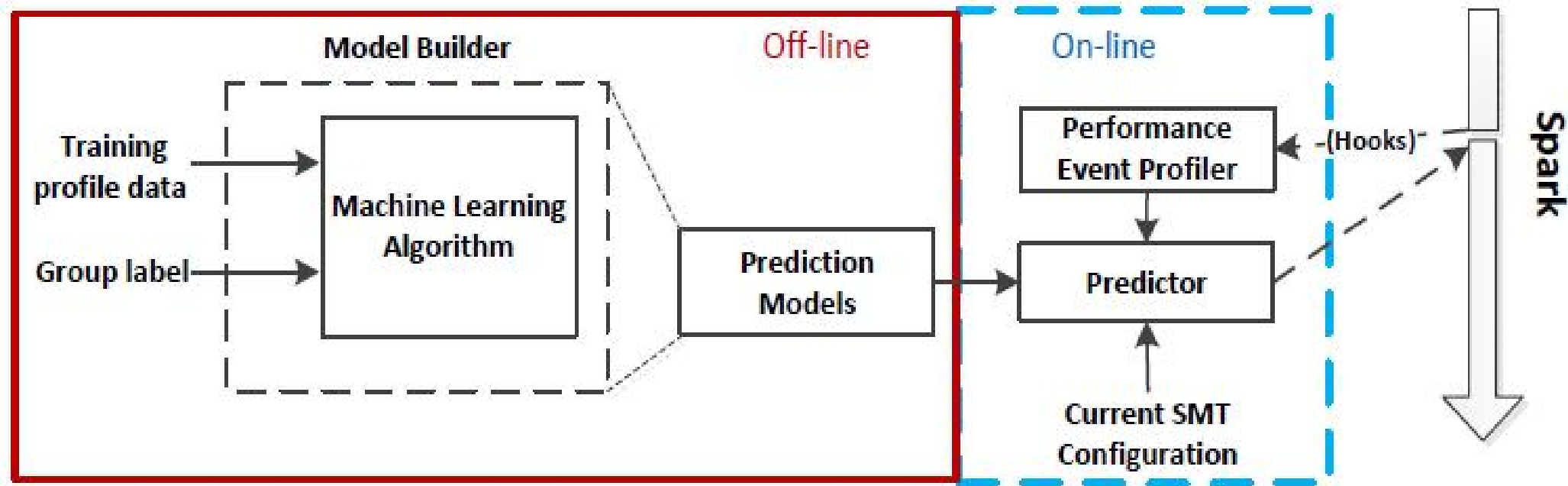
Research Question

- ◆ Different big data workloads have disparate architectural characteristics
- ◆ Identify the most efficient SMT configuration to achieve the best performance
- ◆ Focus on auto-tuning SMT configuration for Spark-based big data workloads on POWER8
- ◆ Methodology could be generalized to Hadoop, Flink

- ◆ Propose a prediction-based dynamic SMT threading(PBDST) framework to adjust the thread count in SMT cores on POWER8 processor by eight machine learning algorithm

- ◆ SMT: simultaneous multithreading has potential to increase the processor resource utilization
- ◆ POWER8: a processor architecture, support up to 8-way multithreading
- ◆ Profile: hardware performance counters
- ◆ Configuration: number of threads

Methodology-PBDS framework

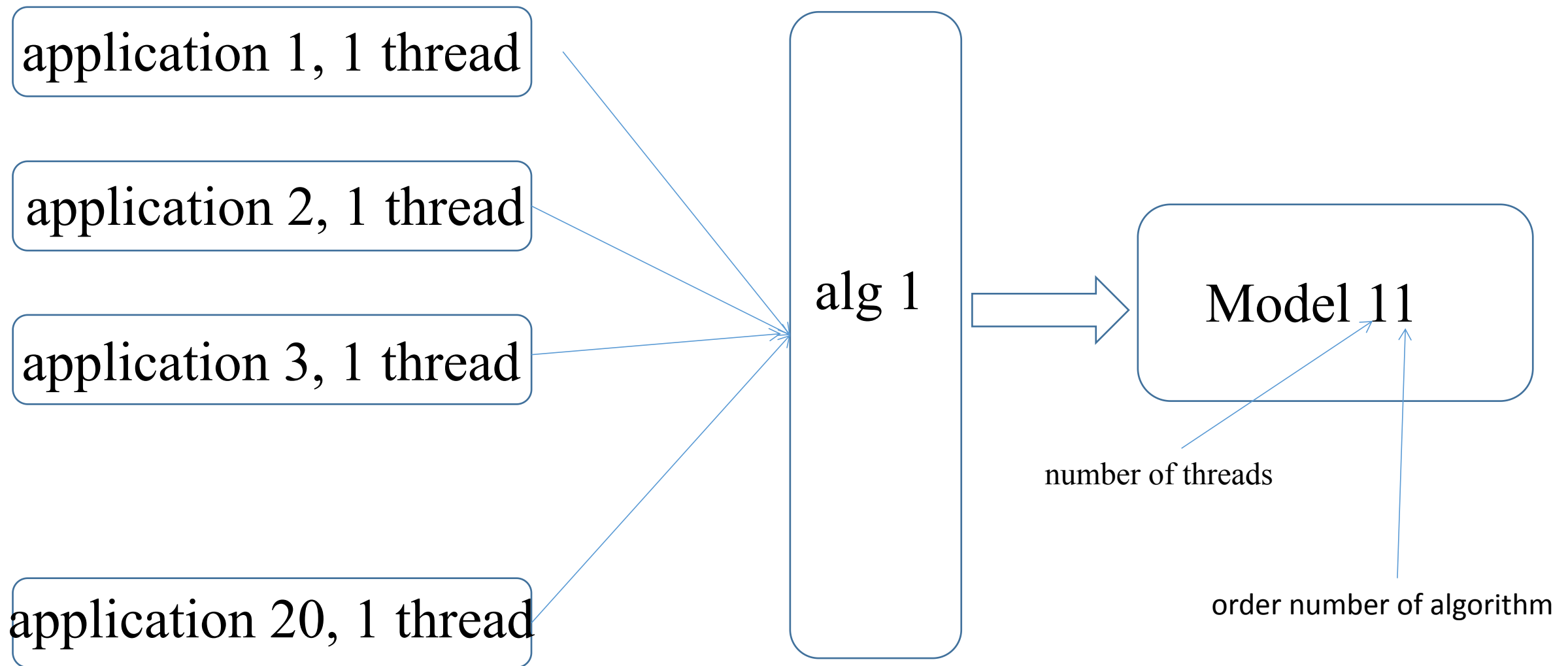


- ◆ 29 Spark application
- ◆ 20 application used for training set
- ◆ 9 application used for testing set

Training Set	No.	Name	Source
	1	Spark-als	Spark-perf
	2	Spark-kmeans	Spark-perf
	3	Spark-spearman	Spark-perf
	4	Spark-chi-sq-feature	Spark-perf
	5	Spark-NaiveBayes	BigDataBench
	6	Spark-terasort	BigDataBench
	7	Spark-grep	BigDataBench
	8	Spark-sort	BigDataBench
	9	Spark-pykmeans	Spark-perf
	10	Spark-pyals	Spark-perf
	11	Spark-pearson	Spark-perf
	12	Spark-pybayes	Spark-perf
	13	Spark-pycout-w-fltr	Spark-perf
	14	Spark-summary-statistics	Spark-perf
	15	Spark-pyspearman	Spark-perf
	16	Spark-count	Spark-perf
	17	Spark-chi-sq-mat:	Spark-perf
	18	Spark-gmm	Spark-perf
	19	Spark-svd	Spark-perf
	20	Spark-WordCount	BigDataBench
Testing Set	21	Spark-ConnectedComponent	BigDataBench
	22	Spark-block-matrix-mult	Spark-perf
	23	Spark-word2vec	Spark-perf
	24	Spark-pca	Spark-perf
	25	Spark-pypearson	Spark-perf
	26	Spark-glm	Spark-perf
	27	Spark-PageRank	BigDataBench
	28	Spark-int-sory-by-key	Spark-perf
	29	Spark-fp-growth	Spark-perf

- ◆ run training set with different SMT configurations, from one thread per core to eight threads per core
- ◆ for each workload, collect the microarchitecture level statistics during the whole lifetime
- ◆ one model for each SMT configuration

Methodology- offline training

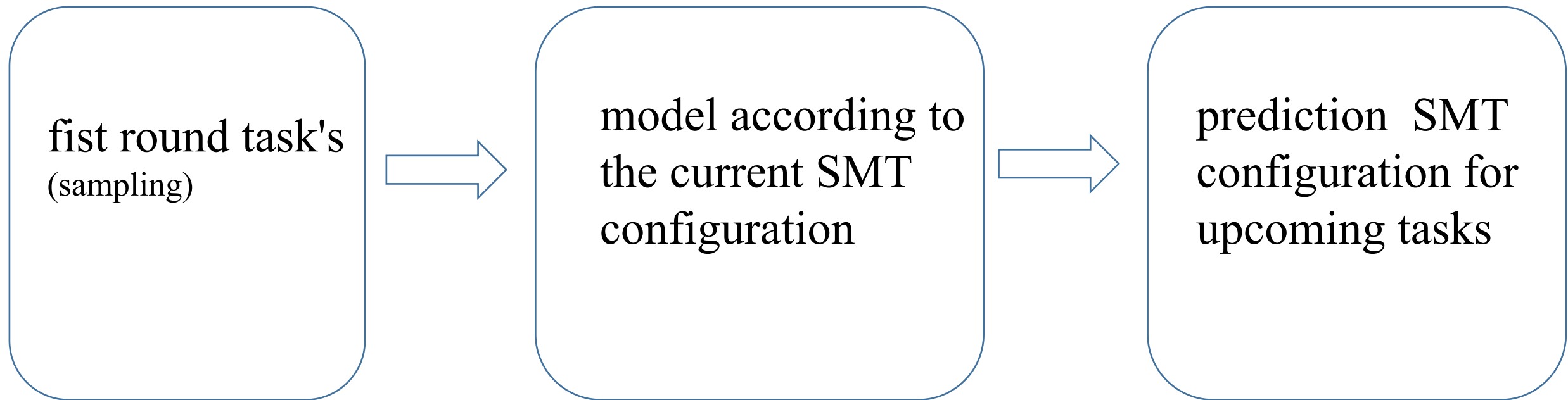


application 1, 1 thread

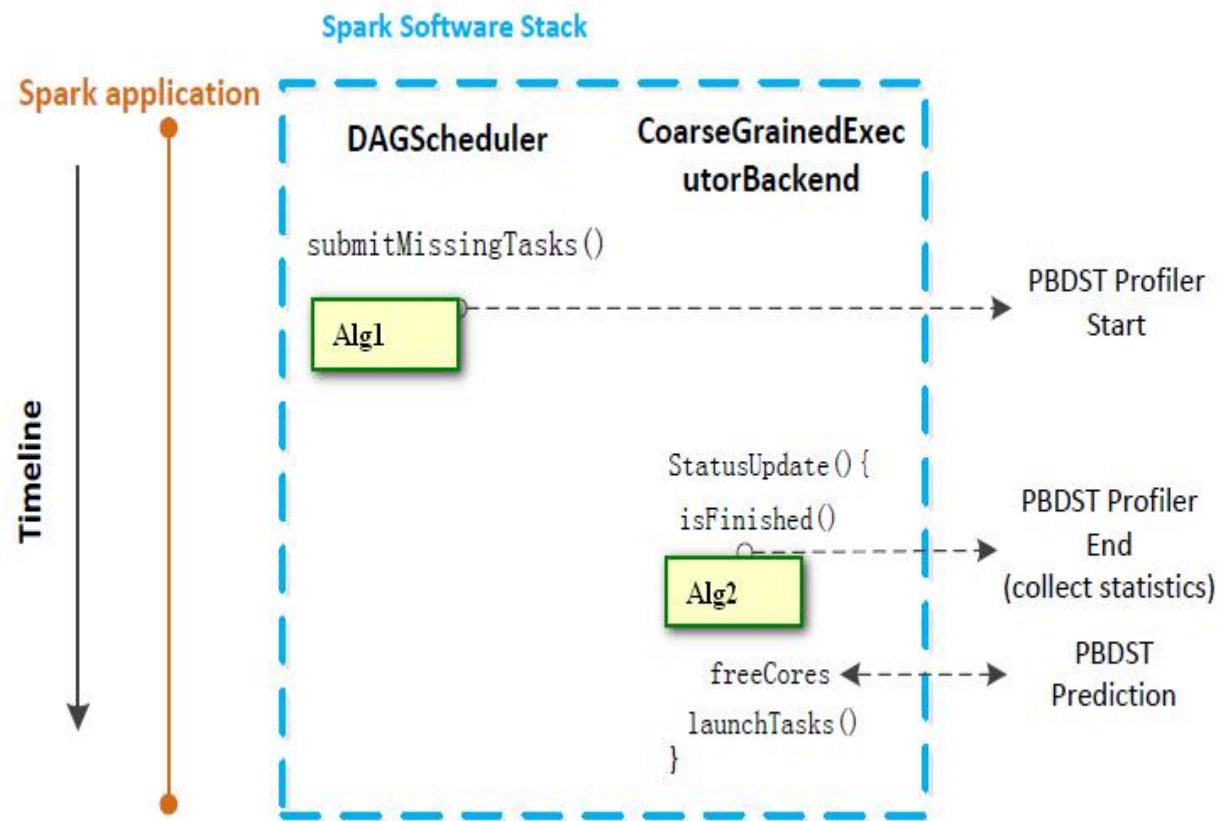
Input data: *traing profile data*, this is a matrix and each row is the micro-architecture level statistics of applications in thread 1

Output data: *group label*, this is an array which records the best SMT configurations for application 1,application 2,...application 20

In 9 testing application, for every application, use the first round of tasks' micro-architecture level profiles to predict the optimal SMT configuration for the upcoming tasks in the same stage



Auto-tuning



Algorithm 1 Algorithm used to trigger the PBDST hooks in *DAGScheduler*

Input: Variables:

Current SMT configurations: *CurSMT*;

The number of cores: *cores*;

The TaskSet information: *taskSet*;

1. **if** *submitTaskSet()* **then**
2. **if** *noStageIsRunningWithProfiler()* **then**
3. **if** *taskSet.size* > *CurSMT* * *cores* **then**
4. *PBDST.profiler.start()*;
5. **end if**
6. **end if**
7. **end if**

Algorithm 2 Algorithm used to trigger the PBDST hooks in *CoarseGrainedSchedulerBackend*

Input: Variables:

Current SMT configurations: *CurSMT*;

Task state: *state*;

Current Stage Id: *curStageId*;

Last Stage Id: *lastStageId*;

1. **if** *isFinished(state)* **then**
2. **if** *lastStageId* != *curStageId* **then**
3. *lastStageId* = *curStageId*
4. **if** *profilerIsRunning()* **then**
5. *PBDST.profiler.stop()*;
6. *predictSMT* = *PBDST.prediction(CurSMT)*;
7. **if** *predictSMT* != *CurSMT* **then**
8. **for** *executor* in *executorMap* **do**
9. *executor.freeCores.updating(predictSMT)*;
10. **end for**
11. **end if**
12. **end if**
13. **end if**
14. **end if**

First Schema: performs multiple predictions, each stage uses its locally optimal SMT configuration

---frequently threading-switching, rebalanceing the thread distribution, prediction. All of those will cause performance penalty

Second Schema: predicts the best SMT configuration in the first stage and keep it unchane in all stages

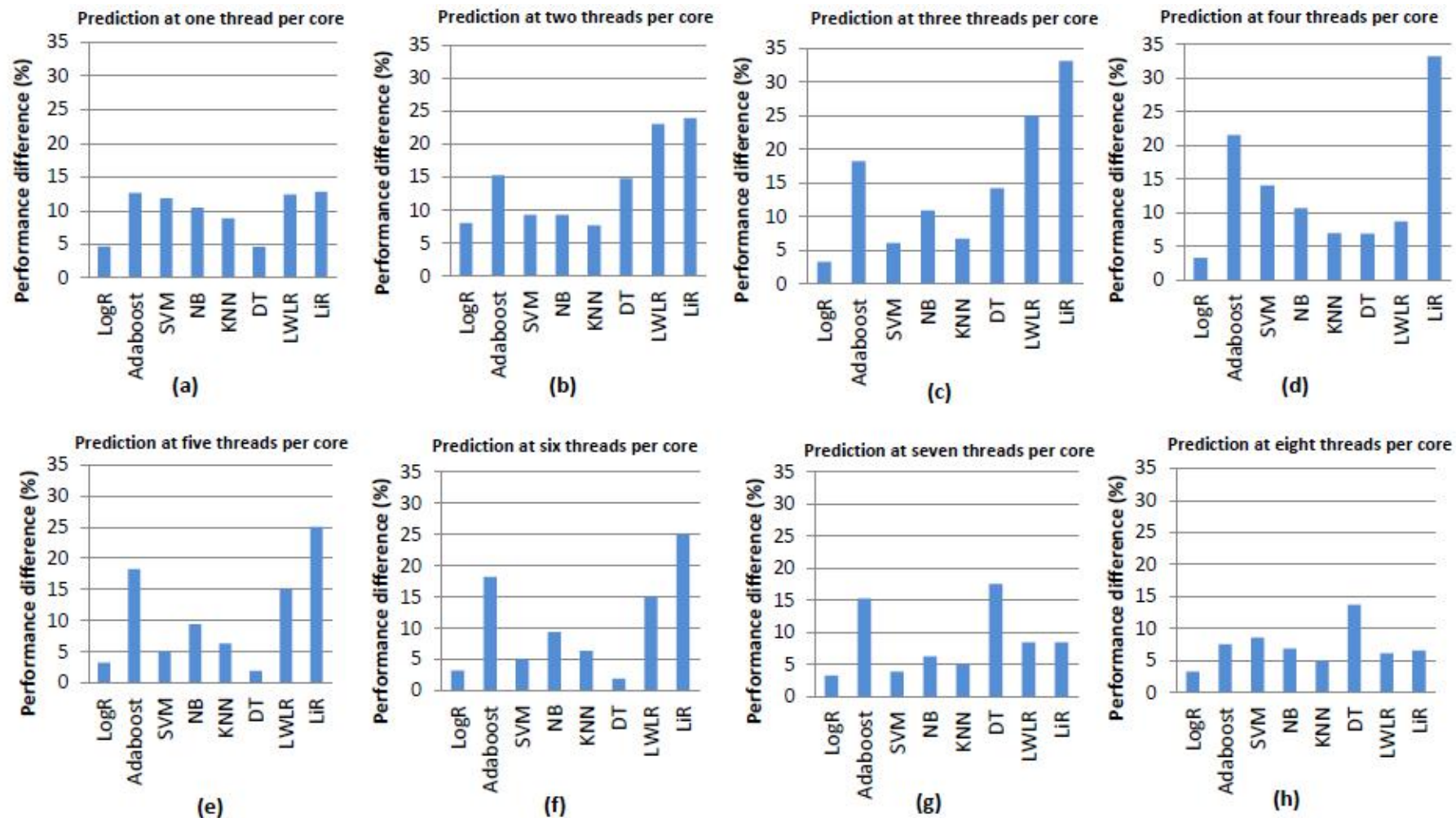


Figure 5: The average static prediction errors of different predictors with diverse SMT configurations.

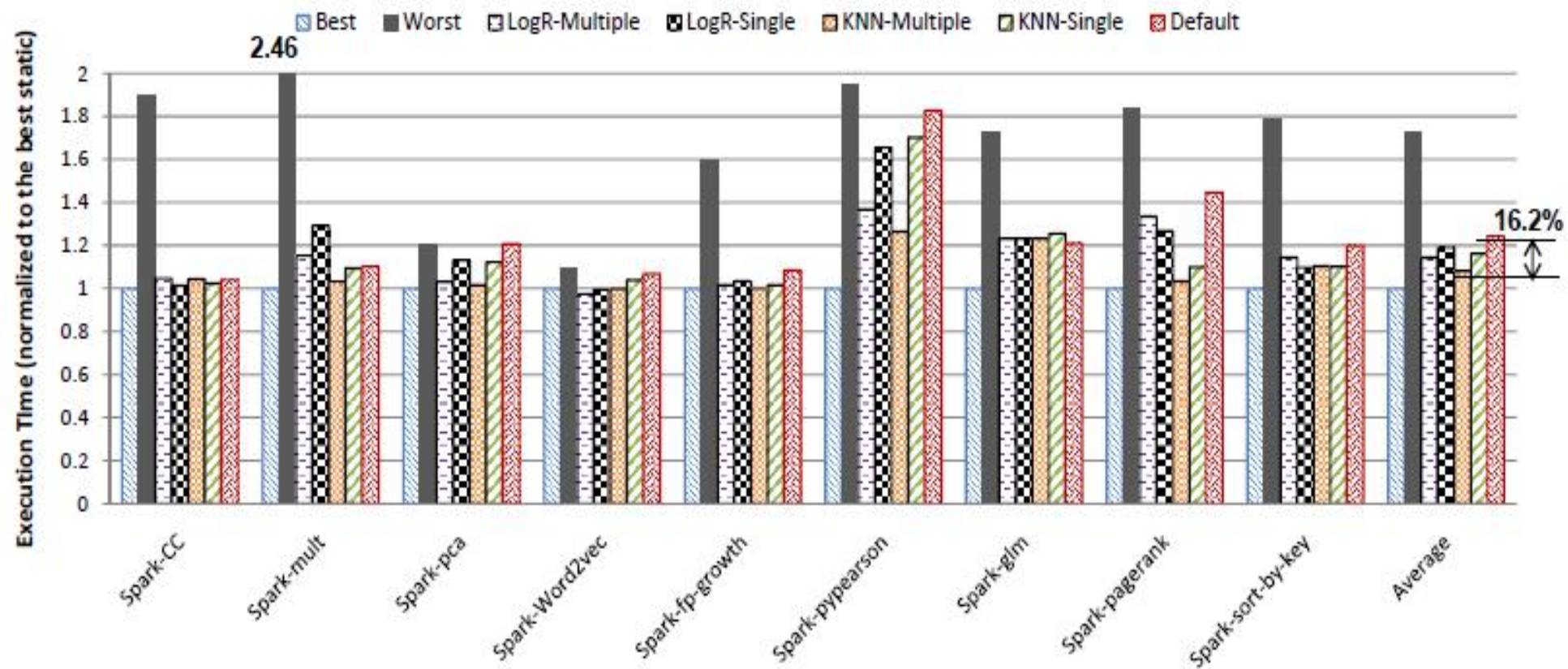


Figure 6: Effects of using different machine learning algorithms when perform automatic predictions.

PBDST framework can achieve avarage perfomance improvement ranging from 5.2% to 16.2% for different predictions with different prediction schemes.

