

Decision Tree Ensembles

INTRODUCTION

Gradient boosting frameworks such as XGBoost, LightGBM and CatBoost, as well as Random Forest algorithms are widely used techniques in recommender systems, search engines and payment platforms.

XGBoost, LightGBM, CatBoost, and Random Forest algorithms are based on learned decision tree ensembles. Such decision trees are fed with training data in order to teach them to ask the right questions about a data set: For example, if the decision tree shall predict whether a user will like a certain movie recommended to him on a website, the tree learns which features of the movie (*i.e.* the data set) are relevant to the user. After the training phase, new data are applied to the decision trees in order to make predictions autonomously without human interaction (**inference phase**).

In applications in which the request rate can amount to thousands of simultaneous predictions, there are two performance metrics in addition to the prediction accuracy:

- the sustained throughput (simultaneous queries)
- the response time of a single query

Xelera provides an accelerator software which offloads the inference phase to data center-grade Field-Programmable Gate Arrays (FPGAs) in order to increase the throughput at a guaranteed query response time.

KEY BENEFITS

- **Acceleration:** One order of magnitude throughput over CPUs and GPUs
- **Cost savings:** One order of magnitude cost savings over CPUs and GPUs
- **Integration:** Integration with standard machine learning frameworks, usable with zero code change

SOLUTION OVERVIEW

The Decision Tree Accelerator speeds up the inference phase of gradient boosting trees and random forest algorithms. It works with models created with **XGBoost**, **LightGBM**, **Scikit Learn**, and **H2O.ai**. The software allows data scientists and engineers to build fast, scalable and cost-efficient machine learning infrastructure, and it does **not require changes in the use of the machine learning frameworks**.

The Decision Tree Accelerator uses the fine-grained parallelism of FPGAs to execute the machine learning models significantly faster than on CPUs or GPUs. The acceleration of the machine learning inference enables more throughput per server node compared to running the frameworks without acceleration.

XELERA TECHNOLOGIES GMBH

RHEINSTR. 40-42
64283 DARMSTADT
GERMANY

FON: +49 6151 162 24 25
CELL: +49 175 268 17 56

MAIL: SALES@XELERA.IO
WEB: WWW.XELERA.IO

MANAGING DIRECTOR:
FELIX WINTERSTEIN

REGISTERGERICHT (LEGAL DOMICILE):
DARMSTADT HRB 97805

IBAN:
DE 17 5085 0150 0000 7709 90
BIC: HELADEFIDAS

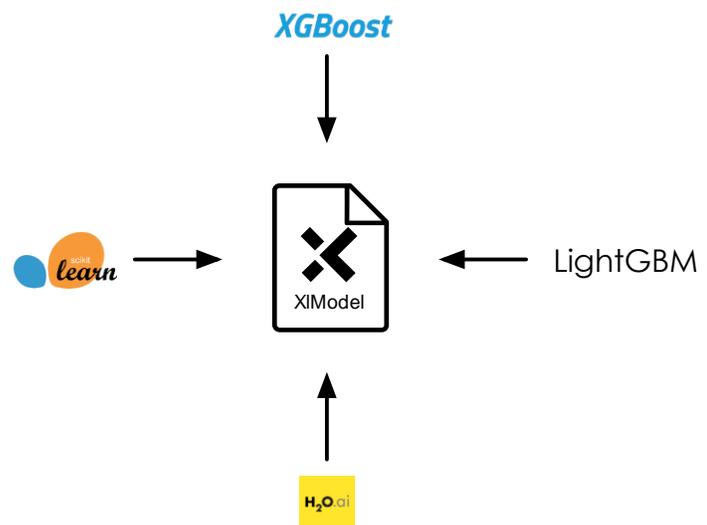
UST-ID-NR.: DE318360521

SOLUTION BRIEF

XELERA

ANALYTICS

- High-throughput Random Forest, XGBoost and LightGBM Inference
- One order of magnitude throughput over CPUs and GPUs
- One order of magnitude cost savings over CPUs and GPUs
- Usable with zero code change

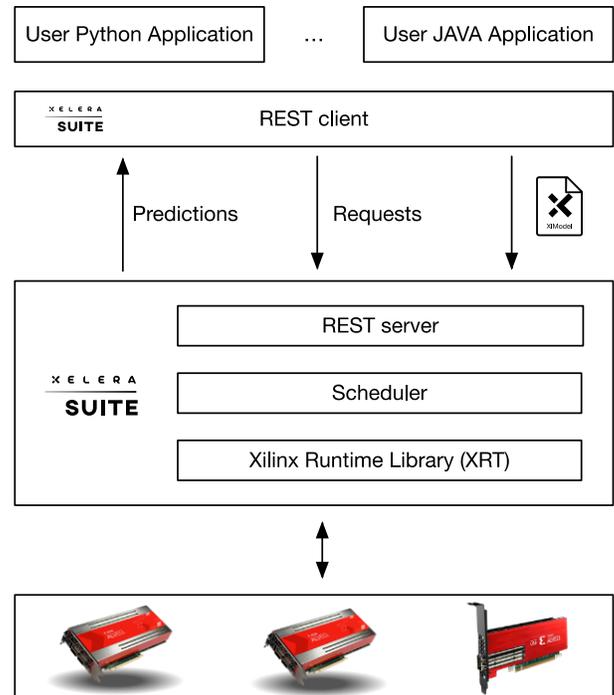


Decision Tree Ensembles

SOLUTION DETAILS

The Decision Tree Inference Accelerator consists of two parts:

1. The **Model Compiler** converts trained gradient boosting machine and trained random forest models from **XGBoost**, **LightGBM**, **Scikit Learn**, and **H2O.ai** automatically into a unified model format (XIModel).
2. The **Acceleration Software**, based on Xelera Suite, loads the compiled XIModel and executes it on Xilinx Alveo U200, U250 and U50 platforms, and on Nimbix and AWS F1 cloud instances. The user can either send inference requests to a REST API and receive the corresponding predictions directly or use a Python or Java interface that is compatible with the previously listed machine learning frameworks and that does make the requests internally, requiring zero code change. The Scheduler dispatches the requests to one or several FPGAs and splits the workload across several server nodes if needed. The scheduling and FPGA access is transparent to the user. Knowledge of FPGAs is not required to use the accelerator.

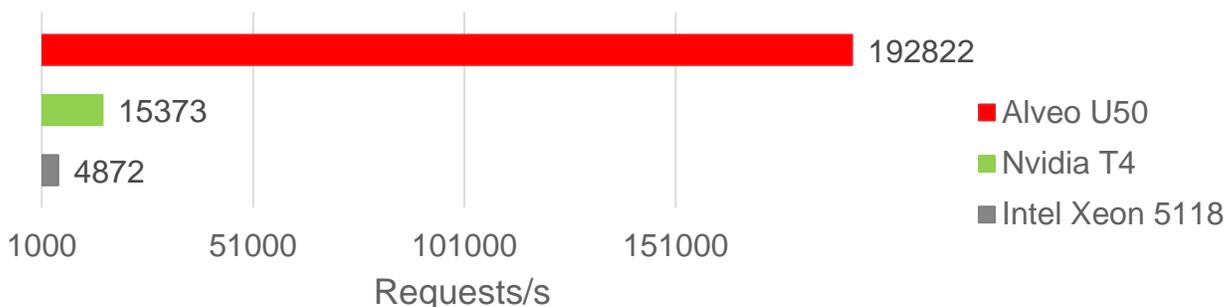


RESULTS

The graph below shows a throughput comparison for an XGBoost regression between FPGA, GPU and CPU platforms, measured using a publicly available data set. The benchmark is exemplary in that the throughput depends on the parameters of the input data set, the algorithm, and the machine learning model.

Data set: Flight (https://www.transtats.bts.gov/OT_Delay/OT_DelayCause1.asp) | Number of features: 10 | Maximum tree depth: 8 levels | Number of trees: 300

- FPGA acceleration software: Xelera Suite
- GPU acceleration software: Nvidia RAPIDS CuML
- CPU: XGBoost (no hardware acceleration)



TAKE THE NEXT STEP

Request a free trial (on-premises or public cloud): <https://xelera.io/product/demo-license-requests>
 Contact us on info@xelera.io to get in touch.