

# Rapids Data

A BORRUI DATA COMPANY

## SPECIAL REPORT



## Parallelized R In-Memory Reduces Execution Time and Analyzes Large Data Sets

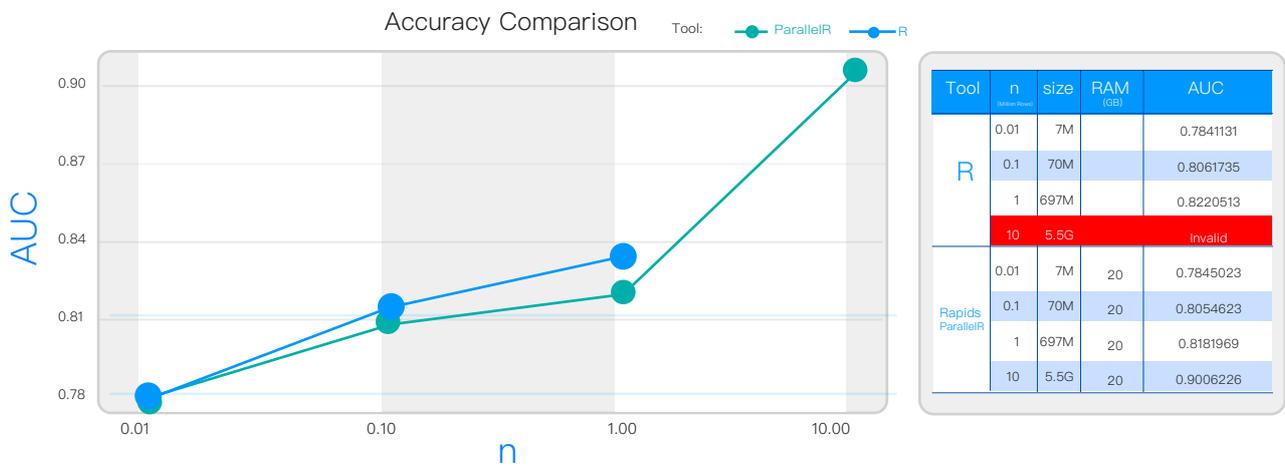
Big data analytics is the process of examining large and complex datasets. Since the beginning of the digital era, data volume has been rising exponentially. As a result, the data that we are dealing with today often exceeds the computational capabilities of some of the big data analytics platforms. The good news is that R, as a leading programming language of data science, consists of powerful functions to tackle problems related to big data processing.

However, R is RAM based. R running on a single node can only perform efficient data analysis on small datasets. The problem is that the computational requirements of analyzing tens or hundreds of GBs of data can exceed the capabilities of a single node. This leads to large up-front investments in high-performance computing infrastructures.

RapidsDB is a fully parallel, distributed, in-memory federated query system designed to support complex analytical SQL queries. Rapids Data's distributed R computing engine breaks through the single-machine restriction commonly encountered in the big data industry. The Rapids ParallelR module of the Rapids Data Platform (RDP) is a distributed, parallel implementation of the R language and the R operating environment. It is integrated with RapidsDB clusters and takes advantage of scalable computing. Memory and computation can be aggregated across hundreds or thousands of nodes whenever a company is ready to take on the horizontal infrastructure expansion.

RapidsDB integrated with ParallelR can be deployed either in Cloud or on premises. With the R module computation happening directly with the RapidsDB database, the usual complicated data uploading or cleansing processes can be eliminated. Rapids Data's internal test results showed that the modeling and the computation performance of ParallelR surpassed that of stand-alone R both on small and large scale datasets. When the data volume reached 5.5GB, the stand-alone R completely stopped while performing the Random Forest algorithm whereas ParallelR could still easily process gigabytes of data. Furthermore, as the data volume continued to increase, the accuracy of the predictive model increased substantially. A different test showed that traditional R took 1,086.90 seconds to load 12GB of data while a 3-node ParallelR cluster only took 105 seconds for the same amount of data.

# Random Forest Performance Comparison



Besides the powerful data analysis functions of the R language, R includes many packages that are exclusively built for simplifying machine learning. Rapids ParallelR also delivers these original R packages to enable users to apply machine learning against data that is managed by RapidsDB. Currently, Rapids ParallelR supports twenty popular algorithms in four categories listed in the table below.

Statistical Analysis	<ul style="list-style-type: none"> <li>Linear Models (GLM)</li> <li>COX (Cox Proportional Hazards)</li> <li>Naive Bayes</li> </ul>
Ensembles	<ul style="list-style-type: none"> <li>Random Forest</li> <li>Distributed Trees</li> <li>Gradient Boosting Machine</li> <li>R Package - Super Learner Ensembles</li> </ul>
Deep Neural Networks	<ul style="list-style-type: none"> <li>Multi-layer Feed- Forward Neural Network</li> <li>Auto-encoder</li> <li>Anomaly Detection</li> <li>Deep Features Clustering</li> <li>Restricted Boltzmann Machines</li> <li>K-Means</li> </ul>
Dimension Reduction	<ul style="list-style-type: none"> <li>Principal Component Analysis</li> <li>Generalized Low Rank Models Solvers &amp; Optimization</li> <li>Generalized ADMM Solver</li> <li>L-BFGS (Quasi Newton Method)</li> <li>Ordinary Least- Square Solver</li> <li>Stochastic Gradient Descent</li> <li>Word2vec</li> </ul>

According to Google, the rise of cloud computing has opened up new opportunities for R. Google recently announced that R scripts can now run on the Google Cloud Platform for modeling and analytics. As an industry leader in big data real-time processing, the launch of Rapids ParalleIR incorporates the philosophy of utilizing R's powerful data wrangling, analysis, modeling, visualization and statistical analysis capabilities to break the restrictions posed by data size. It helps enterprises build large-scale models and gain more real-time insights from their big data business assets in an economical way.

For more information about Rapids ParalleIR, please visit <http://rapidsdb.com/products/rapids-parallelr>



Rapids Data is an industry leader in big data real-time processing and analytics. Rapids Data one-stop big data analysis platform (RDP) comprises the following five application products and provides a full range of real-time big data support for your business.



- **RapidsDB:** standard-based, in-memory big data analysis processing platform
- **Rapids Hadoop:** ultra fast SQL-On-Hadoop technology
- **Rapids StreamDB:** ISO SQL standard-based, in-memory streaming data analytics platform
- **Rapids ParalleIR:** all in-memory real-time big data machine learning algorithm library
- **Rapids DBaaS Cloud System:** composed of a series of physical hosts on the architecture to form an operational resource pool that is centrally managed and regulated by the management server

For more information, please contact: [info@rapidsdata.com](mailto:info@rapidsdata.com)



@RapidsDB