# Statistical & Causal Analysis of the Chimbuko Provenance Database

Author(s) :Serges Love Teutu Talla, Margaret O Ajuwon, Isabelle Kemajou-Brown, Christopher Kelly, Kerstin Kleese Van Dam

## BACKGROUND

Performance data are collected to establish how well Exascale applications are doing with executing their workflow efficiently. Chimbuko collects performance anomalies that are being detected and saved them into its Provenance Database, together with as much contextual information as needed and will be used for our analysis.

## GOAL AND OBJECTIVES

Develop a set of algorithms to query data, perform analysis and visualization using Python to determine if the information collected for each anomaly is sufficient to conduct causal analysis.

## METHODOLOGY

Performed correlation analysis using Theil's U correlation method, applied machine learning by regression for prediction and K-Prototypes for clustering, and ran Random Forest Model and Decision Tree for feature selection


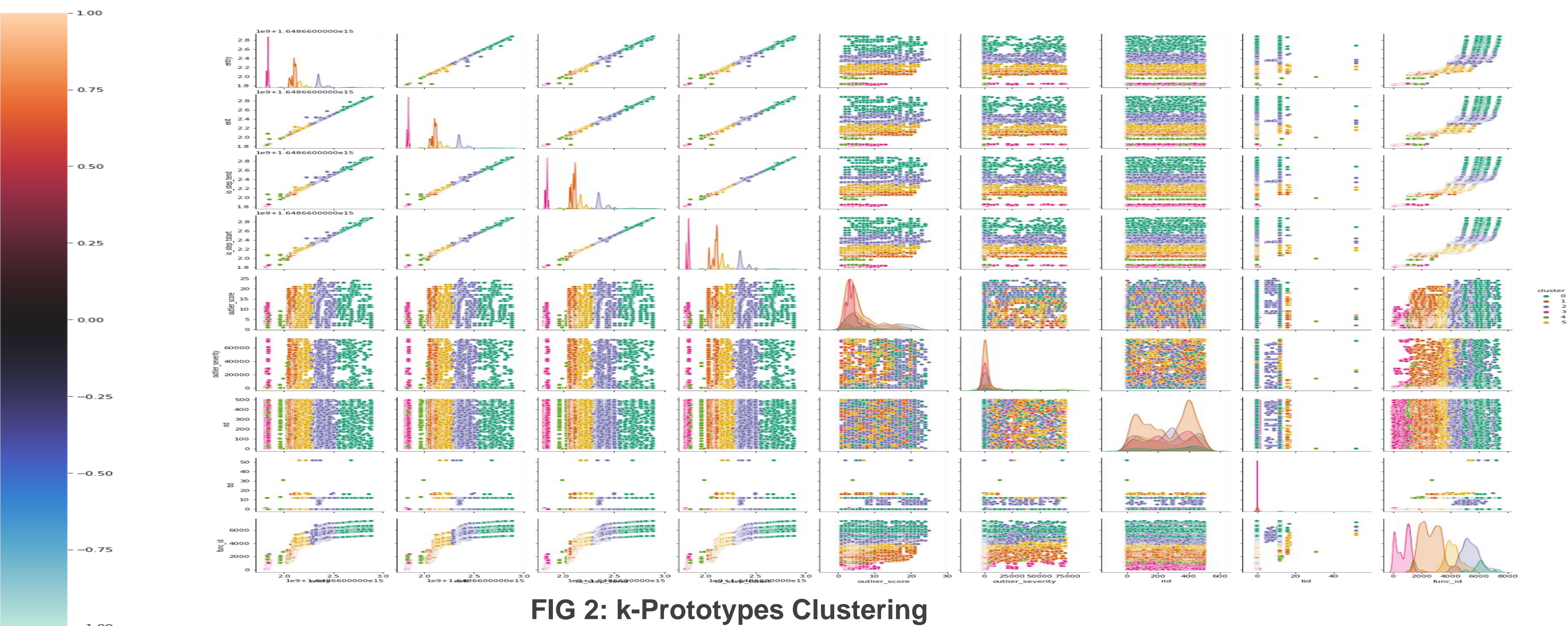
**FIG 1: Correlation Coefficients**



**FIG 2: k-Prototypes Clustering**



**FIG 3: Binary Logistic Regression (BLR) Results**



**FIG 4: BLR With Random Forest Model & Decision Tree**

## RESULTS AND CONCLUSION

- The correlation coefficient of call_stack is above 0.93 with all features except for rid (rank index), and outlier_score with the values of 0.49 and 0.61 respectively. This implies that call_stack is a very important feature of the dataset and could highlight interesting trends.
- From the scatter plot, we can conclude that rid are spread out on the entire dataset. This result suggests there is no apparent inhomogeneity between ranks in the job.
- Basic grouping performed by sum of outlier_severity shows that 'MPI_Allreduce()' has the highest sum of outlier_severity and the top five call_stack with highest sum of outlier_severity were all ending with 'OpenMP_Implicit_Task' . The outlier_sevrerity reflects the importance of an anomaly.
- Regression analysis performed before feature selection (fig 3) indicates that there is a strong evidence of association between the outcome (outlier_score) and predictors: is_gpu_event, tid, runtime_total and outlier_severity.
- Feature selection was performed on selected variables (entry, is_gpu_event, outlier_severity, rid, runtime_total, and tid) using Random Forest and Decision Tree and, showed that outlier_severity and entry are the only important features. The regression analysis results (fig 4) indicate that there is a strong evidence of association between the outlier_score and the outlier_severity.

## ACKNOWLEDGEMENT

## REFERENCE

https://chimbuko-performance-analysis.readthedocs.io/en/ckelly_develop/
https://github.com/margaretajuwon/Statistical-and-Causal-Analysis-of-Chimboku-Provenace-Database

U.S. DEPARTMENT OF ENERGY

ECP EXASCALE COMPUTING PROJECT

www.bnl.gov

SUSTAINABLE HORIZONS INSTITUTE

Brookhaven National Laboratory