# Summary of NSF 1443054: CIF21 DIBBs: Middleware and High-Performance Analytics Libraries for Scalable Data Science project

The NSF 1443054: CIF21 DIBBs: Middleware and High-Performance Analytics Libraries for Scalable Data Science project [1,2] led by Indiana University had two foundational concepts. The first was the big data Ogres work with Indiana University and the NIST Public Big Data Working Group that collected 51 use cases – each with 26 properties [3]. The Ogres were a set of 50 features that categorized applications and allowed one to identify common classes such as Global GML and Local LML Machine Learning. GML is highly suitable for HPC systems while the very common LML and MapReduce categories also perform well on more commodity systems. As another example, "Streaming" [4] was a feature seen in 80% of the applications. The second foundation was the High-Performance Computing enhanced Apache Big Data Stack HPC-ABDS [5-7]) which built systems from the commodity open source Apache software enhanced by HPC high-performance computing as necessary as shown in Fig. 1. We have developed major HPC enhancements to ABDS software including Harp [8, 9] based on Hadoop and Twister2 [10] based on Heron, Spark and Flink for both batch and streaming scenarios.
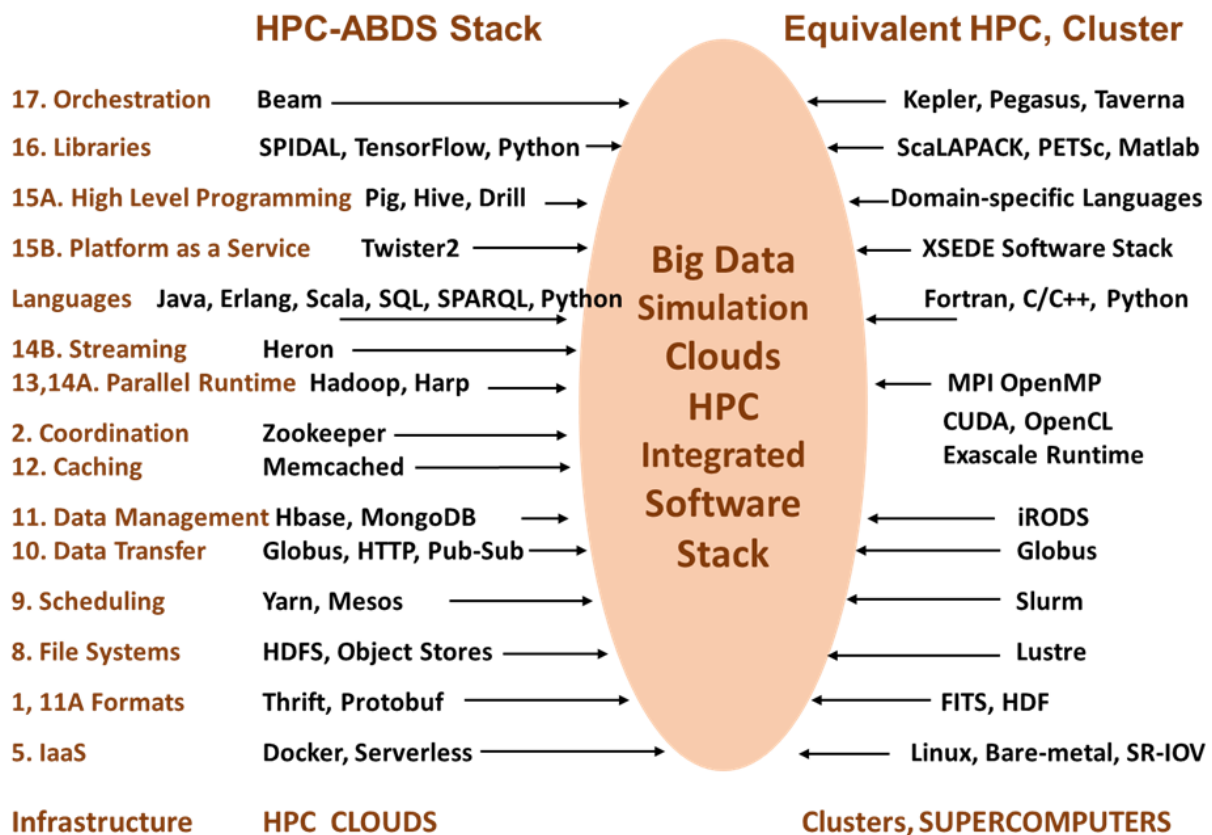


Figure 1: Proposed HPC-ABDS software stack compared to a classic HPC cluster approach. The levels 1-17 correspond to those introduced in [7]. The software systems on the left are standard tools (mainly from Apache augmented by Harp [8] and Twister2 [10] HPC enhancements.

The high-performance, scalable data analytics library SPIDAL has 4 components: a) core library approximating scope of Apache Mahout; b) parallel graph algorithms; c) analysis of biomolecular simulations (high-performance versions of existing libraries from Utah and Arizona State) and d) image processing. The table lists existing routines in the core area, and over next year we will be packaging and adding documentation and tutorials for all components.

| |
|---|
| • DA-MDS Multidimensional scaling Rotate, AllReduce, Broadcast |
| • Directed Force Dimension Reduction AllGather, Allreduce |
| • Irregular DAVS Clustering Partial Rotate, AllReduce, Broadcast |
| • DA Semimetric Clustering Rotate, AllReduce, Broadcast |
| • K-means AllReduce, Broadcast, AllGather DAAL |
| • SVM AllReduce, AllGather |
| • SubGraph Mining AllGather, AllReduce |
| • Latent Dirichlet Allocation Rotate, AllReduce |
| • Matrix Factorization (SGD) Rotate DAAL |
| • Recommender System (ALS) Rotate DAAL |
| • Singular Value Decomposition (SVD) AllGather DAAL |
| • QR Decomposition (QR) Reduce, Broadcast DAAL |
| • Neural Network AllReduce DAAL |
| • Covariance AllReduce DAAL |
| • Low Order Moments Reduce DAAL |
| • Naive Bayes Reduce DAAL |
| • Linear Regression Reduce DAAL |
| • Ridge Regression Reduce DAAL |
| • Multi-class Logistic Regression Regroup, Rotate, AllGather |
| • Random Forest AllReduce |
| • Principal Component Analysis (PCA) AllReduce DAAL |

Table: Current members of core SPIDAL High-Performance Library: DAAL implies integrated with Intel DAAL Optimized Data Analytics Library and so running well on KNL architecture. These use Map-Collective paradigm and collectives used are listed.

Note current target architectures are clusters with either Haswell or Knights Landing (KNL) nodes. We can also extend libraries to GPU's.

**References**
1. Digital Science Center. SPIDAL Home Page: CIF21 DIBBs: Middleware and High-Performance Analytics Libraries for Scalable Data Science - Scalable Parallel Interoperable Data Analytics Library [Internet]. 2015. Available: http://spidal.org/index.html
2. Digital Science Center and SPIDAL Collaboration. DSC-SPIDAL Github repository. 2017 [accessed 2017 April 7]; Available from: https://github.com/DSC-SPIDAL.
3. NIST Big Data Public Working Group: Use Cases and Requirements Subgroup, NIST Big Data Interoperability Framework: Volume 3, Use Cases and General Requirements (NIST Special Publication 1500-3). 2016, NIST: Vol. 3. http://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1500-3.pdf. DOI: http://dx.doi.org/10.6028/NIST.SP.1500-3.

4. Fox G, Jha S, Ramakrishnan L. STREAM2016: Streaming Requirements, Experience, Applications, and Middleware Workshop Workshop Final Report [Internet]. 2016. doi:10.2172/1344785. Also http://streamingsystems.org
5. Geoffrey Fox, Judy Qiu, Shantenu Jha, Saliya Ekanayake, Supun Kamburugamuve. White Paper: Big Data, Simulations, and HPC Convergence. BDEC Frankfurt workshop. 2016. doi:10.13140/RG.2.1.3112.2800
6. HPC-ABDS Kaleidoscope of over 350 Apache Big Data Stack and HPC Technologies [Internet]. Available: http://hpc-abds.org/kaleidoscope/
7. Geoffrey Fox, Judy Qiu, Shantenu Jha, Supun Kamburugamuve, Andre Luckow. HPC-ABDS High-Performance Computing Enhanced Apache Big Data Stack. Invited talk at 2nd International Workshop on Scalable Computing For Real-Time Big Data Applications (SCRAMBL'15) atCCGrid2015, the 15th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing. IEEE; 2015. Available: http://dsc.soic.indiana.edu/publications/HPC-ABDSDescribedv2.pdf
8. Chen L, Peng B, Zhang B, Liu T, Zou Y, Jiang L, et al. Benchmarking Harp-DAAL: High-Performance Hadoop on KNL Clusters. IEEE Cloud 2017 Conference. IEEE; Available: http://dsc.soic.indiana.edu/publications/2017CLOUDResearchTrack_12261.pdf
9. Peng B, Zhang B, Chen L, Avram M, Henschel R, Stewart C, et al. HarpLDA+: Optimizing Latent Dirichlet Allocation for Parallel Efficiency [Internet]. Indiana University, Digital Science Center; 2017 Aug. Available: http://dsc.soic.indiana.edu/publications/HarpLDA%2B%20Optimizing%20Latent%20Dirichlet%20Allocation%20for%20Parallel%20Efficiency.pdf
10. Kamburugamuve S, Fox G. Designing Twister2: Efficient Programming Environment Toolkit for Big Data [Internet]. Digital Science Center; 2017 Aug. Available: http://dsc.soic.indiana.edu/publications/Twister2.pdf