

---

## Large Scale DNA Sequence Analysis and Biomedical Computing using MapReduce, MPI and Threading

---

*Judy Qiu*  
(Indiana U.)

### *Abstract*

Many areas of science are seeing a data deluge coming from new instruments, myriads of sensors and exponential growth in electronic records. We take two examples – one the analysis of gene sequence data (up to 35000 Alu sequences) and other a study of medical information (initially over 100,000 patient records) in Indianapolis and their relationship to Geographic and Information System and Census data available for 635 Census Blocks in Indianapolis. We look at initial processing (such as Smith Waterman dissimilarities), clustering (using robust deterministic annealing) and Multi Dimensional Scaling to map high dimension data to 3D for convenient visualization. We show how scaling pipelines can be produced that depending on data set size, can either use multicore laptop, desktop clients, supercomputers or modest clusters for the computer intensive sections.

We contrast classic concurrency techniques- namely MPI or multi-threading with emerging cloud technologies Dryad (Microsoft) and Hadoop (Yahoo). The latter are especially appropriate at the initial stages when computations are largely independent. However Hadoop and Dryad can prepare data for sophisticated MPI data analysis. We present performance results from Tempest – An Infiniband connected 32 node system running Windows HPCS with each node having 24 cores spread over 4 Intel chips. Such a modest cluster can fully process all stages of the 35,000 element Alu study in less than a day and is suitable for up to 200,000 sequences even though all steps in analysis are of  $O(N^2)$  time complexity. We discuss ease of use and programmability as well as wall clock execution time. We intend to make these capabilities available as services supported on virtual clusters. A summary of our recent work is available at <http://grids.ucs.indiana.edu/ptliupages/publications/CetraroWriteupJune11-09.pdf>.

Dryad and Hadoop are new technologies that as they mature will become essential parts of such systems. MPI naturally supports very efficiently many aspects of the MapReduce model underlying cloud technologies. However it lacks essential fault tolerance and flexibility of MapReduce. We note that extensions of MapReduce can efficiently support all parts of analysis. One needs to support iterative MapReduce with low overhead streaming communication. All these technologies scale to large systems and suggest scalable data analysis will be possible for the rapidly growing field of biomedical computing.

---