

How to use BigDataBench workloads and data sets

Gang Lu

Institute of Computing Technology,
Chinese Academy of Sciences

BigDataBench Tutorial
MICRO 2014 Cambridge, UK



中国科学院
INSTITUTE OF COMPUTING TECHNOLOGY

General steps to use BigDataBench

■ Current release

- Version 3.1 on <http://prof.ict.ac.cn/BigDataBench>

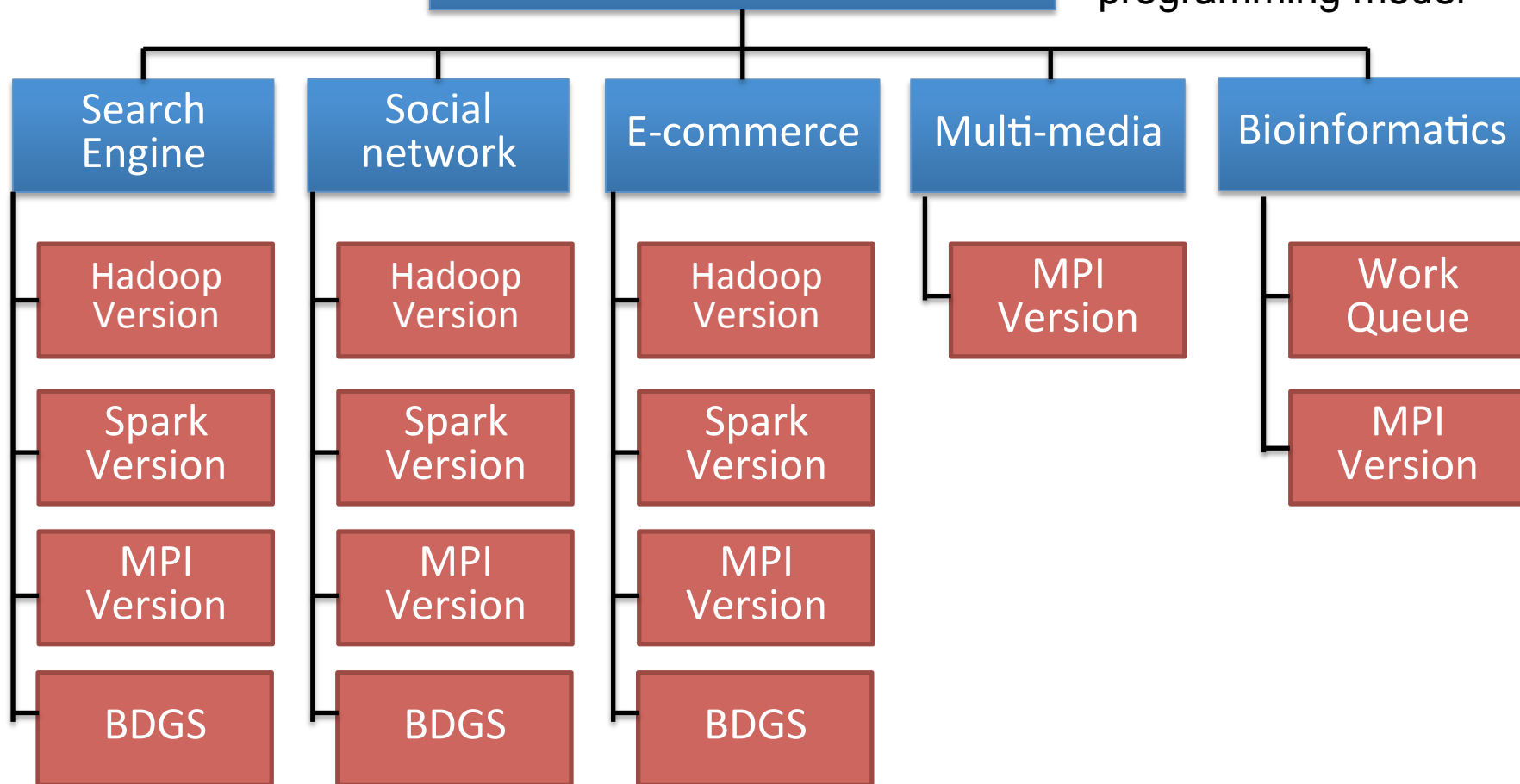
■ General steps to run the benchmarks

- Prepare the package of BigDataBench
- Prepare the environments of the selected software stack
- Generate data sets as you need
 - *You can find a genDate* or a prepare* shell script in each directory of the benchmarks*
- Run the scripts or commands (Handbook!)

A glance of the directory structure

Root directory of BigDataBench

With different programming model



Domain: Search Engine

Mirobenchmarks

Operations or Algorithms	Types	Data Sets	Software Stacks
Read	Cloud OLTP	ProfSearch Resumes	HBase, Mysql
Write			
Scan			
Sort	Offline analytics	Wikipedia	Hadoop, Spark, MPI
Grep			
WordCount			
Index	Offline analytics	Google Web Graph	Hadoop, Spark, MPI
PageRank			
Nutch Server	Online service	SouGou Index	Nutch

Example: Cloud OLTP with HBase (Hadoop Version)

- Target: run “write” operations using HBase
- General steps:
 - Prepare HBase according to the office guide
 - `sh /hbase-0.94.5/bin/hbase shell`
 - `create 'usertable','f1','f2','f3'`
 - Prepare YCSB as the workload generator
 - YCSB is in the directory of BasicDatastoreOperations/ycsb-0.1.4
 - Run YCSB commands like this:
 - `sh bin/ycsb load hbase -P workloads/workloadc -p threads=<thread-numbers> -p columnfamily=<family> -p recordcount=<recordcount-value> -p hosts=<hostip> -s>load.dat`

Example: Cloud OLTP with Hbase (Hadoop Version)

■ Important parameters of running YCSB:

<threadnumber>	The number of client threads, this is often done to increase the amount of load offered against the database.
<family>	In the HBase case, we used it to set database column. You should have database <i>user</i> table with column <i>family</i> before running this command. Then all data will be loaded into database <i>user</i> table with column <i>family</i>
<recordcount-value>	The total records for this benchmark. For example, when you want to load 10G record, you should set it to 10000000.
<hostip>	The IP of the HBase's master node.

Example: PageRank with MPI

- Target: run PageRank using MPI
- General steps:
 - Prepare MPI environments
 - Run the data generation script
 - `cd BigDataBench_MPI_V3.0/SearchEngine/MPI_Pagerank`
 - `sh genData_PageRank.sh`
 - Run the script:
 - `sh run_PageRank.sh <# Iterations of GenGragh >`
 - You can also use the `mpirun` command to run the script on a cluster.
 - `mpirun -x genData_PageRank.sh -n 100 -H inputGraphfile`
 - ><

Steps are almost the same for other programming models.

Refer to the handbook!

Domain: E-commerce

		Operations or Algorithm	Types	Data Sets	Software Stacks
Complex queries	Mirobenchmarks	Bayes	offline analytics	Amazon Movie Review	Hadoop, Spark, MPI
		CF			
		Project	Interactive analytics	CALDA and E-commerce	Hive, Shark, Impala
		Filter			
		Cross Product			
		OrderBy			
	Union				
	Difference				
	Aggregation	Complex queries			
	Join Query				
	Select Query				
	Aggregation Query				

Example: Complex Queries With Shark

■ General steps:

- Prepare Shark environments
- Run the data generation script
 - *cd ./BigDataGeneratorSuite/Table_datagen/*
 - *java -XX:NewRatio=1 -jar pdgf.jar -l demo-schema.xml -l demo-generation.xml -c -s -sf \$number*
 - *Don't forget to upload the output file to HDFS which will be used by Shark tasks*
- Start *Shark* and create three tables which will be used for follow-up queries

Example: Complex Queries With Shark

■ General steps:

- Prepare Shark environments
- Run the data generation script
- Start *Shark* and create three tables which will be used for follow-up queries
 - detailed statements are in the handbook
- Run the queries
 - *sh runQuery.sh*

Domain: Multi-media

Operations or Algorithm	Types	Data Sets	Software Stacks
BasicMPEG	Offline analytics	Stream Data	MPI
SIFT		ImageNet	
Speech Recognition		Audio files	
Ray Tracing		Scene description files	
Image Segmentation		ImageNet	
Face Detection		ImageNet	
DBN		MNIST	

Example: SIFT with MPI

- Target: run SIFT using MPI
- General steps:
 - Prepare MPI environments
 - Run the data generation script
 - *sh getPath \$ ImageNet_1G /BigDataBench_Media*
(The output file will be “ImageNet_1G.path”)
 - Run SIFT using the generated file as input
 - *mpirun -n process_number -f node_file ./siftfeat_mpi*
<input file>

Other domains

Domains	Operations or Algorithm	Types	Data Sets	Software Stacks
Social Network	BFS	offline analytics	Graph500 Data	MPI
	Kmeans		Facebook Social Network	Hadoop, Spark, MPI
	CC		Facebook Social Network	Hadoop, Spark, MPI
Bioinformatics	SAND	offline analytics	Genome sequence Data	MPI
	BLAST		Assembly of the human genome	

Details can be found in the handbook of BigDataBench:
http://prof.ict.ac.cn/BigDataBench_micro_14/



Any
Questions