

SNP detection and Genome Wide Association Study using Hadoop-BAM, CrossBow and Apache HIVE in Hadoop Cluster

Jyotsna Singh – PGDBD201901006

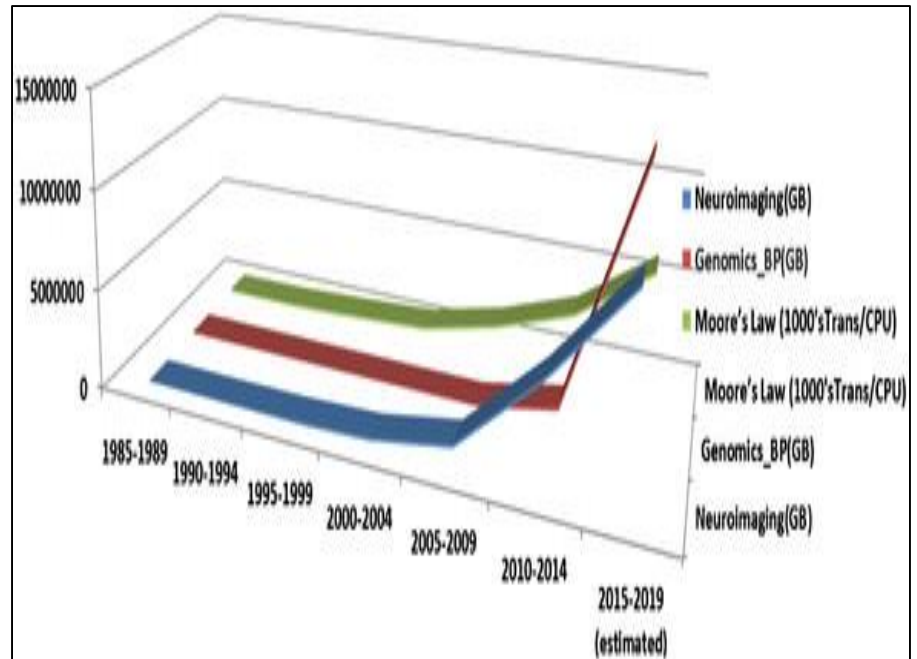
Paulami Das – PGDBD201901009

Poushali Gupta – PGDBD201901010

Institute of Bioinformatics & Applied Biotechnology,
Bangalore

Next Generation Sequencing & Big Data

- The amount of NGS Data worldwide is predicted to double every 5 months which is much faster than Moore's law
- 1000 Genomes project has Petabytes of human genome data sets
- In many GWAS and WGS studies multiple large files have to be processed sequentially



Kryder's law: Exponential growth of neuroimaging and genomics data, relative to increase of number of transistors per chip. By 2025 more than 100 PB of sequenced genome and 1 TB of neuroimaging data will be generated daily.

Different File Formats of Genomic Data

```
@HD VN:1.0 SO:coordinate
@SQ SN:chr20 LN:64444167
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-read-  
edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr  
20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
CCGTGTTTAAAGGTGGATCGGCTACCTTCACAGTAGGCTTAGTGATTCTAGTTGGCCTAGGAATCCAGCTAGTCTGTCTCACTCCCCCTCT
C BBDDCCDDCCDDDDCCDDDDDCDCCDBC?DDDDDDDDDDDDCDCCDDDDDDDDCCCEDDDC?DDDDDDDDDDDDDDDBDHFFFD@@
AS:i:-15 XM:i:3 XO:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0
TGCTGGATCATCTGGTTAGTGCTTCTGACTCAGAGGACCTTCGTCCTCCGGGCAGTGGACCTTCCAGTGATTCCTTGACATAAGGGGCATGGACGA
G DCCCCDEDDDDDDCCDDDDDDCCDDDDDEEC>DFFEJJJJJIGJJJIHGBHHGIJJJJJJGJJJJJJIIHJJJJJHHHHHHFFFFFCCC
AS:i:-16 XM:i:3 XO:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0
GGCTTTATTGGTA AAAAGAAG AATAGC AGATT TAATCAG AAATCC CACTGCC CAGCAGCAC CAACCAG AAAGAAG GGAAGACAGAAAAAACCA
C DDDDDDDCCDDDDDDDEEEEEEEFEFFEGFH HHFGDJJI HJJIJJJ IIIGGFJJ IHHIII JJJJJ IGHHFAHG FHJ FGGHF FDDBB
AS:i:-11 XM:i:2 XO:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0
0 GTGGCTCTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGTGCACTTGGGTCTCCGAAGCAGAAACATCCTCAAATATGACCTCTCG
```

SAM/BAM Files

ReferenceGenome	Homo_sapiens	UCSC	hg19	Sequence	WholeGenome	FASTA
[Regions]						
Name	Chromosome	Start	End	Upstream Probe Length	Downstream Probe Length	
WASHSP-chr1-14363-14829	chr1	14363	14829		0	
WASHSP-chr1-14970-15038	chr1	14970	15038		0	
WASHSP-chr1-15796-15947	chr1	15796	15947		0	
WASHSP-chr1-16607-16765	chr1	16607	16765		0	
WASHSP-chr1-16858-17055	chr1	16858	17055		0	
WASHSP-chr1-17233-17368	chr1	17233	17368		0	
WASHSP-chr1-17606-17742	chr1	17606	17742		0	

BED Files

```
##fileformat=VCF4.2
##INFO=<ID=SVTYPE,Number=1,Type=String,
Description="Type of structure variant">
##INFO=<ID=END,Number=1,Type=Integer,
Description="End position of the variant described in this record">
#CHROM POS ID REF ALT QUAL FILTER INFO
```

1 160929435 rs7520618 G A . SVTYPE=SNP;END=160929436
1 160932043 rs113387749 A . SVTYPE=INS;END=160932043
1 160932206 rs5778188 C . SVTYPE=DEL;END=160932207
1 160932771 rs2256505 A G . SVTYPE=SNP;END=160932772
1 160934077 rs2481074 T A . SVTYPE=SNP;END=160934078
1 160934818 rs1023115 A G . SVTYPE=SNP;END=160934819
1 160935328 . AAA TGC . SVTYPE=SUB;END=160935331
1 160935334 rs75452934 AA TC . SVTYPE=SUB;END=160935336

VCF Files

```
>@HWI-ST216_0180:4:1101:1096:2196#GGCTAC/1
TTTTTCAGNGAATACTGCAAAATCAATAAACTCTTTAG
>@HWI-ST216_0180:4:1101:1158:2236#GGCTAC/1
AAAAGCTCATTTCCTATAGTTAAACAGGACATGCCTT
>@HWI-ST216_0180:4:1101:1448:2211#GGCTAC/1
ATTATATAAGATAGCGGCTTTTTCCGTTAGTTTCCT
>@HWI-ST216_0180:4:1101:1331:2227#GGCTAC/1
CACGTTCTCTGTCCCCAATGGTATTTGCATCCCTGT
>@HWI-ST216_0180:4:1101:1376:2237#GGCTAC/1
GCGTCCCTTAGCTGAAC TACCCAAACGTACGAATGC
```

Fasta Files

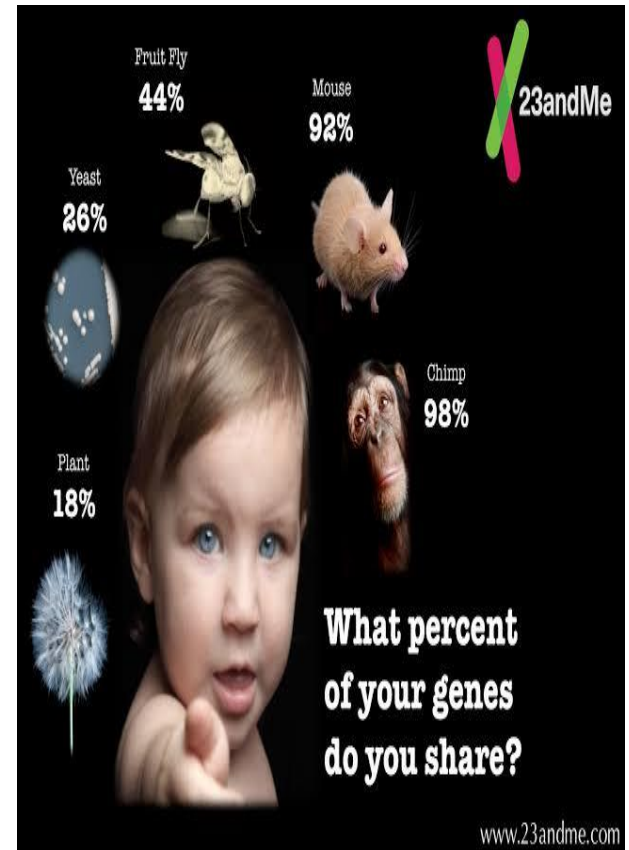
```
@HWI-ST216_0180:4:1101:1096:2196#GGCTAC/1
TTTTTCAGNGAATACTGCAAAATCAATAAACTCTTTAG
+HWI-ST216_0180:4:1101:1096:2196#GGCTAC/1
ceedb]]B[[]]]][fffff\dddddededf fbd
@HWI-ST216_0180:4:1101:1158:2236#GGCTAC/1
AAAAGCTCATTTCTATAGTTAACAGGACATGCCTT
+HWI-ST216_0180:4:1101:1158:2236#GGCTAC/1
gggggggggggggggggggggggfffggggggggggfg
@HWI-ST216_0180:4:1101:1448:2211#GGCTAC/1
ATTATATAAGATAGCGGCTTTTTCCGTTAGTTTCT
```

Fastq Files

Sequence	Source	Feature	Start	End	Score	Strand	Frame	Attributes
chr1	hg19_msk	exon	8386030	8386750	980	-	gene_id	"ENVELE-ef;"
chr1	hg19_msk	exon	152937729	152958046	491	+	gene_id	"MLT1;"
chr1	hg19_msk	exon	242039324	24209678	7249	+	gene_id	"THE1B;"
chr1	hg19_msk	exon	10485712	10489861	8848	-	gene_id	"LTR2C;"
chr1	hg19_msk	exon	53477263	53477536	1005	+	gene_id	"MER2C;"
chr1	hg19_msk	exon	74448477	74449973	9203	-	gene_id	"THE1B-ef;"
chr1	hg19_msk	exon	82872112	82873259	461	-	gene_id	"ManRep1527;"
chr1	hg19_msk	exon	102760404	102760736	1754	+	gene_id	"THE1D;"
chr1	hg19_msk	exon	145751968	145752401	1532	-	gene_id	"LTR47B;"
chr1	hg19_msk	exon	153091484	153092533	4373	-	gene_id	"LTR17E;"
chr1	hg19_msk	exon	160432100	160434254	848	+	gene_id	"LTR85C;"
chr1	hg19_msk	exon	210783724	210784168	1368	-	gene_id	"MLT2B;"
chr1	hg19_msk	exon	211812124	211812364	1068	-	gene_id	"MER21A;"
chr1	hg19_msk	exon	213809477	213910088	1378	-	gene_id	"MLT1F2;"
chr1	hg19_msk	exon	238026629	238027031	816	+	gene_id	"MLT1M;"
chr1	hg19_msk	exon	441772260	441772753	3210	-	gene_id	"LTR24C;"
chr1	hg19_msk	exon	262105	262386	721	+	gene_id	"MER6SD;"
chr1	hg19_msk	exon	3538719	3539068	2001	+	gene_id	"THE1D;"
chr1	hg19_msk	exon	43240484	4325379	1304	-	gene_id	"MLT1B;"
chr1	hg19_msk	exon	51008932	5112709	28994	+	gene_id	"HERVH-ef;"
chr1	hg19_msk	exon	12173599	12174399	808	-	gene_id	"MER2-ef;"
chr1	hg19_msk	exon	12840257	12845900	27617	-	gene_id	"HERVK-ef;"

Properties of Our Data Set

- Semi-Structured
- Reference Genome – Fasta Format
3.2 GB
62743362 bp
- Raw Data – Fastq Format
4.1 GB
27999799 bp
Read length - 150 bp
GC Content – 44%
- Output File – BAM Format, VCF Format
42 GB, 30 KB



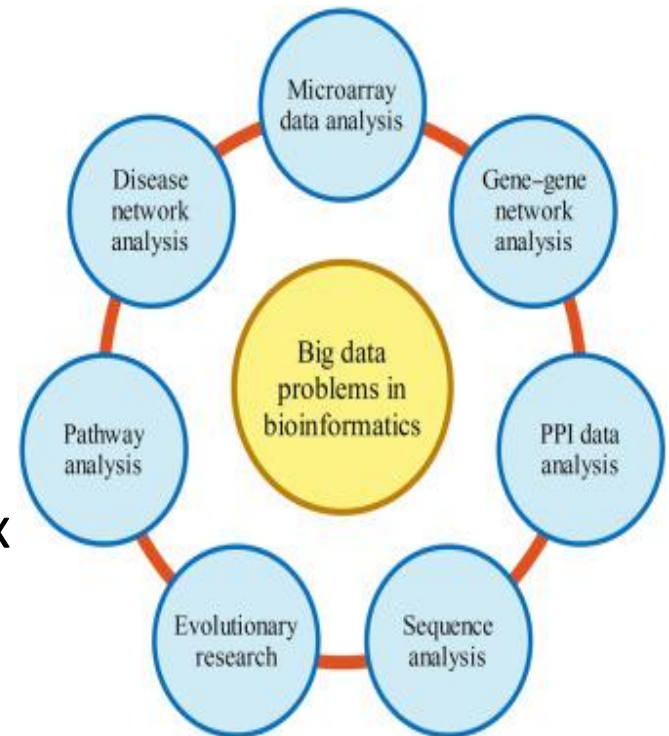
Relationship of Human
Genome with Other
Species

Advantages of Hadoop

- **Hadoop is Open Source distributed data processing system**
- Based on Google's **MapReduce** architecture design
- Cheap commodity hardware for storage
- Fault tolerant distributed filesystems: **HDFS**
- Batch processing systems: **Hadoop MapReduce, Apache Hive, Apache Pig (HDD), Apache Spark (RAM)**
- Parallel SQL implementations for analytics: **Apache Hive, Cloudera Impala, Apache Spark**
- Fault tolerant distributed database: **Hbase**
- **Distributed machine learning libraries, text indexing & search**

Hadoop in Different Biological Aspects

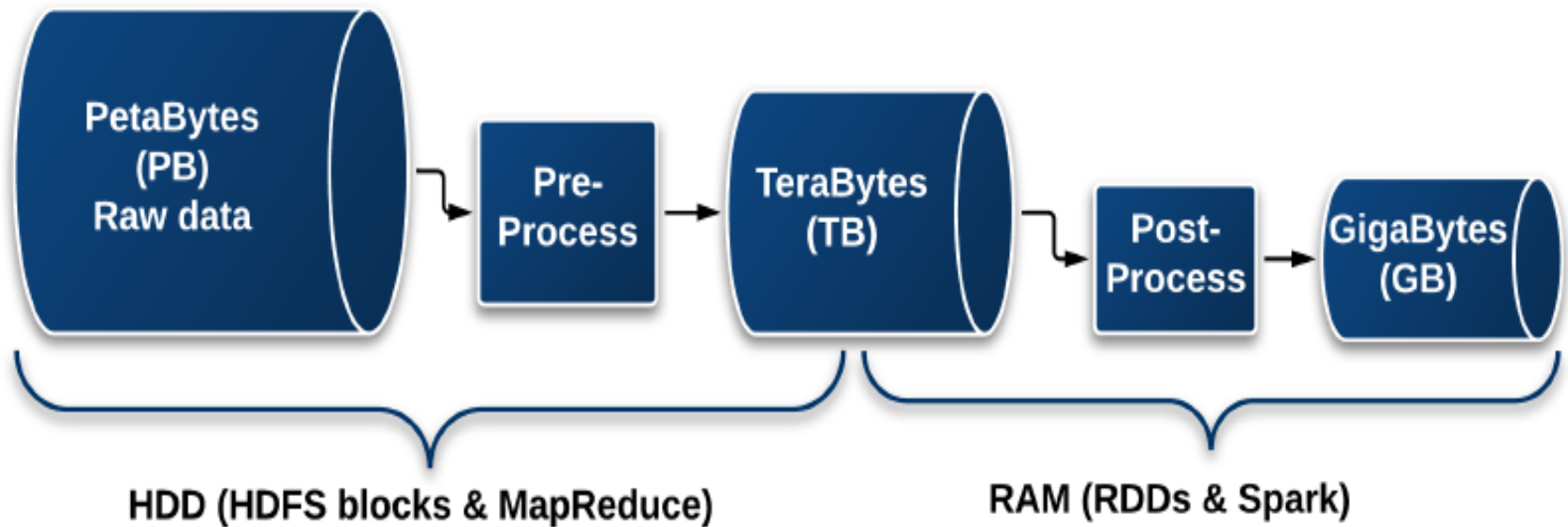
- In Cancer treatments
- In monitoring Patient Vitals
- In the Hospital Network
- In Healthcare Intelligence
- In Structural Bioinformatics –
 - 1) Molecular Docking
 - 2) Clustering of Protein-Ligand complex
 - 3) Structural Alignment
- In Genomic Data Analysis



Tools used in Hadoop For Biological Data Analysis

- **Cloud Burst** – Uses Hadoop as a platform for alignment of short reads.
- **Crossbow** – Uses Hadoop for SNP genotyping from short reads.
- **Contrail** – Uses Hadoop for denovo assembly from short sequencing reads
- **Myrna** – Uses Bowtie and R/Bioconductor for calculating differential gene expression from large RNASeq data sets
- **Cloud Blast** – Uses Gene Set Enrichment Analysis in Hadoop
- **BlueSNP** - Implements GWAS statistical tests in R & executes the calculations with Apache Hadoop using MapReduce formalism.
- **HadoopBAM** – A library for processing NGS data format in parallel with both Hadoop and Spark.
- **Amazon Elastic Compute Cloud & MapReduce.**

Typical Genomics Data Analysis Using HDFS



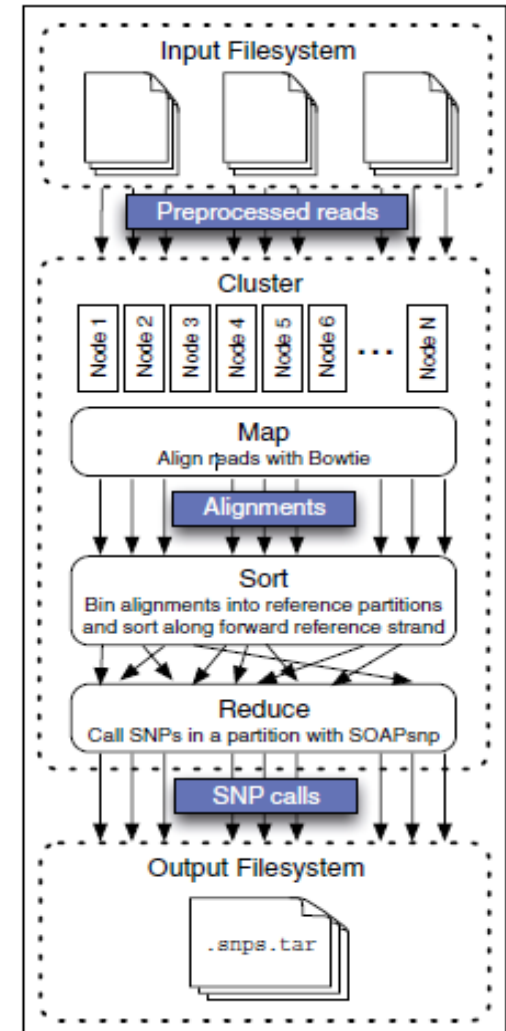
Processing Data in main memory instead of files in hard disks = minimal I/O operations. Map/Reduce data from Petabytes to Gigabytes (million times less in the end)

Project Proposal

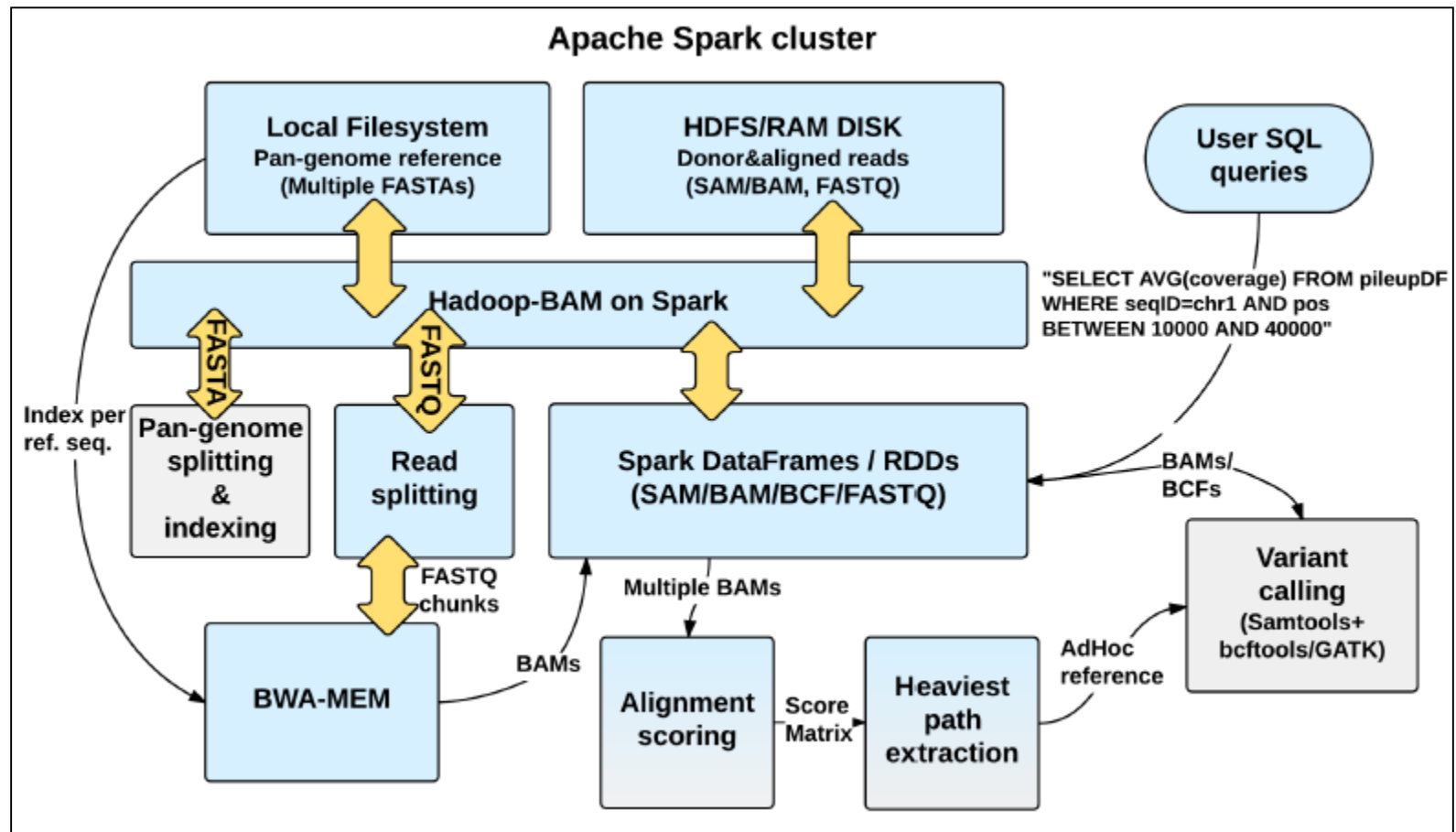
- Genome Alignment using HADOOP-BAM
- SNP Detection using Crossbow (Bowtie+SOAPsnp) and HadoopBAM and comparing both
- Genome Wide Association Study (GWAS) using Apache HIVE across Human Genome of Different Population

SNP Detection using MapReduce Algorithm in Crossbow

- Copying the Fastq raw data and Fasta reference genome from Local File System to HDFS
- Running the Crossbow pipeline in Hadoop Cluster
- Crossbow's Map phase align reads with Bowtie 2 which employs a compact index of reference sequence requiring about 3 GB of memory using HG19
- The index is distributed to all computers in cluster via hadoop file or by instructing each node to independently obtain the index from a shared file
- The reduce phase performs SOAPsnp
- The output of Reduce phase is SNP tuple which stored on the Clustered distributed File System which can be transferred to Local File System.

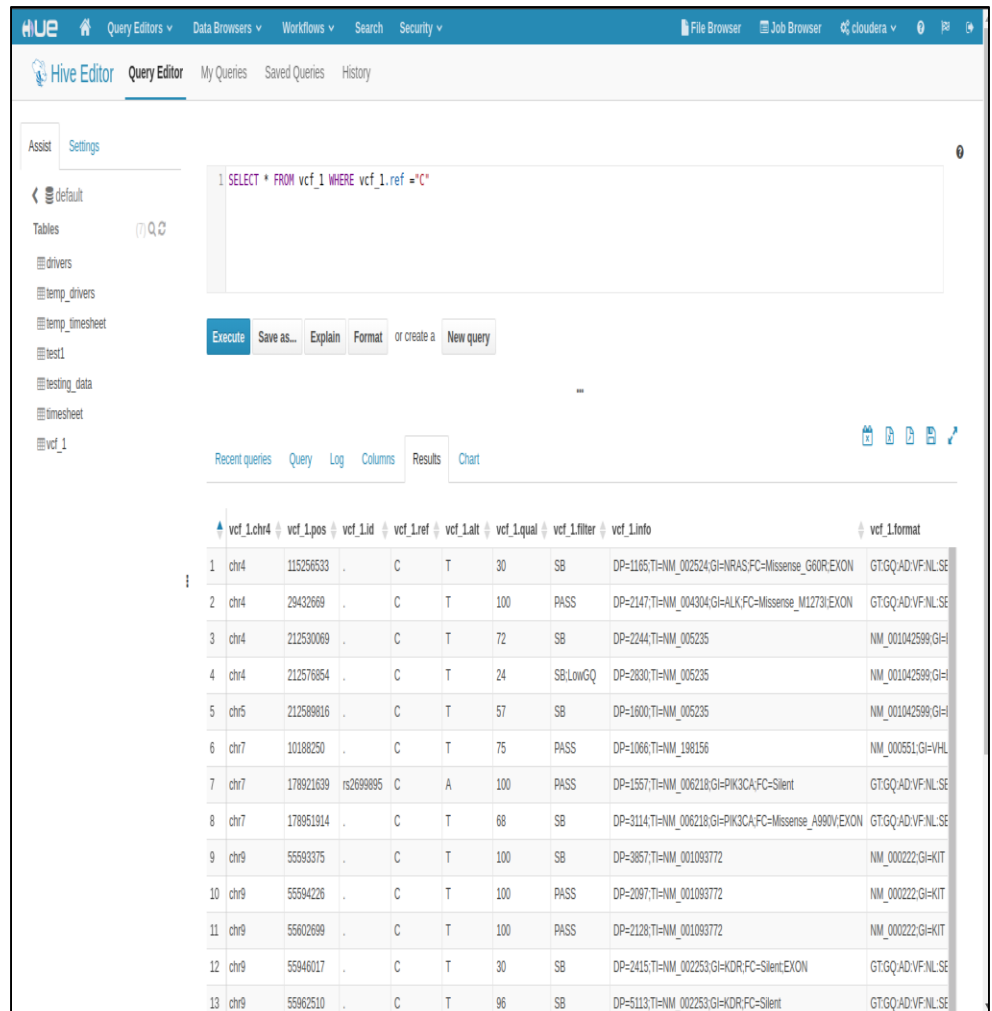


SNP Detection using HadoopBAM



Genome Wide Association Study using Apache HIVE

- Processing of VCF Files in Data Browser and query using Apache HIVE
- Counting the Allele Frequency
- Taking Input Data from different population and finding Genome wide association using Log odds ratio/Likelihood ratio/Chi-square test across different population



The screenshot displays the Apache HIVE Query Editor interface. The top navigation bar includes links for Query Editors, Data Browsers, Workflows, Search, and Security. The main area shows a SQL query: `SELECT * FROM vcf_1 WHERE vcf_1.ref = 'C'`. Below the query editor are buttons for Execute, Save as..., Explain, Format, and a link to create a New query. The results tab is active, showing a table with 13 rows of VCF data. The table columns are: vcf_1_chr4, vcf_1_pos, vcf_1_id, vcf_1_ref, vcf_1_alt, vcf_1_qual, vcf_1_filter, vcf_1_info, and vcf_1_format.

	vcf_1_chr4	vcf_1_pos	vcf_1_id	vcf_1_ref	vcf_1_alt	vcf_1_qual	vcf_1_filter	vcf_1_info	vcf_1_format
1	chr4	115256533	.	C	T	30	SB	DP=1165;TI=NM_002524;GI=NRAS;FC=Missense_G60R;EXON	GT:GQ:AD:VF:NL:SE
2	chr4	29432669	.	C	T	100	PASS	DP=2147;TI=NM_004304;GI=ALK;FC=Missense_M1273I;EXON	GT:GQ:AD:VF:NL:SE
3	chr4	212530069	.	C	T	72	SB	DP=2244;TI=NM_005235	NM_001042599;GI=I
4	chr4	212576854	.	C	T	24	SB,LowGQ	DP=2830;TI=NM_005235	NM_001042599;GI=I
5	chr5	212589816	.	C	T	57	SB	DP=1600;TI=NM_005235	NM_001042599;GI=I
6	chr7	10188250	.	C	T	75	PASS	DP=1066;TI=NM_198156	NM_000551;GI=VHL
7	chr7	178921639	rs2698895	C	A	100	PASS	DP=1557;TI=NM_006218;GI=PIK3CA;FC=Silent	GT:GQ:AD:VF:NL:SE
8	chr7	178951914	.	C	T	68	SB	DP=3114;TI=NM_006218;GI=PIK3CA;FC=Missense_A990V;EXON	GT:GQ:AD:VF:NL:SE
9	chr9	55593375	.	C	T	100	SB	DP=3857;TI=NM_001089772	NM_000222;GI=KIT
10	chr9	55594226	.	C	T	100	PASS	DP=2097;TI=NM_001089772	NM_000222;GI=KIT
11	chr9	55602699	.	C	T	100	PASS	DP=2128;TI=NM_001089772	NM_000222;GI=KIT
12	chr9	55946017	.	C	T	30	SB	DP=2415;TI=NM_002253;GI=KDR;FC=Silent;EXON	GT:GQ:AD:VF:NL:SE
13	chr9	55962510	.	C	T	96	SB	DP=5113;TI=NM_002253;GI=KDR;FC=Silent	GT:GQ:AD:VF:NL:SE

References

- **Searching for SNPs with cloud computing**
Ben Langmead, Michael C Schatz, Jimmy Lin, Mihai Pop and Steven L Salzberg
- **The application of Hadoop in Structural Bioinformatics**
Jamie Alnasir, Hugh P. Shanahan
- **Big Data Processing for Genomics**
Altti Ilari Maarala, Keijo Heljanko, Andre Schumacher, Ridvan Dongelci, Luca Pireddu, Matti Niemenmaa, Aleksi Kallio, Eija Korpelainen and Gianluigi Zanetti

Thank You