

Characterization of Oncogenes in Pretreatment and Post- treatment Sequence Data

Samuel Chatmon
National Cancer Institute

Background

- Obtained Two SRA fastq files from the SRA database.
- Files were labeled as WES Nextgen Sequence data which were obtained from a melanoma patient before and after treatment
- Wanted to examine the sequences for the variant levels in Ocogenes

Project Overview

- Fastq files were obtained from NCBI's SRA database and transfer using the Fastq-dump tool
- The complete hg38 human genomic sequence was downloaded from the UCSC Genomic Browser Site
- The top twenty genes that contains somatic mutations for melanoma was obtained from the Cosmic Cancer Browser site and their chromosome number and position was obtained from the UCSC site
- Both sets of fastq sequences were aligned against the human genome with the BWA aligner

Project Overview (continued)

- Samtools and BCFtools were employed to create VCF files that the contained base substitutions, deletions and insertions in the data.
- Python programming and bash scripting were utilized to score the number of times variants were found in the genes

First Issue: Forgot about HIPAA Compliance

```
manager@bl8vbox[Desktop] fastq-dump SRR6431487 [ 5:08PM]
2018-05-09T17:09:06 fastq-dump.2.9.0 err: query unauthorized while resolving que
ry within virtual file system module - failed to resolve accession 'SRR6431487'
- Access denied - please request permission to access phs001469/HMB-PUB-NPU-MDS
in dbGaP ( 403 )
2018-05-09T17:09:06 fastq-dump.2.9.0 err: item not found while constructing with
in virtual database module - the path 'SRR6431487' cannot be opened as database
or table
manager@bl8vbox[Desktop] █ [ 5:09PM]
```

Fortunately, A Few Files were downloadable

Fortunately, a Few Files were downloadable.

ACCESSION: SRX1515323

☐ [WES of homo sapiens: PBMC: Sample Pt1-normal](#)

325. 1 ILLUMINA (Illumina HiSeq 2500) run: 62.4M spots, 12.5G bases, 4.6Gb downloads

Accession: SRX1515322

☐ [WES of homo sapiens: melanoma: Sample Pt1-baseline](#)

326. 1 ILLUMINA (Illumina HiSeq 2500) run: 90.6M spots, 18.1G bases, 7.5Gb downloads

Accession: SRX1515321

————— ————— —————

Study: WES of patient derived pre-treatment metastatic melanoma on anti-PD-1 antibody treatment. We will call treatment srr8 and pre-treatment srr7 for short.

Second Issue Issue: Did not consider memory or CPU requirements

Bowtie Alignment ~12 hours in Oracles VirtualMachine ~12gb memory and 2 cpu's

	Family	Type	vCPUs	Memory (GiB)	Instance Storage (GB)	EBS-Optimized Available	Network Performance	IPv6 Support
--	--------	------	-------	--------------	-----------------------	-------------------------	---------------------	--------------

Step 2: Choose an Instance Type

<input type="checkbox"/>	GPU graphics	g3.8xlarge	32	244	EBS only	Yes	10 Gigabit	Yes
<input type="checkbox"/>	GPU instances	g2.8xlarge	32	60	2 x 120 (SSD)	-	10 Gigabit	-
<input type="checkbox"/>	GPU compute	p2.8xlarge	32	488	EBS only	Yes	10 Gigabit	Yes
<input type="checkbox"/>	GPU compute	p3.8xlarge	32	244	EBS only	Yes	10 Gigabit	Yes
<input type="checkbox"/>	Memory optimized	r4.8xlarge	32	244	EBS only	Yes	10 Gigabit	Yes
<input type="checkbox"/>	Memory optimized	x1e.8xlarge	32	976	1 x 960 (SSD)	Yes	Up to 10 Gigabit	Yes
<input type="checkbox"/>	Storage optimized	i2.8xlarge	32	244	8 x 800 (SSD)	-	10 Gigabit	Yes
<input type="checkbox"/>	Storage optimized	h1.8xlarge	32	128	4 x 2000	Yes	10 Gigabit	Yes

BWI Alignment ~3 hours in Amazon AWS. But Alas, ,most data on machine. Decided to continue virtual machine and use BWA instead.

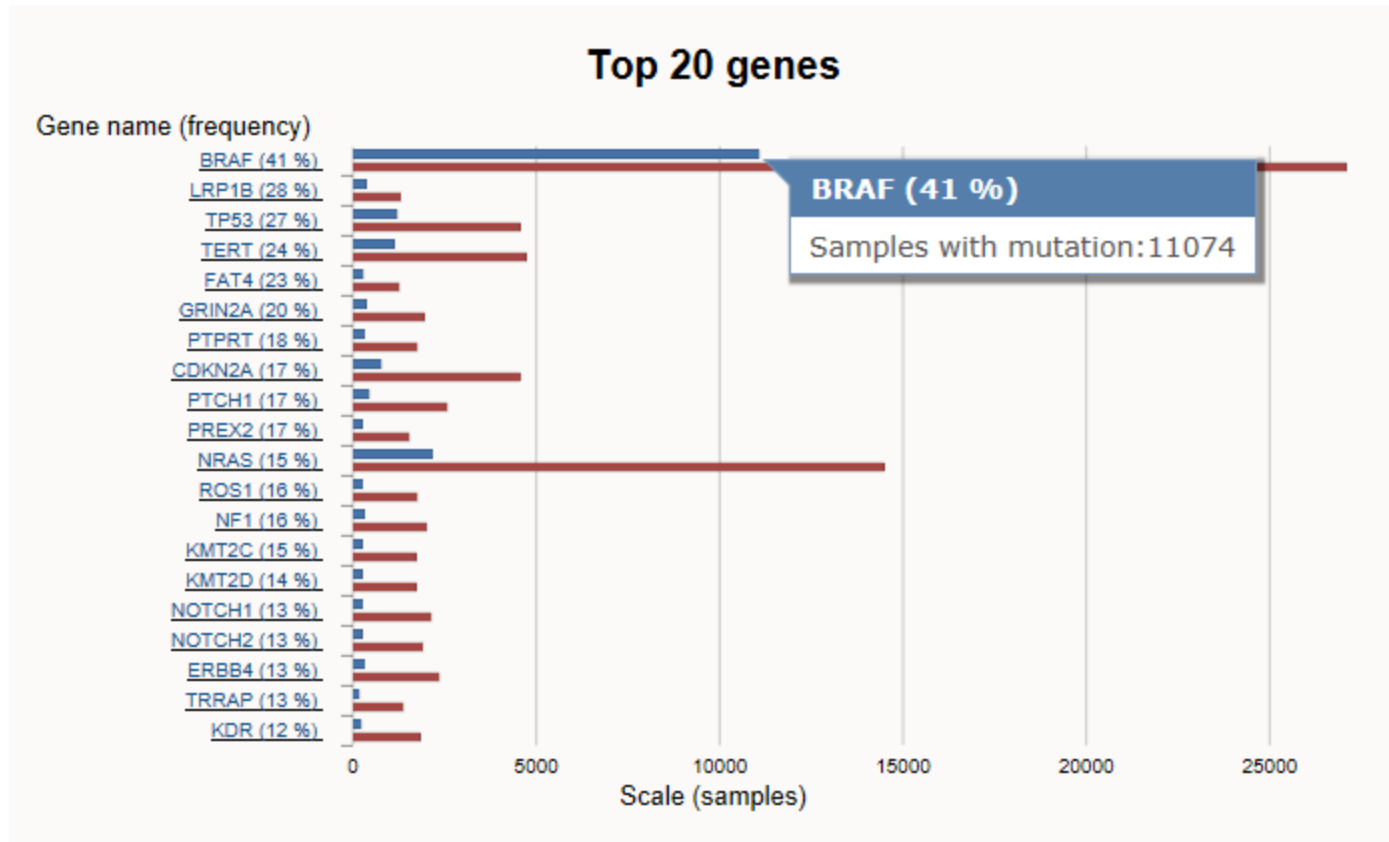
Finally can use Samtools to create bam file of alignments.

[illegible]

and BCFtools to create vcf file which contains all variants.

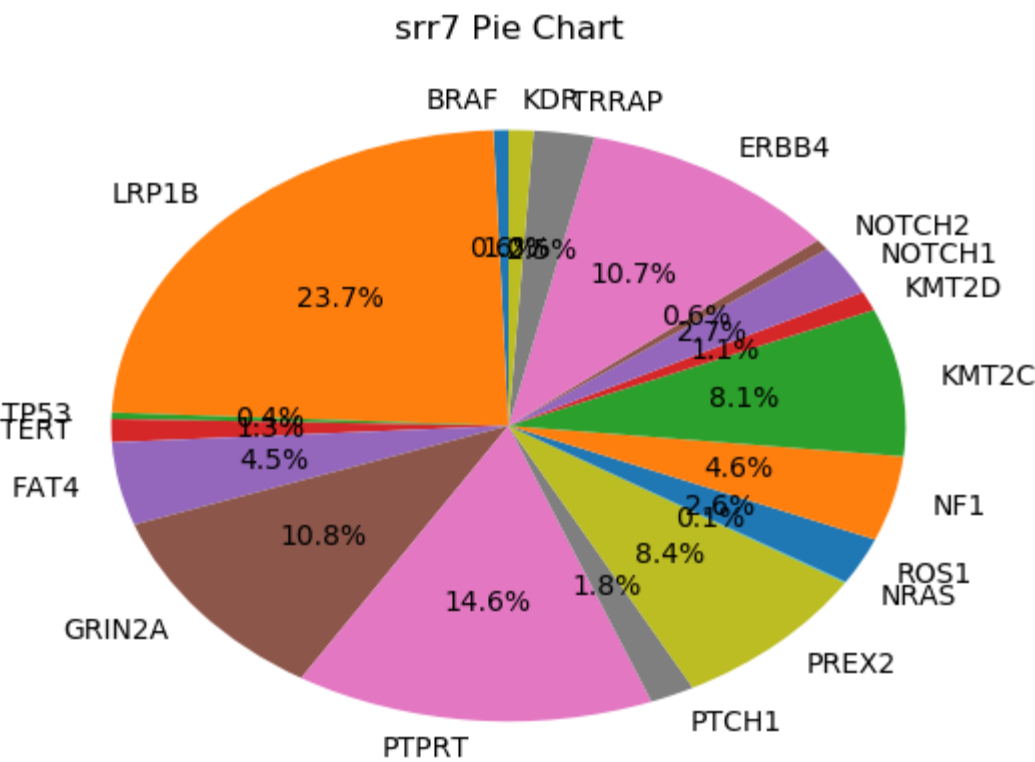
chr1	chr2	chr3	chr4	chr5	chr6	chr7	chr8	chr9	chr10	chr11	chr12	chr13	chr14	chr15	chr16	chr17	chr18	chr19	chr20	chr21	chr22	chr23	chr24	chr25	chr26	chr27	chr28	chr29	chr30	chr31	chr32	chr33	chr34	chr35	chr36	chr37	chr38	chr39	chr40	chr41	chr42	chr43	chr44	chr45	chr46	chr47	chr48	chr49	chr50	chr51	chr52	chr53	chr54	chr55	chr56	chr57	chr58	chr59	chr60	chr61	chr62	chr63	chr64	chr65	chr66	chr67	chr68	chr69	chr70	chr71	chr72	chr73	chr74	chr75	chr76	chr77	chr78	chr79	chr80	chr81	chr82	chr83	chr84	chr85	chr86	chr87	chr88	chr89	chr90	chr91	chr92	chr93	chr94	chr95	chr96	chr97	chr98	chr99	chr100	chr101	chr102	chr103	chr104	chr105	chr106	chr107	chr108	chr109	chr110	chr111	chr112	chr113	chr114	chr115	chr116	chr117	chr118	chr119	chr120	chr121	chr122	chr123	chr124	chr125	chr126	chr127	chr128	chr129	chr130	chr131	chr132	chr133	chr134	chr135	chr136	chr137	chr138	chr139	chr140	chr141	chr142	chr143	chr144	chr145	chr146	chr147	chr148	chr149	chr150	chr151	chr152	chr153	chr154	chr155	chr156	chr157	chr158	chr159	chr160	chr161	chr162	chr163	chr164	chr165	chr166	chr167	chr168	chr169	chr170	chr171	chr172	chr173	chr174	chr175	chr176	chr177	chr178	chr179	chr180	chr181	chr182	chr183	chr184	chr185	chr186	chr187	chr188	chr189	chr190	chr191	chr192	chr193	chr194	chr195	chr196	chr197	chr198	chr199	chr200	chr201	chr202	chr203	chr204	chr205	chr206	chr207	chr208	chr209	chr210	chr211	chr212	chr213	chr214	chr215	chr216	chr217	chr218	chr219	chr220	chr221	chr222	chr223	chr224	chr225	chr226	chr227	chr228	chr229	chr230	chr231	chr232	chr233	chr234	chr235	chr236	chr237	chr238	chr239	chr240	chr241	chr242	chr243	chr244	chr245	chr246	chr247	chr248	chr249	chr250	chr251	chr252	chr253	chr254	chr255	chr256	chr257	chr258	chr259	chr260	chr261	chr262	chr263	chr264	chr265	chr266	chr267	chr268	chr269	chr270	chr271	chr272	chr273	chr274	chr275	chr276	chr277	chr278	chr279	chr280	chr281	chr282	chr283	chr284	chr285	chr286	chr287	chr288	chr289	chr290	chr291	chr292	chr293	chr294	chr295	chr296	chr297	chr298	chr299	chr300	chr301	chr302	chr303	chr304	chr305	chr306	chr307	chr308	chr309	chr310	chr311	chr312	chr313	chr314	chr315	chr316	chr317	chr318	chr319	chr320	chr321	chr322	chr323	chr324	chr325	chr326	chr327	chr328	chr329	chr330	chr331	chr332	chr333	chr334	chr335	chr336	chr337	chr338	chr339	chr340	chr341	chr342	chr343	chr344	chr345	chr346	chr347	chr348	chr349	chr350	chr351	chr352	chr353	chr354	chr355	chr356	chr357	chr358	chr359	chr360	chr361	chr362	chr363	chr364	chr365	chr366	chr367	chr368	chr369	chr370	chr371	chr372	chr373	chr374	chr375	chr376	chr377	chr378	chr379	chr380	chr381	chr382	chr383	chr384	chr385	chr386	chr387	chr388	chr389	chr390	chr391	chr392	chr393	chr394	chr395	chr396	chr397	chr398	chr399	chr400	chr401	chr402	chr403	chr404	chr405	chr406	chr407	chr408	chr409	chr410	chr411	chr412	chr413	chr414	chr415	chr416	chr417	chr418	chr419	
------	------	------	------	------	------	------	------	------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--

Obtained 20 more frequent Mutated Genes from Cosmic Site



Utilized Python and BASH scripting to score variants in each data set : List of Genes and Pie Chart for srr7

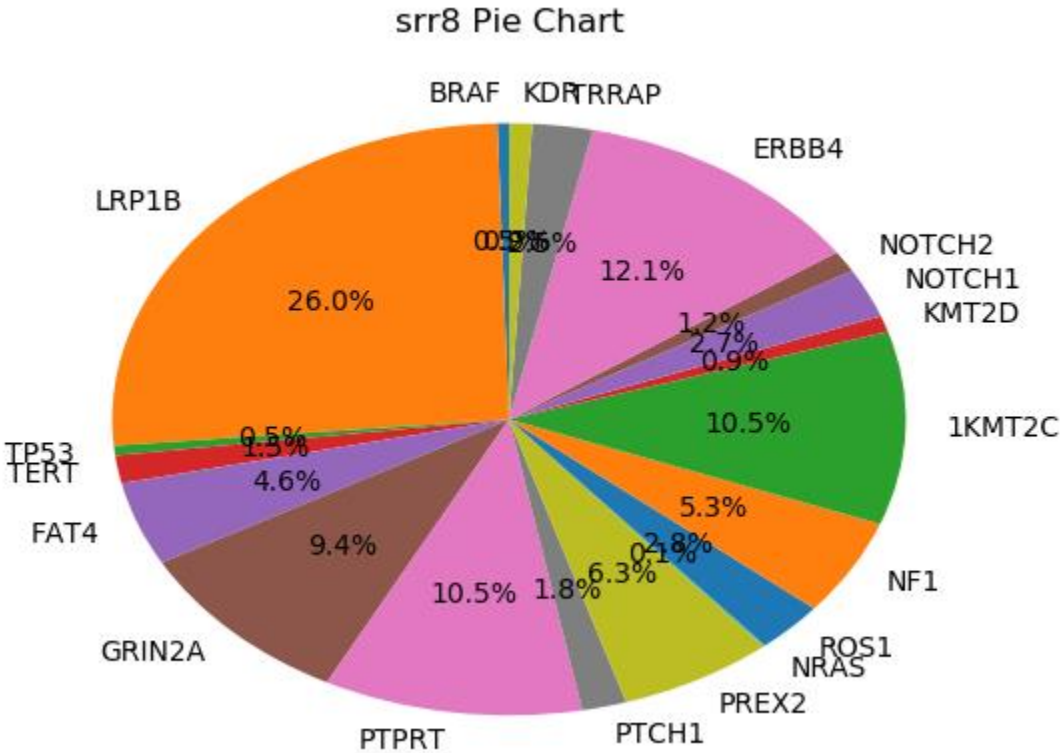
srr7	
Number	Gene Name
1	BRAF
2	ERBB4
3	FAT4
4	GRIN2A
5	KDR
6	KMT2C
7	KMT2D
8	LRP1B
9	NF1
10	NOTCH1
11	NOTCH2
12	NRAS
13	PREX2
14	PTCH1
15	PTPRT
16	ROS1
17	TERT
18	TP53
19	TRRAP



Total number of variants = 1985

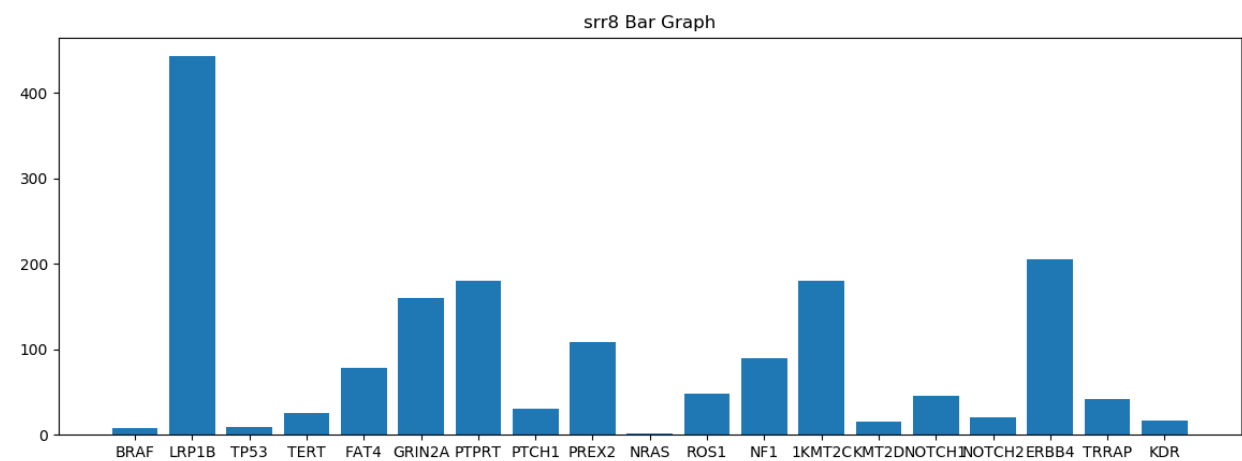
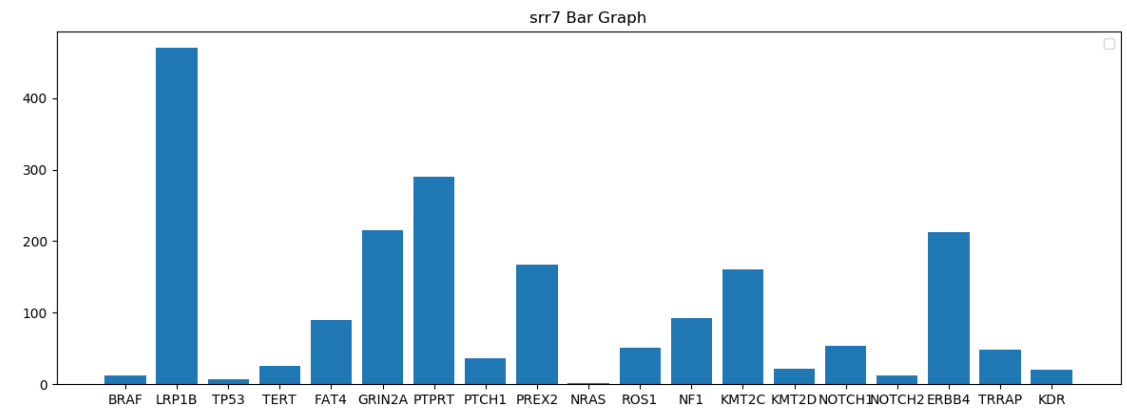
Utilized Python and BASH scripting to score variants in each data set : List of Genes and Pie Chart for srr8

srr8	
Number	Gene Name
1	BRAF
2	ERBB4
3	FAT4
4	GRIN2A
5	KDR
6	KMT2C
7	KMT2D
8	LRP1B
9	NF1
10	NOTCH1
11	NOTCH2
12	NRAS
13	PREX2
14	PTCH1
15	PTPRT
16	ROS1
17	TERT
18	TP53
19	TRRAP



Total number of variants = 1707

Utilized Python and BASH scripting to score variants in each data set : Bar Graphs for srr7 and srr8



Conclusion

- Can't really hypothesize anything by just using simple visualization tools.
- Most likely would need to look at much more data and using Transcriptomics tools and star alignment tool

