

Method

1 Experimental Procedure

1.1 Evaluation of DNA quality

The quality of isolated genomic DNA was verified by using these three methods in combination:

- 1) DNA degradation and contamination were monitored on 1% agarose gels;
- 2) DNA concentration was measured by Qubit® DNA Assay Kit in Qubit® 3.0 Fluorometer (Invitrogen, USA).

1.2 Library Preparation

The exome sequences were efficiently enriched from 0.4 µg genomic DNA using Agilent liquid capture system (Agilent SureSelect Human All Exon V6) according to the manufacturer's protocol. Firstly, qualified genomic DNA was randomly fragmented to an average size of 180-280bp by Covaris S220 sonicator. Remaining overhangs were converted into blunt ends via exonuclease polymerase activities. Secondly, DNA fragments were end repaired and phosphorylated, followed by A-tailing and ligation at the 3' ends with paired-end adaptors. DNA fragments with ligated adapter molecules on both ends were selectively enriched in a PCR reaction. After PCR reaction, libraries hybridize with liquid phase with biotin labeled probe, then use magnetic beads with streptomycin to capture the exons of genes. Captured libraries were enriched in a PCR reaction to add index tags to prepare for sequencing. Products were purified using AMPure XP system (Beckman Coulter, Beverly, USA), DNA concentration was measured by Qubit®3.0 Fluorometer (Invitrogen, USA), libraries were analyzed for size distribution by NGS3K/Caliper and quantified by real-time PCR (3 nM). At last, DNA library were sequenced on Illumina for pairedend 150bp reads

1.3 Clustering & Sequencing

2 Bioinformatics Analysis Pipeline

2.1. Data Quality Control

2.1.1. Raw data

The original fluorescence image files obtained from Illumina platform are transformed to short reads (Raw data) by base calling and these short reads are recorded in FASTQ format, which contains sequence information and corresponding sequencing quality information.

2.1.2. Evaluation of data (Data quality control)

Sequence artifacts, including reads containing adapter contamination, low-quality nucleotides and unrecognizable nucleotide (N), undoubtedly set the barrier for the subsequent reliable bioinformatics analysis. Hence quality control is an essential step and applied to guarantee the meaningful downstream analysis.

All the downstream bioinformatics analyses were based on the high quality clean data, which were retained after these steps. At the same time, QC statistics including total reads number, raw data, raw depth, sequencing error rate and percentage of reads with Q30 (the percent of bases with phred-scaled quality scores greater than 30) were calculated and summarized.

2.2. Reads Mapping to Reference Sequence

Valid sequencing data is mapped to the reference genome (GRCh37/hg19/ GRCh38) by BurrowsWheeler Aligner (BWA) software to get the original mapping result in BAM format. Subsequently, Samtools and Sambamba are respectively utilized to sort bam files, do duplicate-marking to generate final bam file. If one or one pair read(s) has multiple mapping positions, the strategy adopted by BWA are to select the best one, if there are multi best mapping position, we randomly pick one. Mapping step is very difficult due to mismatches, including true mutation and sequencing error, and duplicates resulted from PCR amplification. These duplicate reads are uninformative and shouldn't be considered as evidence for variants. Sambamba is employed to mark these duplicates so that we will ignore them in the following analysis.

2.3. Variant detection

SAMtools mpileup and bcftools were used to do variant calling and identify SNP, InDels.

2.4. Somatic mutation calling

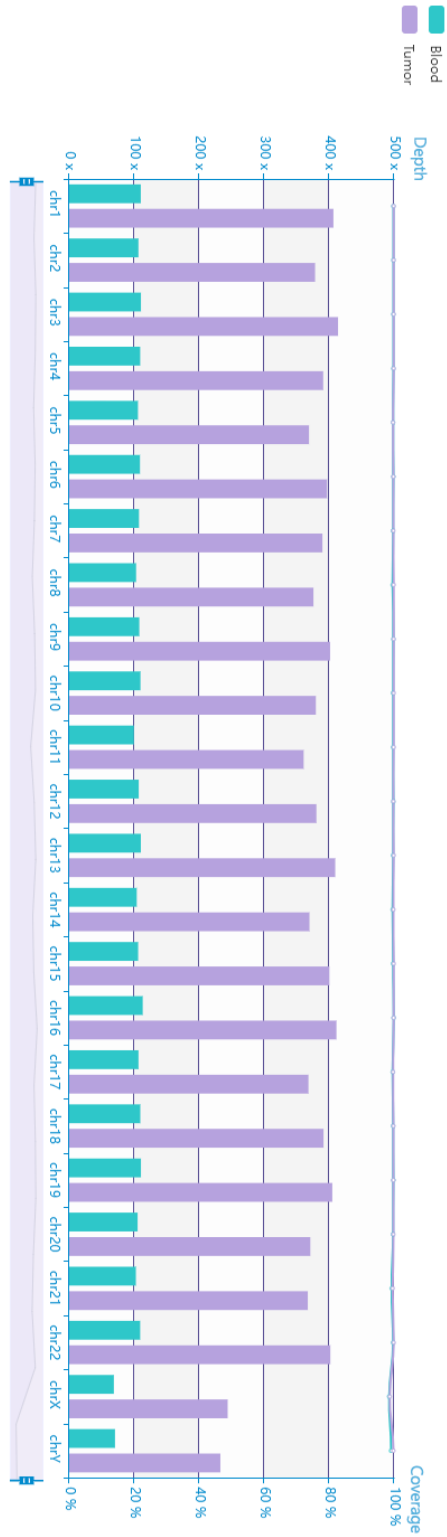
The somatic SNV was detected by muTect, the somatic InDel by Strelka, and Control-FREEC was used to detect somatic CNV

2.5. Annotation

ANNOVAR is performed to do annotation for VCF file obtained in the previous step. The variant position, variant type, conservative prediction and other information are obtained at this step through a variety of databases, such as dbSNP, 1000 Genome, esp6500, GnomAD, CADD, HGMD and COSMIC, and so on. Since we are interested in exonic variants, gene transcript annotation databases, such as Consensus CDS, RefSeq, Ensemble and UCSC, are also applied for annotation to determine amino acid alternation. Functional annotation is very important because the link between genetic variation and cancer can be found. In this step, the databases mainly include GO, KEGG, Reactome, Biocarta, PID, and so on.

Result

1 Sequencing depth and coverage



2 The outcome of somatic mutation

2.1 Single nucleotide variant(SNV)

Sample	Xuezheng_Tumor
CDS	3
synonymous_SNP	1
missense_SNP	2
stopgain	0
stoploss	0
unknown	0
intronic	53
UTR3	6
UTR5	1
splicing	2
ncRNA_exonic	7
ncRNA_intronic	10
ncRNA_UTR3	0
ncRNA_UTR5	0
ncRNA_splicing	1
upstream	3
downstream	2
intergenic	23
Others	0
Total	111

2.2 Insertion and Deletion(InDel)

Sample	Xuezheng_Tumor
CDS	0
frameshift_deletion	0
frameshift_insertion	0
nonframeshift_deletion	0
nonframeshift_insertion	0
stopgain	0
stoploss	0
unknown	0
intronic	0
UTR3	0
UTR5	0
splicing	0
ncRNA_exonic	0
ncRNA_intronic	0
ncRNA_UTR3	0
ncRNA_UTR5	0

ncRNA_splicing	0
upstream	0
downstream	0
intergenic	0
Others	0
<hr/>	
Total	0

SNV list

Priority	L	L	L	L	L	L	L	L	L	L	L	L	L
CHROM	chr6	chr1	chr7	chr9	chr15	chr12	chr22	chr2	chr7	chr15	chr17	chr17	chr11
POS	5575 9124	1615 1393 3	6576 3276	4164 8075	2040 8438	15934	2465 7860	12998 0059	4996 980	10197 5947	41105 795	41105 831	1097926
ID	.	rs409 763	rs368 0944 33	.	rs112 3557 17	rs797 03804 9	rs200 1116 41	rs675 0772	rs963 9891	rs1998 37286	rs756 54580 1	rs7539 2608	.
REF	T	G	T	T	A	T	A	G	C	C	A	A	C
ALT	C	A	C	C	G	G	G	C	G	G	G	C	G
QUAL
FILTER	PASS	PASS	PASS	PASS	PASS	PASS	PASS	PASS	PASS	PASS	PASS	PASS	PASS
GeneName	BMP5	FCGR2A	CCT6P1	FAM27E2	HERC2P3	LOC100288778	POM121L10P	RAB6C-AS1	RNF216P1	WASH3P	KRTA4-9	KRTAP4-9	MUC2
Description		Fc fragment of Bone IgG, morphologic protein-5 receptor for (CD32)											Mucin 2, intestinal /tracheal
Func	splicing	splicing	ncRNA_A_splicing	ncRNA_A_exonic	ncRNA_A_exonic	ncRNA_A_exonic	ncRNA_A_exonic	ncRNA_A_exonic	ncRNA_A_exonic	ncRNA_exonic	exonic	exonic	exonic
ExonicFunc	missense SNV	missense SNV	synonymous SNV
Gene	NM_021073	NM_036219, NM_021642	NR_003110	NR_003071	NR_003643	NR_0130745	NR_0024593	NR_006537	NR_009023384, NR_0023385	NR_003659	NM_001146041	NM_001146041	NM_002457

	rs409763	rs368094433	rs112355717	rs797038049	rs200111641	rs6750772	rs9639891	rs199837286	rs756545801	rs75392608
avsnp150
clinvar_20190305
gwasCatalognew
1000g2015aug_Chinese	0.24	.	.	.
1000g2015aug_eas	0.25	.	.	.
1000g2015aug_all	0.14	.	.	.
esp6500siv2_all
GnomAD_exome_AF_popmax	0	#####	0.001	.
GnomAD_exome_AF_eas	0	3E-04	3E-04	0.001
GnomAD_genome_AF_popmax	0	0	0.31	.	0	0.119	0.11	0.028	0.001	0.005
GnomAD_genome_AF_eas	0.01	0	0.3	.	0	0.136	0.26	0.007	0.005	0.005
NovoDb_WES	0.02
NovoDb_WGS	.	.	0.2	.	.	0.109	0.25	.	0.004	0.004
dbSNP_SCORE	0.0003	1.0000
Interpro_domain
SIFT	1.0,T	0.899,T
Polyphen2_HDIV	0.0,B	0.0,B
Polyphen2_HVAR	0.0,B	0.0,B
LRT	0.421,N	0.002,N
MutationTaster	0.999	1,N	1,N
MutationAssessor	1.29,N	3.205,N
FATHMM	2.2,T	6.59,T

CADD_raw	1.75										2.089	2.107		
	2,14										,0.00	,0.00		
	.71										1	1		
SiPhy_29way_logOdds	7.87										7.549	2.384		
gerp++gt2	2.16	2.1			2.51									
REVEL											0.194	0.013	0.029	
cosmic70														
ExAC_AL	0										2E-04	0.001		
ExAC_EAS	0.01										0.001	4E-04		
INFO	DB;S OMA TIC; VT= SNP	DB;S OMA TIC; VT= SNP	SOM ATIC ;VT= SNP	SOM ATIC ;VT= SNP	DB;S OMA TIC; VT= SNP	SOMA TIC;V T=SN P	DB;S OMA TIC; VT=S NP	DB;S OMAT IC;VT =SNP	DB;S OMA TIC; VT= SNP	DB;SO MATIC ;VT= SNP	SOMA TIC;V T=SN P	SOMA TIC;V T=SN P	SOMATI C;VT=S NP	
FORMAT	GT:A D:BQ :DP: FA:S S	GT: AD: BQ: :DP:F A:S S	GT:A D:BQ :DP: FA:S S	GT:A D:BQ :DP: FA:S S	GT: AD:B Q:DP :FA: SS	GT:A D:BQ :DP:F A:SS	GT:A D:BQ :DP: FA:S S	GT:A D:BQ :DP:F A:SS	GT:A D:BQ :DP: FA:S S	GT:A D:BQ :DP:FA :SS	GT:A D:BQ :DP:F A:SS	GT:A D:BQ :DP:FA :SS	GT:AD: BQ:DP:F A:SS	
Xuezheng_Blood	0:27 ,0: :27: 0.00 :0	0:37 ,1: :38: 0.02 6:0	0:58 ,1: :59: 0.01 7:0	0:8, 0: 8:0. 00:0	0:14 ,0: :14: 0.00 :0	0:12, 0: :14: 0:0.0 0:0	0:15 ,0: :15: 0.00 :0	0:40, 1: :40: 1:0.0 24:0	0:9, 0: 9:0. 00:0	0:10, 0: :10: :0.00 :0	0:35, 1: :35: 7:0.0 28:0	0:40, 1: :40: 1:0.0 24:0	0:25,0: : :16:0. 00:0	
Xuezheng_Tumor	0/1: 187, 22:3 7:21 0:0. 105: 2	0/1: 163, 18:3 7:18 2:0. 099: 2	0/1: 115, 12:3 7:12 7:0. 094: 2	0/1: 0,15 :37: 15:1 :00: 2	0/1: 22,5 :37: 27:0 :185 :2	0/1:1 3,5:3 7:18: 0.278 :2	0/1: 35,8 :37: 43:0 :186 :2	0/1: 164,1 3:37: 177:0 :073: 2	0/1: 29,1 :39: 39: 0.25 6:2	0/1:4 3,9:3 5:52: 0.173 :2	0/1:1 12,13 :36:1 25:0. 104:2	0/1:1 34,22 :36:1 56:0. 141:2	0/1:46, 7:37:53 : :0.132: 2	
Ori_REF	T	G	T	T	A	T	A	G	C	C	A	A	C	
Ori_ALT	C	A	C	C	G	G	G	C	G	G	G	C	G	
shared_hom	0	0	0	0	0	0	0	0	0	0	0	0	0	0
shared_het	0	0	0	0	0	0	0	0	0	0	0	0	0	0

ID=CO
SM393
2020;
OCCU
RENCE
=3(uri
nary_tr
act)

	1991		
	5573		
	1536		
	7919		
	2565		
	886 2		
	9653		
	89 24		
	1421		
	5 168		
	1158 3644		9722984
	0864 1901		1980995
	2263 1614		2703501
	636 1 2260		1985113
	0677 8500		1675487
	507 1 9843		7 170580
PubMedID	4279 982 1		67 22658
	04 13 3011		29 89757
	3931 70 25		11 15081
	6 863 2173		123 2407
	0036 2 213		2822 786
	5319 9735		4825 106
	761 1700		36731 18
	5690		85763 11
	8636		872843
	449 2		
	9716		
	15 19		
	9658		
	03 10		
	6753		
	63 26		
	4362		
GO_BP			DIGESTI ON; SYS TEM_PR ROTEIN ACEOUS _EXTRA CELLULA R_MATRI X; EXTR CELLUL AR_REGI ON; EXT RACELL ULAR_RE GION_PA STRUCT URAL_M OLECULE _ACTIVI TY; EXT RACELL ULAR_M ATRIX_S
GO_CC			
GO_MF			

	KEGG	
	FC	
	GAM	
	MA_	
	R_ME	
KEGG_PA	KEGG	
THWAY	_HED	DIAT
	GEHO	ED_P
	G_SI	HAG
	GNAL	OCY
	ING_	TOSI
	PATH	S;KE
	WAY	GG_L
	;KEG	EISH
	G_TG	MAN
	F_BE	IA_I
	TA_S	NFEC
	IGNA	TION
	LING	;KEG
	_PAT	G_S
	HWA	YSTE
	Y	MIC_
		LUPU
		S_ER
		YTHE
		MAT
		OSU
		PID_
		PTP1
		BPAT
PID_PATH		HWA
WAY		Y;PI
		D_IN
		TEGR
		IN2_
		PATH
		WAY
BIOCART	BIOC	
A_PATH	ART	
WAY	A_AL	
	K_PA	
	THW	
	AY	
		KEGG_VI
		BRIO_CH
		OLERAE_
		INFECTI
		ON

REACTOM
E_PATHW
AY

REACTO
ME_O_LI
NKED_GL
YCOSYL
ATION_
OF_MUC
INS;REA
CTOME_
TERMIN
ATION_
OF_O_G
LYCAN_
BIOSYNT
HESIS;R
EACTOM
E_META
BOLISM_
OF_PRO
TEINS;R