# GaeaAnnotator 说明文档(v0.1)

# 大数据计算组

## 黄志博

# 2017年3月2日

# 目录

1	概还	2
	1.1 开发目的	2
	1.2 开发人员	2
2	程序特性说明	2
	2.1 reference 版本	2
	2.2 注释变异类型	2
	2.3 输入输出格式	2
	2.4 运行环境及开发语言	2
	2.5 性能	3
3	注释数据库	3
	3.1 基因及转录本信息数据库	3
	3.2 其他可用数据库	3
4	使用说明	3
	4.1 配置文件	3
	4.2 运行	4
	4.3 经甲层子	5



### 1 概述

#### 1.1 开发目的

随着高通量测序技术的发展及其越来越广泛的应用,当今世界正在累积巨量的 NGS 数据,而对 NGS 变异数据的注释和解读已经成为快速理解这些测序数据的一个瓶颈。为了实现对大量变异数据的快速注释,我们决定基于 hadoop 平台开发并行化的注释程序。

#### 1.2 开发人员

黄志博张勇李胜康石泉肖鹏

### 2 程序特性说明

#### 2.1 reference 版本

目前暂只支持 hg19/GRCh37 版本,后续会导入 GRCh38 版本数据库。

- hg19/GRCh37
- o GRCh38

#### 2.2 注释变异类型

Туре	What is means	Example						
SNP	Single-Nucleotide Polymorphism	Reference = 'A', Sample = 'C'						
Ins	Insertion	Reference = 'A', Sample = 'AGT'						
Del	Deletion	Reference = 'AC', Sample = 'C'						
MNP	Multiple-nucleotide polymorphism	Reference = 'ATA', Sample = 'GTC'						
MIXED	Multiple-nucleotide and an InDel	Reference = 'ATA', Sample = 'GTCAGT'						

#### 2.3 输入输出格式

• 输入文件: VCF

• 输出文件: TSV, Excel<sup>1</sup>

#### 2.4 运行环境及开发语言

• 运行环境: Hadoop 2.6.0 + Hbase 1.2.0

• 开发语言: Java

 $<sup>^1</sup>$ 使用工具/ifs4/ISDC\_BD/GaeaProject/software/GaeaAnnotator/toExcel.sh 将 TSV 格式结果转换成 Excel, 注意 TSV 文件不能过大



#### 2.5 性能

• 在 BGI hadoop50 集群注释单个 WGS 变异数据(NA2878 样本 4987513 条变异数据)用时约 11 分钟。

### 3 注释数据库

对 NGS 变异信息进行 Gene 注释、功能预测和致病性预测等。

### 3.1 基因及转录本信息数据库

- UCSC refgene
- ENSEMBL Gene
- o UCSC Known Gene
- o CCDS

#### 3.2 其他可用数据库

详细信息见AnnotationDatabase.md,

注释条目列表见/ifs4/ISDC\_BD/huangzhibo/Data/database/header/。

- ESP6500
- G1000
- EXAC
- dbSNP
- HGNC
- gwasCatalog
- CLINVAR
- dbNSFP
- HGMD

## 4 使用说明

#### 4.1 配置文件

配置文件中需要指明以下几项内容:

ref: reference 版本



GeneInfo: 指定基因及转录本信息文件

GeneInfoType: 指定基因及转录本信息数据库版本

{dbName}.fields: 根据注释 header List 中<sup>2</sup>的条目设置各数据库的注释字段。不设置或值为空则视为不对该数据库进行注释。

#### 配置文件示例: /ifs4/ISDC\_BD/GaeaProject/software/GaeaAnnotator/config.properties

```
ref
 2
        GeneInfo
                                   = /ifs4/ISDC_BD/GaeaProject/software/GaeaAnnotator/hg19_refGene.txt
       GeneInfoType = refGene
 3
       # dbName.fields = Header Name
 6
                                             = EFFECT, IMPACT, GENE, GENEID, TRID, BIOTYPE, RANK, HGVS_DNA, HGVS_P
       HGMD.fields
                                               = HGMD_disease, HGMD_gene, HGMD_chrom, HGMD_genename, HGMD_gdbid, HGMD_omimid,
                  HGMD_amino, HGMD_deletion, HGMD_insertion, HGMD_codon, HGMD_codonAff, HGMD_descr, HGMD_hgvs,
                  HGMD_hgvsAll, HGMD_dbsnp, HGMD_chromosome, HGMD_startCoord, HGMD_endCoord, HGMD_tag, HGMD_dmsupport,
                  HGMD_author, HGMD_fullname, HGMD_allname, HGMD_vol, HGMD_page, HGMD_year, HGMD_pmid, HGMD_reftag,
                  {\tt HGMD\_comments}, {\tt HGMD\_acc\_num}, {\tt HGMD\_new\_date}, {\tt HGMD\_base}
        ESP6500.fields
                                              = MAF, HGVS_CDNA_VAR, HGVS_PROTEIN_VAR, GWAS_PUBMED, GTC, GTS, EA_AC, AA, AA_AC, AA_AGE,
                  AA_GTC,CDS_SIZES,GL
10
        G1000.fields = AF, EUR_AF, AFR_AF, AMR_AF, EAS_AF, SAS_AF, VT
        EXAC.fields
                                              = AC, AF
11
        dbSNP.fields
                                              = RS, DBSNP_CAF, DBSNP_COMMON, dbSNPBuildID
12
                                             = CLNACC, CLNDBN, CLNDSDBID, CLNDSDB, CLNSIG, CLNSRCID, CAF, CLNHGVS, GENEINFO, SAO, VC, PM,
       CLINVAR.fields
13
                  OM, MTP, R3, R5, ALT
       #gwasCatalog.fields
                                                                 = riskAllel,riskAlFre,title,pubMedID,trait
        HGNC.fields
                                            = hgnc_id,symbol,name,locus_group,locus_type,status,location,location_sortable,
                 date_symbol_changed,date_name_changed,date_modified,entrez_id,ensembl_gene_id,vega_id,ucsc_id,
                  ena,refseq_accession,ccds_id,uniprot_ids,pubmed_id,mgd_id,rgd_id,lsdb,cosmic,omim_id,mirbase,
                  \verb|homeodb|, \verb|snornabase|, \verb|bioparadigms_slc|, or phanet|, \verb|pseudogene.org|, \verb|horde_id|, \verb|merops|, imgt|, iuphar|, or paradigms_slc|, or phanet|, \verb|pseudogene.org|, borde_id|, \verb|merops|, imgt|, iuphar|, or paradigms_slc|, or phanet|, or paradigms_slc|, or paradigms_s
                  kznf_gene_catalog,mamit-trnadb,cd,lncrnadb,enzyme_id,intermediate_filament_db
                                           = SIFT_pred,SIFT_score,Polyphen2_HDIV_pred,Polyphen2_HDIV_score,Polyphen2_HVAR_pred,
                  Polyphen2_HVAR_score, LRT_pred, LRT_score, MutationTaster_pred, MutationTaster_score,
                  MutationAssessor_pred, MutationAssessor_score, FATHMM_pred, FATHMM_score
```

#### 4.2 运行

执行方式: hadoop jar GaeaAnnotator.jar [options]

程序路径: /ifs4/ISDC\_BD/GaeaProject/software/GaeaAnnotator/GaeaAnnotator.jar

示例脚本: /ifs4/ISDC\_BD/GaeaProject/software/GaeaAnnotator/run.sh

目前只开放 hadoop50 集群作为注释计算平台,使用方法见示例脚本。程序具体参数如下:

```
-c,--config <arg>
                           config file.
2
       --cacheref
                           DistributedCache reference sequence file list
3
       --debug
                           for debug.
    -h,--help
                           help information.
    -i,--input <arg>
                           input file(VCF).
    -m,--mapperNum <arg>
                         mapper number.
    -o,--output <arg>
                           output file of annotation results
   -r,--reference <arg>
                           indexed reference sequence file list [null]
                           display verbose information.
       --verbose
```

 $<sup>^2</sup>$ 见/ifs4/ISDC \_BD/huangzhibo/Data/database/header/



说明:

- -i 须指定非压缩 VCF 文件。
- -o 本地输出文件(gz 压缩)。
- -r vcf 对应的 reference 序列 Gaea 索引文件,见/ifs4/ISDC\_BD/GaeaProject/reference/
- -c 配置文件
- -m mapper 数目,根据集群情况设定。

#### 4.3 结果展示

#CHROM	POS	REF	ALT	EFFECT	IMPACT	GENE	GENEID	TRID	BIOTYPE	RANK	HGVS_DNA	HGVS_P	gwasCato	log_ris	<allel< th=""><th>gwasCato</th><th>alog_ris</th></allel<>	gwasCato	alog_ris
22	37690808	3			intron_	variant	MODIFIER		CYTH4	CYTH4	NM_013385	protein_	coding		c.167+	43C>T	
22	37692024	1			intron_	variant	MODIFIER		CYTH4	CYTH4	NM_013385	protein_	coding		c.168-	16G>A	
22	37692092				missens	e_varian		MODERATE		CYTH4	CYTH4 NM_0133	85	protein_	coding		c.220A>0	ĵ.
22	37692120	9	G	Α	intron_	variant	MODIFIER	:	CYTH4	CYTH4	NM_013385	protein_	coding	4	c.234+	14G>A	

图 1: TSV 格式注释结果

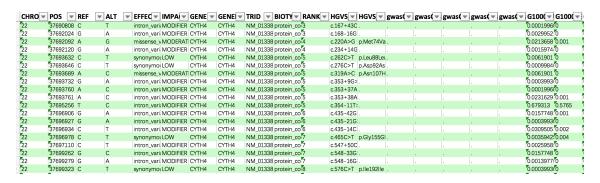


图 2: Excel 格式注释结果

5