

# VCF\_creator : module of DiscoSnp++

## Mapping and VCF Creation features

### Table of contents

User part.....	1
VCF_creator : one module in DiscoSnp++ Mapping and VCF Creation.....	1
Quick starting.....	1
Running VCF_creator.....	1
Options.....	2
Output.....	2
Examples for the filter fields.....	5

### User part

#### **VCF\_creator : one module in DiscoSnp++ Mapping and VCF Creation**

VCF\_creator is designed to map on a genome the output of DiscoSnp++ the Single Nucleotide Polymorphism (SNP) and small indels. This module create a Variant Calling Format (VCF) from the output of DiscoSnp++ or from an alignment obtained with the software BWA\* .

#### **Quick starting**

- Download and uncompress the DiscoSnp++
- Install BWA (you can add it to your PATH )

#### **Running VCF\_creator**

- The main script *run\_VCF\_creator.sh* has three modes according to your needs :
  - **MODE 1** : You don't have a reference genome but you want to create a vcf (it will summarize the DiscoSnp++ informations and will give the position of the variant on the upper path of DiscoSnp++ variant). The module will create a vcf from the output of DiscoSnp++ :
    - `./run_VCF_creator.sh -p <disco_file> -o <output>`
  - **MODE 2** : You have a reference genome and you want to align your variants against a reference genome. The module will run BWA to make an alignment between your reference genome and the output of DiscoSnp++ :
    - `./run_VCF_creator.sh -G <ref> -p <disco_file> -o <output> [-B <path_bwa>] [-w]`
  - **MODE 3** : You already have an alignment (.sam file) and you want to create a vcf file :
    - `./run_VCF_creator.sh -f <sam_file> -o <output>`

#### **Options**

General options :

- -p : DiscoSnp++ output file (<file>.fasta) (**Mandatory** unless MODE 3)
- -o : ouput (<file>.vcf) (**Mandatory**)
- -G : reference genome (<file.fasta>) (Only in MODE 2)
- -B : bwa path ( /home/me/my\_programs/bwa-0.7.12/) (note that bwa must be pre-compiled) (Only in MODE 2) (not necessary if BWA is in the path)

---

\* : (Li H. and Durbin R. (2010) *Fast and accurate long-read alignment with Burrows-Wheeler Transform*. *Bioinformatics*, Epub. [PMID: 20080505]).

- -f : alignment already done (<file>.sam) (Only in MODE 3)
- -h : show help
- -w : remove waste tmp files (bwa index files)
- -I : create of an output (VCF format) specific to IGV (Integrative Genomics Viewer) : sorting VCF by mapping positions and removing unmapped variants

## Output

- **Final results are in your <file>.vcf.** It's a Variant Call Format that will summarize all the mapping information and the header of DiscoSnp++ informations. Example :
- **MODE 1 :**

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	G1	G2	G3	G4	G5
SNP_higher_path_13736	30	13736	A	G	.	.	.	.	.	.	.	.	.
Ty=SNP;Rk=0.99909;UL=.;UR=.;CL=.;CR=.									GT:DP:PL:AD 0/0:53447:1947,160066,1067676:53425,22				
0/1:3:40,9,20:1,2 1/1:94862:1895477,284398,3575:30,94832 0/1:5:34,9,54:3,2									1/1:65389:1306661,196101,2493:19,65370				
SNP_higher_path_5442	30	5442_1	A	G	.	.	.	.	.	.	.	.	.
Ty=SNP;Rk=0.42995;UL=.;UR=.;CL=.;CR=.									GT:DP:PL:AD 0 1:19:98,14,198:12,7 0				
1:18:138,12,138:9,9 0 1:8:66,10,66:4,4 0 1:130:257,124,1813:104,26 0 0:224:14,598,4324:220,4									SNP_higher_path_5442 31 5442_2 G A . . . . .				
Ty=SNP;Rk=0.42995;UL=.;UR=.;CL=.;CR=.									GT:DP:PL:AD 0 1:19:98,14,198:12,7 0				
1:18:138,12,138:9,9 0 1:8:66,10,66:4,4 0 1:130:257,124,1813:104,26 0 0:224:14,598,4324:220,4									INDEL_higher_path_586 30 586 GGAC G . . . . .				
Ty=INS;Rk=0.40534;UL=.;UR=.;CL=.;CR=.									GT:DP:PL:AD 0/1:127:837,17,977:67,60				
1/1:387:6056,518,408:52,335 0/1:70:372,21,651:42,28 0/1:126:1395,53,477:40,86									0/1:102:671,16,791:54,48				

- **MODE 2 and MODE 3 :**

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	G1	G2	G3	G4	G5
Pseudomonas	109860	15103	C	G	.	.	PASS	.	.	.	.	.	.
Ty=SNP;Rk=0.65101;DT=1;UL=.;UR=.;CL=.;CR=.;Genome=C;Sd=-1									GT:DP:PL:AD 1/1:4:84,16,4:0,4				
1/1:1:24,7,4:0,1 1/1:1:24,7,4:0,1 0/1:6:17,15,97:5,1 0/1:10:124,14,44:3,7									Pseudomonas 4416118 1416_1 A G . . . . .				
Ty=SNP;Rk=0.41638;DT=1;UL=.;UR=.;CL=.;CR=.;Genome=A;Sd=1									GT:DP:PL:AD 1 1:35:606,71,28:3,32				
0 1:26:403,41,43:4,22 1 1:34:617,79,18:2,32 0 1:46:316,14,356:24,22 0									1:43:547,36,128:11,32				
Pseudomonas	4416119	1416_2	A	C	.	.	PASS	.	.	.	.	.	.
Ty=SNP;Rk=0.41638;DT=1;UL=.;UR=.;CL=.;CR=.;Genome=A;Sd=1									GT:DP:PL:AD 1 1:35:606,71,28:3,32				
0 1:26:403,41,43:4,22 1 1:34:617,79,18:2,32 0 1:46:316,14,356:24,22 0									1:43:547,36,128:11,32				
Pseudomonas	1312837	14361	C	.	.	.	PASS	.	.	.	.	.	.
Ty=DEL;Rk=0.26896;DT=0;UL=.;UR=.;CL=.;CR=.;Genome=.;Sd=-1									GT:DP:PL:AD				
0/0:792:47,1978,15015:771,21 0/0:745:34,1905,14223:728,17 0/0:697:37,1765,13268:680,17									0/0:783:32,2016,14979:766,17 0/1:849:1303,868,12339:701,148				
Pseudomonas	1753272	6164	T	G	.	.	PASS	.	.	.	.	.	.
Ty=INS;Rk=0.13677;DT=0;UL=.;UR=.;CL=.;CR=.;Genome=.;Sd=1									GT:DP:PL:AD				
0/0:1086:34,2817,20829:1064,22 0/0:1158:34,3011,22226:1135,23									0/0:1003:32,2608,19250:983,20 0/0:1486:124,3555,27823:1437,49				
0/0:1132:585,1967,19224:1033,99													

- In this example, the three types of DiscoSnp++ variants : in red : simple SNP ; in green close SNPs ; in black INDEL.
- The VCF file has the following fields :

- **CHROM** : chromosome id where the prediction is mapped, or allele id of the upper path if the variant is unmapped or if no reference genome is provided.
- **POS** (1-based leftmost position):
  - If a reference genome is provided and if the variant is mapped on a unique position : the mapping position of the variant
  - If a reference genome is provided and if the variant is not uniquely mapped : one of the positions of the variant (1-based leftmost position)
  - Else (no reference genome provided or unmapped variant) : position of the variant on the upper path of the discoSnp++ prediction (including the left extension)
- **ID** : identification of the variant (used by DiscoSnp++). For the close SNPs, the SNP number is added to the ID. Example : *10388\_2*
- **REF** :
  - If one of the two predicted allele maps this position : the corresponding variant
  - Else, or if no reference genome provided : the lexicographically smallest of the two variants
  - In case of close SNPs : the first is defined as previously described. The following SNPs are those located on the same path
- **ALT** : The variant non reported as the “REF” variant
- **QUAL** : “.” (unused)
- **FILTER** :
  - PASS if the variant is mapped at unique position
  - MULTIPLE if the variant is mapped on multiple positions
  - “.” : if the variant is unmapped or if no reference genome is provided
- **INFO** :
  - **Ty** : Type of variant
    - SNP : If the variant is a simple SNP or close SNPs
    - INS : If the variant mapped corresponds to the longest path ; the alt carries the deletion
    - DEL : If the variant mapped corresponds to the shortest path ; the alt carries the deletion
  - **Rk** : Rank of the prediction computed by DiscoSnp++ (if several datasets are used in DiscoSnp++, ranks the predictions according to their read coverage in each condition favoring SNPs that are discriminant between conditions value between 0 and 1)
  - **UL** : Length of the left unitig (“.” if not computed)

- **UR** : Length of the right unitig (“.” if not computed)
- **CL** : Length of the left contig (“.” if not computed)
- **CR** : Length of the right contig (“.” if not computed)
- **Genome** : Applies only for SNPs when a reference genome is provided (“.” for INDELs and when no reference genome provide or if the variant is unmapped). Reference nucleotide (!!nucleotide in the reference genome !! In general it is correspond to the REF field ; could be different for close snps). **Important Remark** : If one of the two predictions matches the reference : equal to the “REF” field, else equal to the nucleotide of the reference genome.
- **Sd** : Applies only when a reference genome is provided (“.” if no reference genome provided or if the variant is unmapped). Strand of the prediction mapping. “1” : Forward ; “-1” : Reverse. **Important Remark** : Fields “REF” , “ALT” and “Genome” are based on the mapped predictions. If Sd is 1 then these fields correspond to the DiscoSnp++ prediction, else if the Sd is “-1”, then they correspond to the reverse complement of the DiscoSnp++ predictions.
- **FORMAT** : Description of the genotype fields (G1, G2, G3 ...)
- **GT** : genotype, encodes as allele values with 0 corresponding to the reference and 1 to the alternative. About genotypes :
  - If the separator is a “/” the genotypes are unphased ( INDEL, Simple SNP)
  - If the separator is a “|” the genotypes are phased (Close SNPs with the same ID )
- **DP** : Cumulated depth across samples (sum)
- **PL** : Phred-scaled Genotype Likelihood (given by DiscoSnp++)
- **AD** : Depth of each allele by sample
- **HQ** : Haplotype Quality (Q in DiscoSnp++ header)
- **Genotype** : G1, G2 , G3 and so on. Each fields contains all the informations of the format fields for each samples.

## Examples

- **FILTER Fields:**
  - PASS if the variant is mapped at unique position
  - MULTIPLE if the variant is mapped on multiple positions
  - “.” : if the variant is unmapped or if no reference genome is provided

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	G1	G2	G3	G4	G5
Pseudomonas	1945315		C	G	8816	.		<b>PASS</b>					

