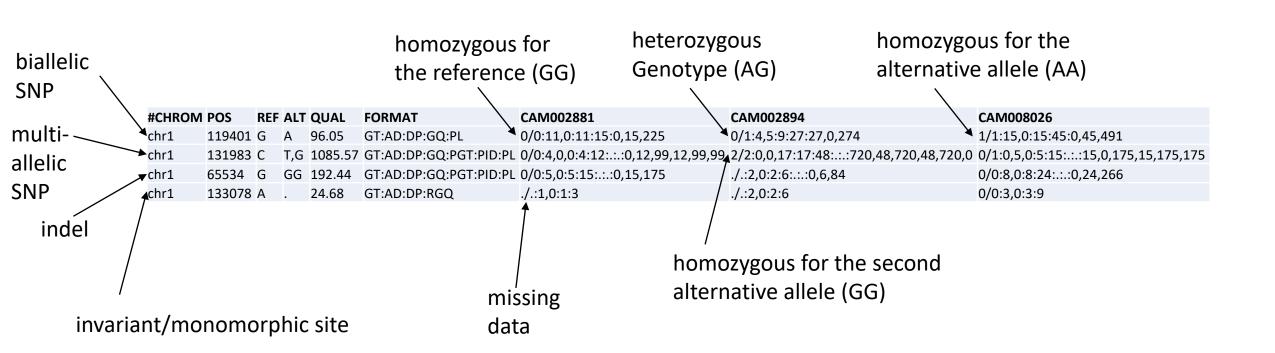
## Vcf file format



## Vcf filtering

| #CHROM | POS    | REF | ALT | QUAL    | FORMAT                 | CAM002881                         | CAM002894                              | CAM008026                            |
|--------|--------|-----|-----|---------|------------------------|-----------------------------------|--|--------------------------------------|
| chr1   | 119401 | G   | Α   | 96.05   | GT:AD:DP:GQ:PL         | 0/0:11,0:11:15:0,15,225           | 0/0:9,0:9:27:0,27,274                  | 0/0:15,0:15:45:0,45,491              |
| chr1   | 131983 | С   | T,G | 1085.57 | GT:AD:DP:GQ:PGT:PID:PL | 0/0:4,0,0:4:12:::0,12,99,12,99,99 | 0/0:17,0,0:17:48:::0,48,720,48,720,720 | 0/0:5,0,0:5:15:::0,15,175,15,175,175 |
| chr1   | 65534  | G   | GG  | 192.44  | GT:AD:DP:GQ:PGT:PID:PL | 0/0:5,0:5:15:::0,15,175           | ./.:2,0:2:6:::0,6,84                   | 0/0:8,0:8:24:::0,24,266              |
| chr1   | 133078 | Α   |     | 24.68   | GT:AD:DP:RGQ           | 0/0:15,0:15:25                    | 0/0:6,0:6:13                           | 0/0:3,0:3:9                          |

## **Genotype** filtering

minGQ 30 (remove genotypes with quality below 30) minDP 10 (remove genotypes with less than 10 reads)

| #CHROM | POS    | REF | ALT | QUAL    | FORMAT                 | CAM002881      | CAM002894                              | CAM008026               |
|--------|--------|-----|-----|---------|------------------------|----------------|--|-------------------------|
| chr1   | 119401 | G   | Α   | 96.05   | GT:AD:DP:GQ:PL         | ./.            | ./.                                    | 0/0:15,0:15:45:0,45,491 |
| chr1   | 131983 | С   | T,G | 1085.57 | GT:AD:DP:GQ:PGT:PID:PL | ./.            | 0/0:17,0,0:17:48:::0,48,720,48,720,720 | ./.                     |
| chr1   | 65534  | G   | GG  | 192.44  | GT:AD:DP:GQ:PGT:PID:PL | ./.            | ./.                                    | ./.                     |
| chr1   | 133078 | Α   |     | 24.68   | GT:AD:DP:RGQ           | 0/0:15,0:15:25 | ./.                                    | ./.                     |

## **Site** filtering

max-missing 0.33 (remove sites with more than 33% missing data) minQ 30 (remove sites with quality below 30)

| #CHROM          | POS               | REF | ALT            | QUAL               | FORMAT                 | CAM002881               | CAM002894                              | CAM008026               |
|-----------------|-------------------|-----|----------------|--------------------|------------------------|-------------------------|--|-------------------------|
| chr1            | 119401            | G   | Α              | 96.05              | GT:AD:DP:GQ:PL         | 0/0:11,0:11:15:0,15,225 | ./.                                    | 0/0:15,0:15:45:0,45,491 |
| <del>chr1</del> | <del>131983</del> | €   | <del>T,G</del> | <del>1085.57</del> | GT:AD:DP:GQ:PGT:PID:PL | <del>./.</del>          | 0/0:17,0,0:17:48:::0,48,720,48,720,720 | <del>./.</del>          |
| <del>chr1</del> | 65534             | G   | GG             | <del>192.44</del>  | GT:AD:DP:GQ:PGT:PID:PL | <del>./.</del>          | <del>./.</del>                         | <del>./.</del>          |
| <del>chr1</del> | <del>133078</del> | A   | <del>.</del>   | <del>24.68</del>   | GT:AD:DP:RGQ           | 0/0:15,0:15:25          | <del>./.</del>                         | <del>./.</del>          |