



# SNPs and Small INDELs Calling From NGS Data

李莉

2013.3

# 工作内容

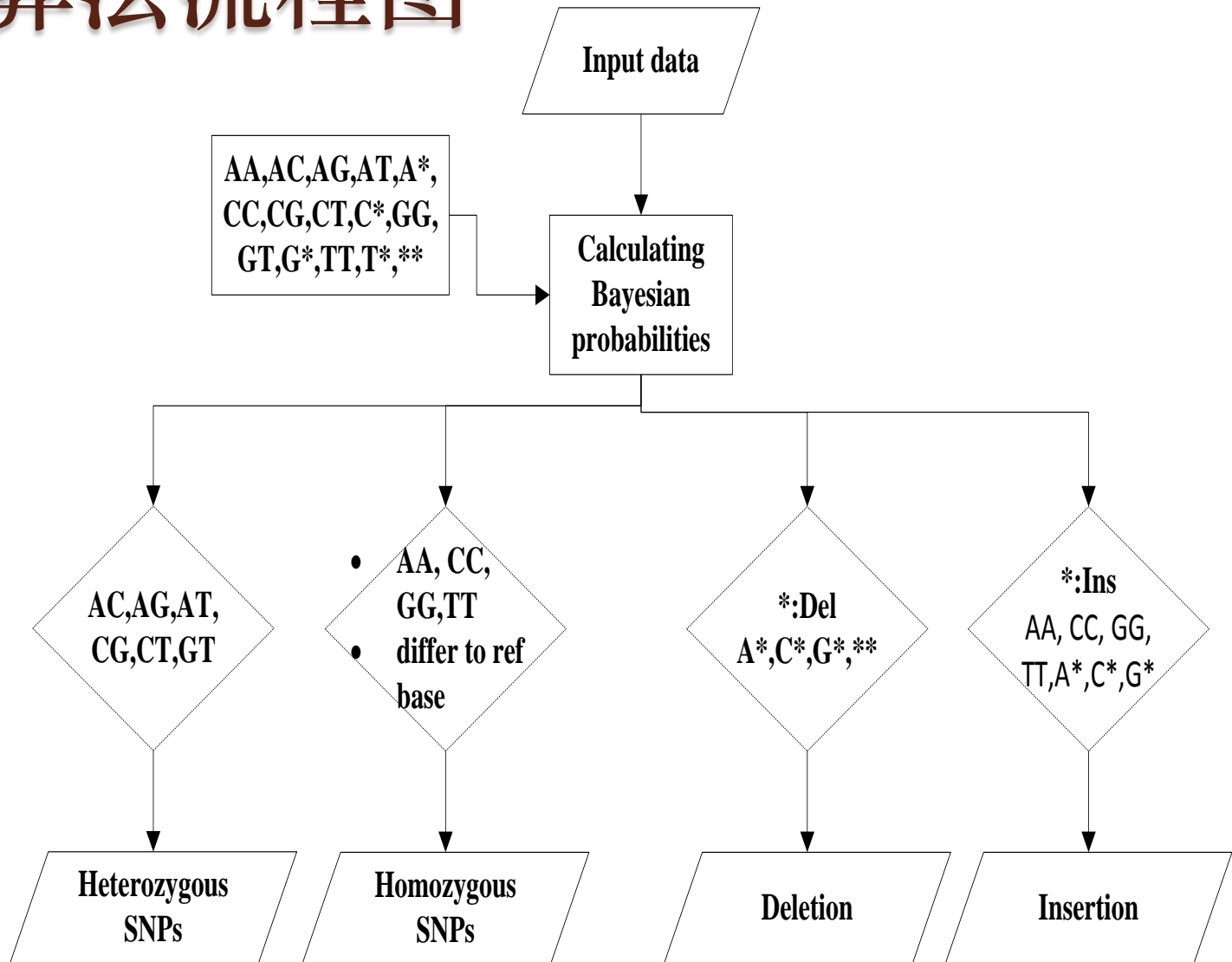
- 基于NGS数据的genomic variation检测(算法及流程)
  - SNP以及small indels(<10bp)的检测
  - large indels检测
  - 其他structural variations检测: inversion、tandem duplication、intra-chromosomal translocation、inter-chromosomal translocation

- SNPs常用检测软件：GATK、SAMtools
- INDELS常用检测软件：
  - Small indels检测：GATK, Dindel, Samtools
  - Large indels检测：Delly(del), Breakdancer
  - Both: Pindel, SOAPindel

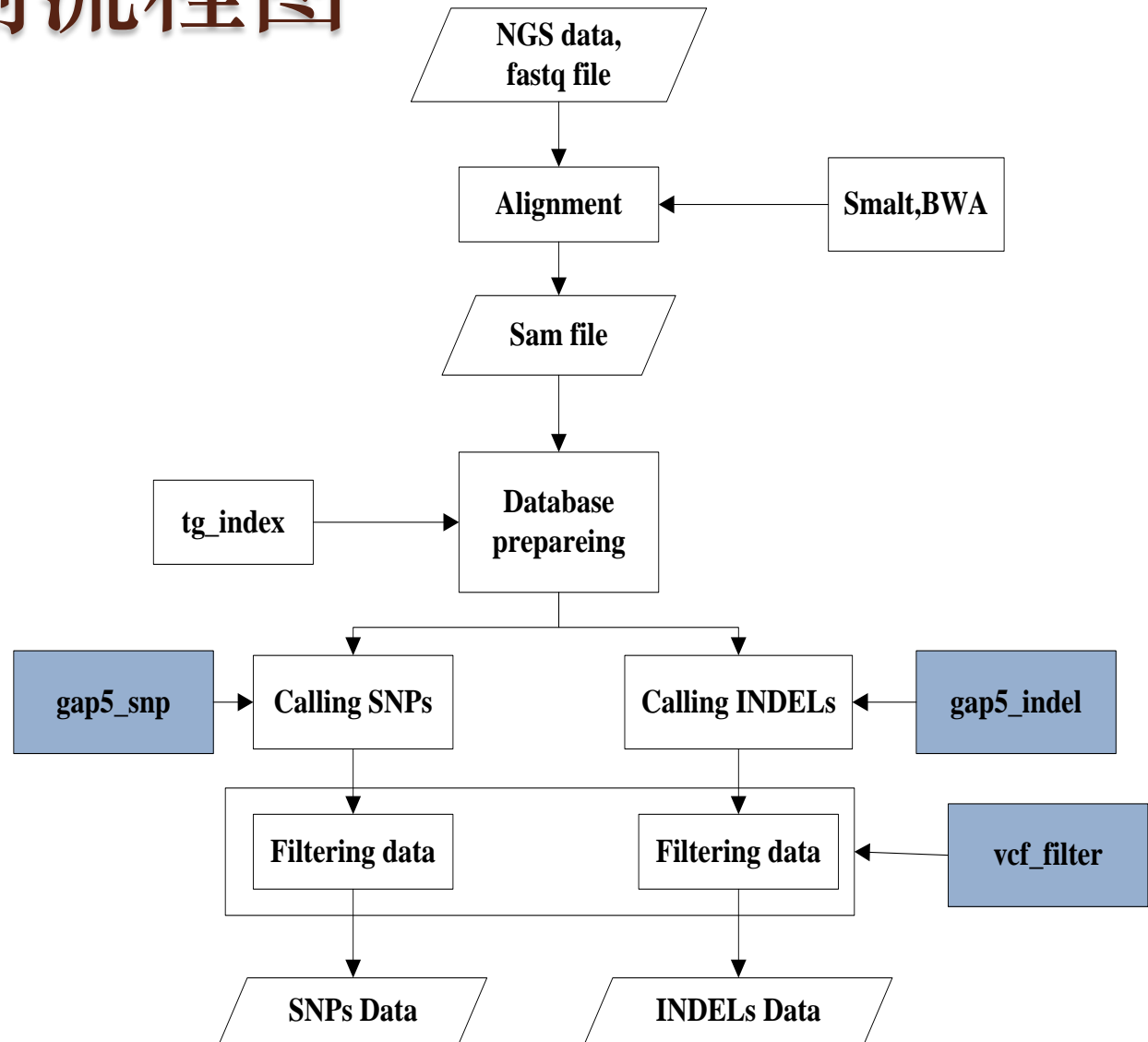
# 基于NGS数据的SNPs和small indels 检测算法流程

- 目的： 尽可能检测出存在的SNPs和small indels
- 算法： 基于Bayesian公式

# 算法流程图



# 检测流程图



# 输出结果例图

```
1 ##fileformat=VCFv4.1
2 ##INFO=<ID=DP,Number=1,Type=Integer,Description="Total number of reads depth">
3 ##INFO=<ID=HP,Number=1,Type=Integer,Description="Reference homopolymer tract length">
4 ##INFO=<ID=BaseQRankSum,Number=1,Type=Float,Description="Z-score from Wilcoxon rank sum test of
5 ##INFO=<ID=MQ,Number=1,Type=Float,Description="RMS Mapping Quality">
6 ##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
7 ##FORMAT=<ID=AD,Number=.,Type=Integer,Description="Allelic depths for the ref and alt alleles :
8 ##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype quality">
9 ##FORMAT=<ID=PL,Number=G,Type=Integer,Description="Normalized, Phred-scaled likelihoods for gen
. n">
0 ##contig=chromosome04
1 ##reference=chr4_ref.fasta
2 #CHROM POS ID REF ALT QUAL FILTER INFO FORMAT VALUE
3 chromosome04 8 . G C 176 PASS DP=40;HRun=0;MQ=54.88 GT:AD:GQ:PL 1/1:0,40:99.00:1780,182,0
4 chromosome04 26 . A G 220 PASS DP=55;HRun=0;MQ=54.96 GT:AD:GQ:PL 1/1:0,55:99.00:2171,225,0
5 chromosome04 34 . T C 204 PASS DP=64;HRun=0;MQ=55.10 GT:AD:GQ:PL 1/1:1,63:99.00:2171,204,0
6 chromosome04 162 . C T 208 PASS DP=80;HRun=0;MQ=54.55 GT:AD:GQ:PL 1/1:0,71:99.00:2171,264,0
7 chromosome04 170 . A G 285 PASS DP=80;HRun=0;MQ=54.36 GT:AD:GQ:PL 1/1:0,80:99.00:3076,290,0
8 chromosome04 176 . G A 275 PASS DP=78;HRun=1;MQ=54.79 GT:AD:GQ:PL 1/1:0,77:99.00:2171,282,0
9 chromosome04 184 . A G 276 PASS DP=78;HRun=0;MQ=55.18 GT:AD:GQ:PL 1/1:0,78:99.00:2171,282,0
0 chromosome04 192 . C T 255 PASS DP=70;HRun=1;MQ=55.66 GT:AD:GQ:PL 1/1:0,70:99.00:2171,260,0
```

# 结果比较(与GATK)

- SNPs calling

Test data: GLA4, chromosome04,  
100bp, coverage of 70x;

Total number of SNPs reference data: 140726

SNPs calling tools: GAP5 and GATK

| <b>Case</b> | <b>Total number</b> | <b>TP</b>     | <b>FP</b>     |
|-------------|---------------------|---------------|---------------|
| <b>GAP5</b> | <b>323588</b>       | <b>118495</b> | <b>205093</b> |
| <b>GATK</b> | <b>351910</b>       | <b>120820</b> | <b>231090</b> |



# 结果比较(与GATK)

- INDELs calling

Test data: GLA4, chromosome04, 100bp, coverage of 70x;

Total number of indels reference data: 19231 ( $\leq 10$ bp)

Indels calling tools: GAP5 and GATK

Insertion: the addition of one or more bases  
into the reference sequence, on the  
contrary is deletion.

| Case | Total number | TP    | FP    |
|------|--------------|-------|-------|
| GAP5 | 39356        | 17219 | 22137 |
| GATK | 36417        | 17927 | 18490 |

# 存在问题

- 算法对杂合的敏感度较高，导致SNPs检测假阳性或者将纯合SNP检测为杂合SNP。
- 不适合平均覆盖度低(低于5)的检测。
- 假阳性偏高

# 过滤方法

- 目的：平衡准确性与检出率，在减少假阳性的同时，保持较高检出率。
  - Raw data (reads) filtering
    - Alignment quality
    - Mapping quality(MQ)
  - Coverage depth(DP)
  - Calling score(QUAL)

# DP对结果的影响

## GAP5

| DP      | total | TP    | FP    |
|---------|-------|-------|-------|
| 0-10    | 1955  | 444   | 1511  |
| 10-20   | 3014  | 955   | 2059  |
| 20-30   | 3514  | 1470  | 2044  |
| 30-40   | 4127  | 1995  | 2132  |
| 40-50   | 4919  | 2722  | 2197  |
| 50-60   | 5708  | 3447  | 2261  |
| 60-70   | 4408  | 2926  | 1482  |
| 70-80   | 1950  | 1314  | 636   |
| 80-90   | 548   | 308   | 240   |
| 90-100  | 152   | 72    | 80    |
| 100-110 | 79    | 26    | 53    |
| 110-120 | 50    | 20    | 30    |
| 120-130 | 52    | 22    | 30    |
| 130-    | 274   | 84    | 190   |
| total   | 30750 | 15805 | 14945 |

## GATK

| DP      | total | TP    | FP    |
|---------|-------|-------|-------|
| 0-10    | 857   | 58    | 799   |
| 10-20   | 2247  | 318   | 1929  |
| 20-30   | 2791  | 643   | 2148  |
| 30-40   | 3561  | 1199  | 2362  |
| 40-50   | 5024  | 2104  | 2920  |
| 50-60   | 6924  | 3656  | 3268  |
| 60-70   | 7240  | 4307  | 2933  |
| 70-80   | 4493  | 2841  | 1652  |
| 80-90   | 1754  | 1067  | 687   |
| 90-100  | 587   | 282   | 305   |
| 100-110 | 294   | 112   | 182   |
| 110-120 | 200   | 67    | 133   |
| 120-130 | 159   | 58    | 101   |
| 130-    | 723   | 226   | 497   |
| total   | 36854 | 16938 | 19916 |

# QUAL对结果的影响

## GAP5

| QUAL    | total | TP    | FP    |
|---------|-------|-------|-------|
| 0-1     | 2049  | 1140  | 909   |
| 1-10    | 4234  | 2312  | 1922  |
| 10-50   | 6298  | 3278  | 3020  |
| 50-100  | 6291  | 2818  | 3473  |
| 100-150 | 5005  | 2450  | 2555  |
| 150-200 | 3262  | 1724  | 1538  |
| 200-250 | 2165  | 1245  | 920   |
| 250-300 | 755   | 465   | 290   |
| 300-350 | 254   | 141   | 113   |
| 350-400 | 136   | 80    | 56    |
| 400-450 | 74    | 47    | 27    |
| 450-500 | 57    | 32    | 25    |
| 500-    | 170   | 73    | 97    |
| total   | 30750 | 15805 | 14945 |

## GATK

| QUAL      | total | TP    | FP    |
|-----------|-------|-------|-------|
| 0-500     | 3982  | 441   | 3541  |
| 500-1000  | 4621  | 1073  | 3548  |
| 1000-1500 | 5930  | 2353  | 3577  |
| 1500-2000 | 7556  | 3905  | 3651  |
| 2000-2500 | 6480  | 3771  | 2709  |
| 2500-3000 | 4090  | 2575  | 1515  |
| 3000-3500 | 1984  | 1320  | 664   |
| 3500-4000 | 1016  | 710   | 306   |
| 4000-4500 | 451   | 317   | 134   |
| 4500-5000 | 256   | 181   | 75    |
| 5000-     | 488   | 292   | 196   |
| total     | 36854 | 16938 | 19916 |

# QUAL/DP对结果的影响

## GAP5

| QUAL/DP | total | TP    | FP   |
|---------|-------|-------|------|
| 0-2     | 18659 | 10493 | 8166 |
| >=2     | 15233 | 7117  | 8116 |
| >=3     | 12091 | 5312  | 6779 |
| >=4     | 8782  | 3441  | 5341 |
| >=5     | 4721  | 1580  | 3141 |
| >=6     | 2790  | 817   | 1973 |
| >=7     | 1935  | 488   | 1447 |
| >=8     | 1399  | 329   | 1070 |
| >=9     | 1038  | 227   | 811  |
| >=10    | 860   | 193   | 667  |
| >=15    | 0     | 0     | 0    |

## GATK

| QUAL/DP | total | TP    | FP    |
|---------|-------|-------|-------|
| 0-2     | 143   | 13    | 130   |
| 2-10    | 1977  | 257   | 1720  |
| >=10    | 34691 | 16665 | 18026 |
| >=20    | 31022 | 15734 | 15288 |
| >=30    | 21380 | 11019 | 10361 |
| >=40    | 9573  | 5222  | 4351  |
| >=50    | 3800  | 2170  | 1630  |
| >=60    | 1881  | 1128  | 753   |
| >=70    | 874   | 542   | 332   |
| >=80    | 482   | 322   | 160   |
| >=90    | 271   | 188   | 83    |
| >=100   | 142   | 104   | 38    |

- 总结： GATK在这两个过滤条件上均表现出一定的规律，即低阈值得到的FP会大于TP。GAP5这方面表现出异常，有待修正。

# 不同过滤结果比较 (GATK)

- SNPs

| Filter             | total         | TP            | FP            | TP:FP      |
|--------------------|---------------|---------------|---------------|------------|
| --                 | <b>351910</b> | <b>120820</b> | <b>231090</b> | <b>0.5</b> |
| <i>Mismatch</i> ≤2 | 165778        | 101828        | 63950         | 1.6        |
| <i>Mismatch</i> ≤1 | 100379        | 75565         | 24814         | 3.0        |
| MQ≥60              | 88711         | 72697         | 16014         | 4.5        |
| Cluster filter     | 87557         | 72616         | 14941         | 4.9        |
| QUAL/DP≥20         | 82333         | 70300         | 12033         | 5.8        |
| QUAL ≥1500         | 20016         | 18917         | 1099          | 17.2       |



# 不同过滤结果比较(GATK)

- small indels

| Filter                | total        | TP           | FP           | TP:FP       |
|-----------------------|--------------|--------------|--------------|-------------|
| --                    | <b>36417</b> | <b>17927</b> | <b>18490</b> | <b>0.9</b>  |
| <i>Mismatch=0</i>     | <b>14447</b> | <b>11304</b> | <b>3143</b>  | <b>3.6</b>  |
| <b>Cluster filter</b> | <b>12489</b> | <b>10955</b> | <b>1534</b>  | <b>7.1</b>  |
| <b>QUAL/DP&gt;=20</b> | <b>11424</b> | <b>10280</b> | <b>1144</b>  | <b>9</b>    |
| <b>QUAL &gt;=1500</b> | <b>3778</b>  | <b>3503</b>  | <b>275</b>   | <b>12.7</b> |
| <b>MQ&gt;=60</b>      | <b>13584</b> | <b>11420</b> | <b>2164</b>  | <b>5.3</b>  |
| <b>FP pool</b>        | <b>14400</b> | <b>11795</b> | <b>2605</b>  | <b>4.5</b>  |
| <i>Mismatch&lt;=1</i> | <b>22324</b> | <b>14583</b> | <b>7741</b>  | <b>1.9</b>  |
| <i>Mismatch&lt;=2</i> | <b>26956</b> | <b>15807</b> | <b>11149</b> | <b>1.4</b>  |
| <i>Mismatch&lt;=3</i> | <b>30038</b> | <b>16289</b> | <b>13749</b> | <b>1.2</b>  |



## Examples of FP and undetected indels in gap5\_indel

- **Example of undetected indels**

INDEL reference:

| CHROM        | POS   | REF    | ALT |
|--------------|-------|--------|-----|
| chromosome04 | 40792 | ACGTTC | -   |

```

}TATTGCATTTTCACGTTTCACGTTTCGTTTCGCTCGTTTCACGTTCCCTAGCTCGTTCAI
 4078Q      4079Q      4080Q      4081Q      4082Q      408:
          CGTTCATGTTTCGTTTCGCTCCTTCACGTTCCCTAGCTCGTTCAI
          CGTTCATGTTTCGTTTCGCTCCTTCACGTTCCCTAGCTCGTTCAI

```

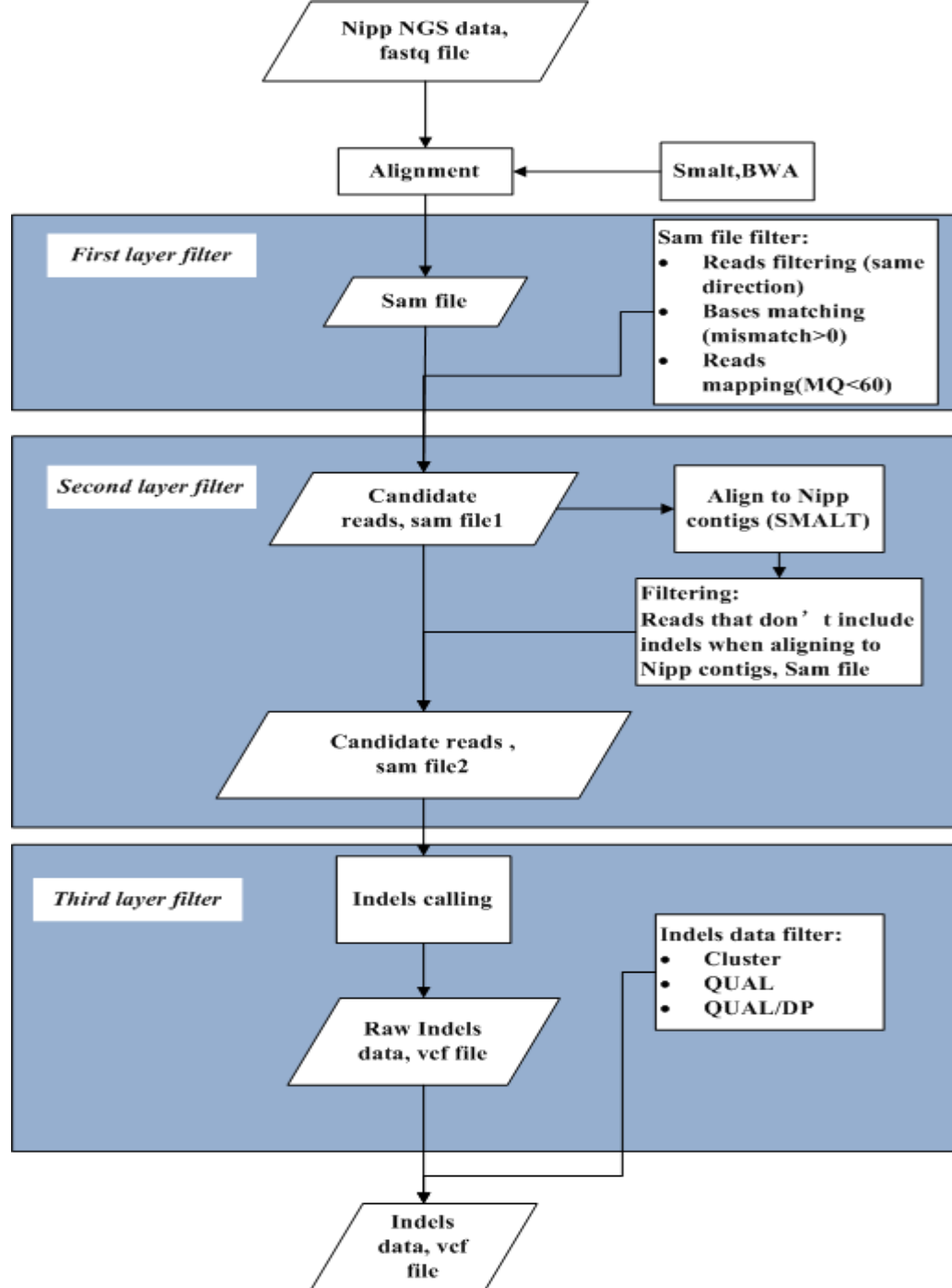
```

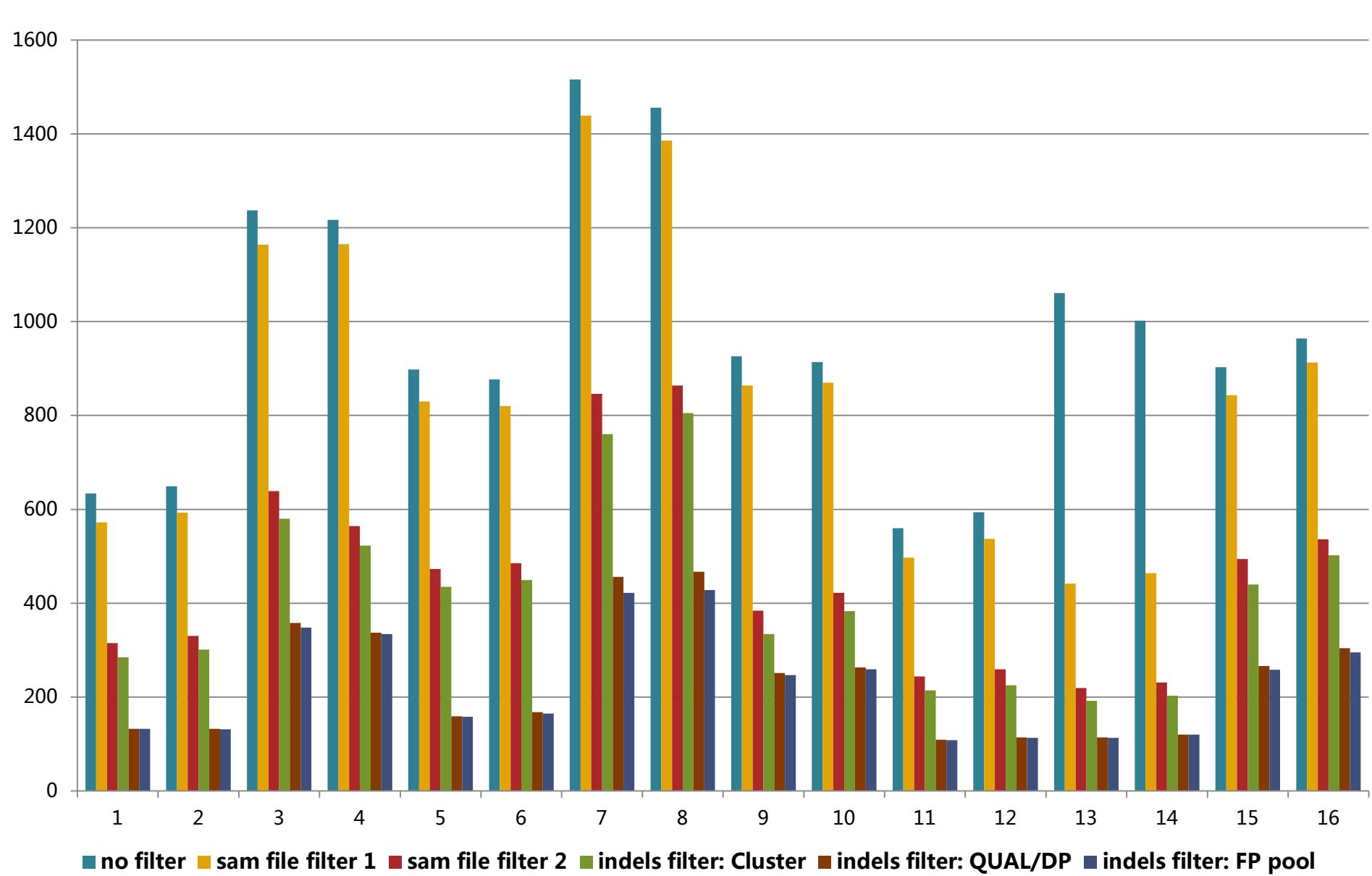
}
}TATT
}TATTGCATTTTCAC
}TATTGCATTTTCACGT
}TATTGCATTTTCACGTTCA
}TATTGCATTTTCACGTTTCAC
}TATTGCATTTTCACGTTTCACG
}TATTGCATTTTCACGTTTCACGTTTC
}TATTGCATTTTCACGTTTCACGTTTCGTTTC
}TATTGCATTTTCACGTTTCACGTTTCGTTTC
}TATTGCATTTTCACGTTTCACGTTTCGTTTC
}TATTGCATTTTCACGTTTCACGTTTCGTTTC
}TATTGCATTTTCACGTTTCACGTTTCGTTTCGCTCGTTCA
}TATTGCATTTTCACGTTTCACGTTTCGTTTCGCTCGTTTCAC
}TATTGCATGTTTCACGTTTCACGTTTCGTTTCGCTCGTTTCACGTTTC
}TATTGCATTTTCACGTTTCACGTTTCGTTTCGCTCGTTTCACGTTTC

```

# 日本晴small indels检测流程

- **目的：** 尽可能准确地检测small indels





**X-coordinate: Sample ID**

**Y-coordinate: Numbers of INDELS**

# 下一步工作

- Detective method of large indels (Structural Variations) from PE read
  - 方法：
    - STEP1: identifying breaking point
    - STEP2: assembly & alignment
  - 重点: breaking point detection
  - 参考算法：
    - Clustering-based algorithm: PINDEL, SOAPindel
    - Distribution-based algorithm: MoDIL
    - Combination algorithm: BreakDancer



谢谢!