

CNV calling, and a set of other useful notes

Vincent Plagnol

Inivata- Head of Computational Biology
UCL- Reader in Statistical Genetics

Longer reads exist if you need them

Query= ch442_file7.twodirections
(1880 letters)

>Descan_9.11.6.27.4.7
Length = 2207

Score = 1552
Identities = 1750/1878 (93%), Gaps = 101/1878 (5%)
Strand = Minus / Plus

Query: 1871 TCGGATCCATTATCTCCCGGAGCTCCATGTGCGAGCGG--TTGTCAATCACTTCATGAG 1813

|||||
Sbjct: 307 TCGGATCCATTATCTCCCGGAGCTCCATGTGCGAGCGGTTGTCAATCACTTCATGAG 366

Query: 1812 GCGGAGATTATGACCGG--GTATGTTATCATAGGGTAAATCTGCCGTTCTGAAAGTGTTCAT 1754

|||||
Sbjct: 367 GCGGAGATTATGACCGAGTATGTTATCAAG--GGTAAATCTGCCGTTCTGAAAGTGTTCAT 425

Query: 1753 CCGCTCCTGCGAGCTCCATTGTGTCGAGTGAAGTCTGGATCATGAGAGGCGAAGT 1694

|||||
Sbjct: 426 CCGCTCCTGTCGTG--CTGATTTTGTTCGAGTGAAGTCTGGATCATGAGAGGCGAAGT 484

Query: 1693 TCTTGTCTTTTCGACAAATTCAG--ATGGAAGATACCTGGTATGTGCTCTGGAGAGCTG 1635

|||||
Sbjct: 485 TCTTGTCTTTTCGACAAATTCAGAGTGAAGATACCTGGTATGTGCTCTGGAGAGCTG 544

Query: 1634 CACATCCGTGAGGTTG--ACCGGAGGAGC--ATACAGAGCTACCAAGTCCGACCAACAAT 1577

|||||
Sbjct: 545 CACATCCGTGAGTGTGACCGGAGGAGCGATACAGAGCTACCAAGTCCGACCAACAAT 604

Query: 1576 GTCTGACCGGAGAAACCCGATCTTAAATGCCAAAAGAGCATGTGTCATCAG--AA 1518

|||||
Sbjct: 605 GTCTGACCGGAGAAACCCGAT--TAAATGCCAAAAGAGCATGTGTCATCAGAGAA 662

Query: 1517 CCTGTGG--CAATGTGCGCTCAATGACAGTGAAGAAATGCTCTGAAA--TSCAGAGCT 1460

|||||
Sbjct: 663 CCTGTGGGAGTGTGCGGCTCAAGATCAGGTGAAGAAATGCTCTGAAAATGCGAGGCT 722

Query: 1459 CGAATTTGCTTCCA--TTACCACTCTTCT--GCCTTGCTCAAGCTTATCAGTACCTCTCT 1403

|||||
Sbjct: 723 CGAATTT--GTTCTCAGCTTCAACCTCTTTTGGCTTGCTCAGGCTTATCAGTACCTCTCT 781

Query: 1402 TAG--ATGTGACA--GTCTCAGAGGAACACCCGTAAA--AGCTCTTGCTGCTCTAAATG 1347

|||||
Sbjct: 782 TAGATGTGTACAAAGTCTCAGAGGAACACCCGTAAACAGCTGT--GTTGCT--AAATG 838

Query: 1346 AC---GTGAGCAGGTCTCTGTGATCTCTCGGCGCCCTATTAAAGAGCGAGTAGGTGCGAG 1290

|||||
Sbjct: 839 ACCGTGTGAGCAGGTCTCTGTGATCTCTCA-----TTATTAAAGAGCGAGTAG--TCGAGG 892

Query: 1289 ACTTCGGAAGGTACTGTGCGGTGTAAACATTTCCGTGGAGGCGGAAGGCGTTGAGACT 1230

|||||
Sbjct: 893 ACTTCGGAAGAA--TACTGTGCGGTGTGAACAAATTTCCGTGGAGGCGGAAGGCGTTGAGACT 951

Query: 1229 GTGCTTACCGTCACTGCGGCC--TATCGGCTAAGATGATGCTTCCGACACAGAGAGTTG 1109

|||||
Sbjct: 952 GTGCTTACCGTCACT--TGCGCCCTTATCGG--TAAAGTCTGATCTTCCGACACAGAGAGTTG 1171

Query: 1170 ACTTCGAGCGCCCGCGCTATTACT--GCAGGTACATCGGAATTCATCAAGACGTTA 1112

Query: 937 GCCTTCCAGGAAGAGACCATGGAACCGGACCAAGTGTATTCTCAAGTGCCTCCCGGC 878

|||||
Sbjct: 1246 GCCTTCCAGGAAGAGACCATGGAACCGGACCAAGTGTATTCTCAAGTGCCTCCCGGC 1305

Query: 877 GGAACCCCAACGCCGAAAT--CTGTG--AACTGCAGCGCACAGAAATGCGCGCGCAAC 820

|||||
Sbjct: 1306 GGAACCCCAACGCCGAAATTTCTGTGGAAGTGCAGGCA--AGAAATGCGC---AACA 1360

Query: 819 ACATGCTATCATAGTGTGAGCATGTATGAGGCTCAAGGAGATGTGGTTCTCTA--CT 761

|||||
Sbjct: 1361 ACATGCTCATCATAGTGTGAGCATGTATG--GAGGCTCAAGGAGATGTGGTTCTCTATCT 1418

Query: 760 CAACATATACCTCGTCCAGCCAAAGATGCG--T--TCTA--AGTGCAC--GCTCAAGAGC 708

|||||
Sbjct: 1419 GAA--CATTAACCTCGTCCAGCCCAAGCATGCGCTCTCTACAGTGCATAGC--CAAGAGC 1476

Query: 707 AAAGTGG--GTGGCGGAACACTGGGGTACATTTGAATGTATATGAGTGCCT--TACATCG 650

|||||
Sbjct: 1477 AAAGTGGGCGTGGCGGAACACTGCGCGCAAGTTGAATGTATATGAGTGCCTTACATCG 1536

Query: 649 CAATGGGAGAGAAAGCAATTTGCTGGAGAACCTGATTTGTCACCATATGCTCTGTA 590

|||||
Sbjct: 1537 CAATGGGAGAGAGAGCAATTTGTTCTGGAGAACCTGATTTGTCAC--ATATGCTCTGTA 1593

Query: 589 CTGAGTACCCCATGAGCATGATGATTTCTGGAGGAG--ATACCGTGGCGCTGCCAT 531

|||||
Sbjct: 1594 GCTGAGTACCCCATCGA--TTGATTTGCTGGAGGAGGATAACCGTGGCGCTGCCAT 1650

Query: 530 CAACCGCAACAGAGAGTCTTCCCAACGGAAGCGTATTTATCGAGATGTGAGACGAG 471

|||||
Sbjct: 1651 CAACCGCAACAGAGAGTCTTCCCAACGGAAGCGTATTTATCGAGATGTGAGAGCG--GA 1709

Query: 470 ACTCAGACCAAGCGACTACA--TGCCTTGCAGGAATCAGAGAGAT--CTCTC--CGCG 415

|||||
Sbjct: 1710 ACTCAGACCAAGCGACTACACTTGCCTTGCAGGAATCAGAGAGATCTCTCTCGAG 1769

Query: 414 GATCTCTGGAAGTGCAGGTCAAG--TGCTCTCAAAATTCGACCTTTGACTCTCGGAGC 356

|||||
Sbjct: 1770 GATCTCTGGAAGTGCAGGTCAAGTGTGCTCTCAAAATTCGACCTTTGACTCTCGGAGC 1829

Query: 355 GAGGCTTCCAACTCAGGCGAAGCATGGGCTTCCTCGGTGACAGGTATCA--GGGAGATCT 297

|||||
Sbjct: 1830 GAGGCTTCCAACTCAGGCGAAGCATGGGCTTC--AGTGCAGGTATCTAAGGGGAGATCT 1888

Query: 296 GCCCATCAATATAGCT--GCTTCAACCAAACTCACACCTGGAACAGTGGTGACTTAGAC 239

|||||
Sbjct: 1889 GCCCATCAATATAGCTGGGCTTC--TAAATATCACACCTGGAACAGTGGTGACTTAGAC 1947

Query: 238 --GTGATCGTCAAGTGTCCAGCAAGTGCAGACCCGGAATATCATCACTACGCGGC 181

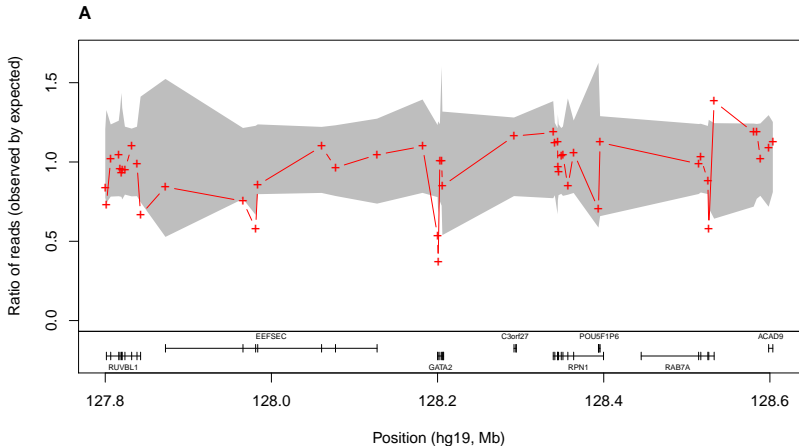
BEDtools: your swiss army knife for all issues

- BEDtools is one of the most widely used tools in bioinformatics.
- It does tons of things, and while many look trivial, together they are very impressive.
- Well worth looking at what it can do, because it may solve many practical questions.
 - In R, most routines are implemented within `GenomicRanges`, a bioconductor package.

Calling CNVs

- There are many tools to call CNVs from sequence data, and I think you have already covered some of them.
- One option is read depth: excess of reads mark a duplication, too few reads mark deletions.
- But there are also specific read patterns, like reads mapping further apart than they should, that can mark a deletion.
- Split read is another way to go about it.

An example of read depth based call in GATA2



Copy number variant analysis

