

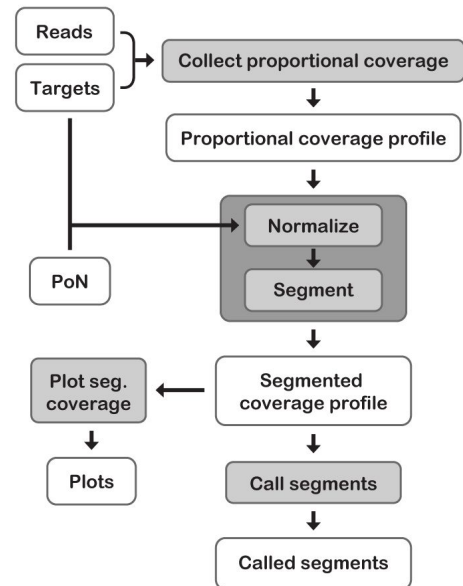
# Call somatic copy number variants using GATK CNV

## About the workflow

This workflow detects *somatic* copy number variation using a panel of normals (PoN). The workflow is optimized for Illumina short-read whole exome sequencing (WES) data. Our team is working on optimizing application to whole genome sequencing (WGS) data and also a separate CNV workflow for germline calling.

The underlying algorithms and current workflow options, e.g. syntax, may change during development. The presented basic approach and general concepts will still be germane. Please check the GATK website (<https://software.broadinstitute.org/gatk/>) for updates. For detailed examples of workflows, see

(<http://gatkforums.broadinstitute.org/gatk/discussion/6791>). For the mathematics behind the workflow, see (<https://github.com/broadinstitute/gatk-protected/blob/master/docs/CNVs/CNV-methods.pdf>).



The workflow uses the following tools from the *gatk-protected-1.0.0.0-alpha1.2.3* pre-release (<https://github.com/broadinstitute/gatk-protected/releases/tag/1.0.0.0-alpha1.2.3>). Note other tools in this release may be unsuitable for analyses.

GATK CNV tools		page
1. <b>CalculateTargetCoverage</b>	This is pre-computed for you in this tutorial.	2
2. <b>CombineReadCounts and CreatePanelOfNormals</b>	Creates the Panel of Normals (PoN) using data from (1).	3
3. <b>NormalizeSomaticReadCounts</b>	Uses data from (1) and (2).	4
4. <b>PerformSegmentation</b>	Uses data from (3).	4
5. <b>PlotSegmentedCopyRatio</b>	Optional visualization step uses data from (3) and (4).	5
6. <b>CallSegments</b>	Uses data from (3) and (4).	6

## About the example data

We have whole exome capture sequence data for chromosomes 1–7 of matched normal and tumor samples. Because the data is from real cancer patients, we have anonymized them at multiple levels. The anonymization process preserves the noise inherent in real samples. The data is representative of Illumina sequencing technology from five years ago (2011).

# 1. Collect proportional coverage using target intervals and read data.

This first step is done for you ahead of time and the command is here for your reference.

Process each BAM, whether normal or tumor. The tool collects coverage per read group at each target and divides these counts by the total number of reads per sample.

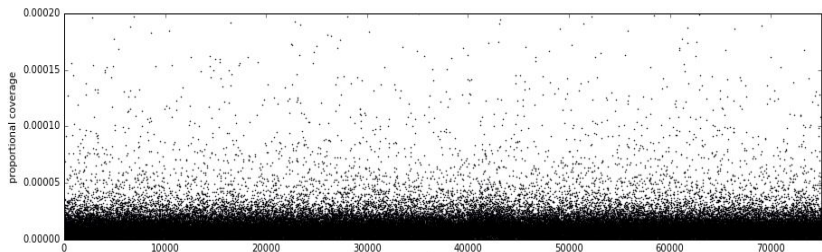
```

VNDM dVM`SMVWEVM` ,MOaXMIQMSQ` [bQ^MSQ`Æ
 ¼( `i UZ\`a` QNMVCRUXQ` `Æ
·
 ¼( `i UZ\`a` Ç` MSQ` Ç` _bi` `Æ
·
 ¼` ^M_R[ ^Y` $` #*` `Æ
·
 ¼S^[ a\`e`` `! $zi` `Æ
·
 ¼` MSQ` ŁZR[ ` /) žž` `Æ
·
 ¼WQPa\` `Æ
·
 ¼#` i [ a\` \a` ÇQ[ bCRUXQ`

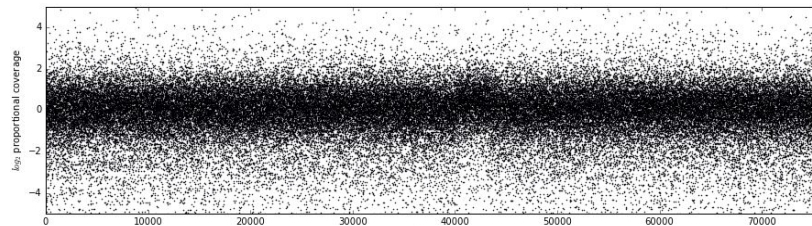
```

The target file (¼()) is a *padded* intervals list of the baited regions. You can add padding to a target list using the PadTargets tool. For us, padding each exome target 250bp on either side increases sensitivity. The ¼` MSQ` ŁZR[ ` /) žž` option keeps the original target names from the target list. The ¼WQPa\` option asks the tool to include alignments flagged as duplicate.

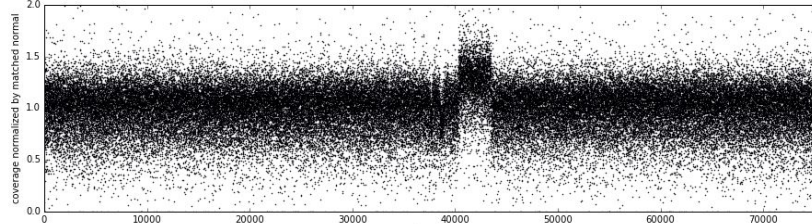
The top plot shows the raw proportional coverage for our tumor sample for chromosomes 1–7. Each dot represents a target. The y-axis plots proportional coverage and the x-axis targets.



The middle plot shows the data after a median-normalization and log2-transformation.



The bottom plot shows the tumor data after normalization against its matched-normal.



→ For each of these progressions, how certain are you that there are copy-number events? How many copy-number variants are you certain of? What is contributing to your uncertainty?

## 2. Create the CNV panel of normals (PoN) with two commands.

The normals should represent the same sequencing technology, e.g. sample preparation and capture target kit, as that of the tumor samples under scrutiny. The PoN is meant to encapsulate sequencing noise and may also capture common germline variants. Like any control, you should think carefully about what sample set would make an effective panel. At the least, the PoN should consist of ten normal samples that were ideally subject to the same batch effects as that of the tumor sample, e.g. from the same sequencing center. Our current recommendation is 40 or more normal samples. Depending on the coverage depth of samples, adjust the number.

→ *What is better, tissue-matched normals or blood normals of tumor samples?*

→ *What makes a better background control, a matched normal sample or a panel of normals?*

The first step combines the proportional read counts from the multiple normal samples into a single file. The `cnvtools a`zU`` parameter takes a file listing the relative file paths, one sample per line, of the proportional coverage data of the normals.

```
cnvtools a`zU` [YNUZQ&QMP] [aZ`_`Æ
cnvtools a`zU` Z[ ^YMK_E`d`Æ
##`_NMPN[d³Q YNUZQPZ[ ^YMK_E`_b`
```

The second step creates a single CNV PoN file. The PoN stores information such as the median proportional coverage per target across the panel and projections of systematic noise calculated with PCA (principal component analysis). Our tutorial's PoN is built with 39 normal blood samples from cancer patients from the same cohort (not suffering from blood cancers).

```
cnvtools a`zU` SMWUEVM` ^QMO$MZOXR` [ ^YMK`Æ
##`_NMPN[d³Q YNUZQPZ[ ^YMK_E`_b`Æ
##`_NMPN[d³Z[ ^YMK_E\Z`Æ
Z[ %`Æ
##PU_MMXQ \MWÆ
##YUZUYaY(MSQ/MO [ ^SQ^OOZ` UXQ(T^Q_T[XP`U`
```

This results in two files, the CNV PoN and a `MSQ cccUST`_E`d`` file that typical workflows can ignore. Because we have a small number of normals, we include the `Z[ %`` option and change the `##YUZUYaY(MSQ/MO [ ^SQ^OOZ` UXQ(T^Q_T[XP`` to 5%.

→ *Based on what you know about PCA (<http://setosa.io/ev/principal-component-analysis/>), what do you think are the effects of using more normal samples? A panel with some profiles that are outliers?*

### 3. Normalize a raw coverage profile using the PoN.

We normalize the tumor coverage against the PoN's target medians and against the principal components of the PoN.

```
WVbM #VM`SMWWEVM` "[^YVXUFQ [YMUO&QMP, [aZ`_`Æ
#` Q[b^` aY[ ^E`_b`Æ
#` "_MPN[d^` Z[ ^YVX`E\`Z`Æ
#`("` "_MPN[d^` aY[ ^E` ZE`_b`Æ
#`("` "_MPN[d^` aY[ ^E` ZE`_b`
```

This produces the *pre-tangent-normalized* file (`#`("``) and the *tangent-normalized* file (`#`("``), respectively. Resulting data is log2-transformed.

Denosing with a PoN is critical for calling copy-number variants from WES coverage profiles. It can also improve calls from WGS profiles that are typically more evenly distributed and subject to less noise. Furthermore, denosing with a PoN can greatly impact results for (i) samples that have more noise, e.g. those with lower coverage, lower purity or higher activity, (ii) samples lacking a matched normal and (iii) detection of smaller events that span only a few targets.

---

### 4. Segment the normalized coverage profile.

Segmentation groups contiguous targets with the same copy ratio.

```
WVbM #VM`SMWWEVM` $Q^R[ ^Y` OSYQZ` MU[Z`Æ
#`("` "_MPN[d^` aY[ ^E` ZE`_b`Æ
#` "_MPN[d^` aY[ ^E` OS`Æ
#`#fi`
```

For our tumor sample, we reduce the ~73K individual targets to 14 segments. The `#fi`` parameter tells the tool that the input coverages are log2-transformed.

---

This command will error if you have not installed R and certain R components. If you need, take a few minutes to install R from <https://www.r-project.org/>. Then install the components with the following command.

```
&_O^U` `UZ`_`MXX&C\MVMSQ`E&`
```

We include `UZ`_`MXX&C\MVMSQ`E&`` in the tutorial data bundle. Alternatively, download it from ([https://github.com/broadinstitute/gatk-protected/blob/master/scripts/install\\_R\\_packages.R](https://github.com/broadinstitute/gatk-protected/blob/master/scripts/install_R_packages.R)).

---

View the resulting file with `cat *_MNP[d3`aY[^E_OS`.

Sample	Chromosome	Start	End	Num_Probes	Segment_Mean
tumor	1	1	249238921	17585	0.99647119
tumor	2	1	243193925	13159	0.995918783
tumor	3	1	21484649	1139	0.999307093
tumor	3	21502423	22485593	54	0.798851916
tumor	3	22503367	92855531	3740	0.997231251
tumor	3	92873305	96090725	174	0.816146473
tumor	3	96108499	109371107	705	0.994470169
tumor	3	109388881	115519763	326	0.816995479
tumor	3	115537537	141901787	1390	0.994263396
tumor	3	141919561	198008273	3016	1.146312186
tumor	4	1	191131706	6818	0.998268633
tumor	5	1	180899061	7849	0.999514915
tumor	6	1	171097489	8860	0.995090745
tumor	7	1	159125275	8166	0.996402123

→ Which chromosomes have events?

## 5. (Optional) Plot segmented coverage.

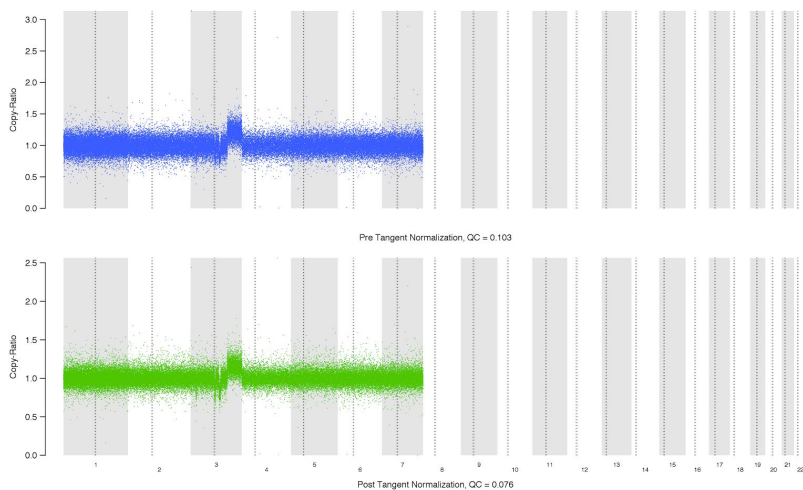
This command requires XQuartz <<https://www.xquartz.org/>> installation. If you do not have this dependency, then view the results in the precomputed\_results folder instead. Currently plotting only supports human assembly b37 autosomes. Going forward, this tool will accommodate other references and the workflow will support calling on sex chromosomes.

```

VMM $VM`SMWWEVM`$X[ ` ` OSYQZ`QP. [\e&MU[ `Æ
  #(" `*_MNP[d3`aY[^E`ZE`_b`Æ
  #("$(" `*_MNP[d3`aY[^E\`ZE`_b`Æ
  #`*_MNP[d3`aY[^E_OS`Æ
  ##`*_MNP[d`Æ
  #^Q`aY[^`Æ
  #z#fi`
  .
  
```

The `##`` defines the output directory, and the `#^Q`` defines the basename of the files. Again, the `#z#fi`` parameter tells the tool that the inputs are log2-transformed. The output folder contains seven files--three PNG images and four text files.

- `OR[^Q`R`Q^EVS` (shown above) plots copy-ratios pre (top) and post (bottom) tangent-normalization across the chromosomes.



The plot automatically adjusts the y-axis to show all available data points. Dotted lines represent centromeres.

- `~/QC/chr1/chr1.qc` shows the same but fixes the y-axis range from 0 to 4 for comparability across samples.
- `~/QC/chr1/chr1.qc.colors` colors differential copy-ratio segments in alternating blue and orange. The horizontal line plots the segment mean. Again the y-axis ranges from 0 to 4.

→ Open each of these images. How many copy-number variants do you see?

Each of the four text files contain a single quality control (QC) value. This value is the median of absolute differences (MAD) in copy-ratios of adjacent targets. Its calculation is robust to actual copy-number variants and should decrease post tangent-normalization.

- `~/QC/chr1/chr1.qc` gives the QC value before tangent-normalization.
- `~/QC/chr1/chr1.qc.post` gives the post-tangent-normalization QC value.
- `~/QC/chr1/chr1.qc.diff` gives the difference between pre and post QC values.
- `~/QC/chr1/chr1.qc.fraction` gives the fraction difference  $(preQc - postQc)/(preQc)$ .

## 6. Call segmented copy number variants.

This final step makes one of three calls for each segment--neutral (0), deletion (-) or amplification (+). These deleted or amplified segments could represent somatic events.

```
~/QC/chr1/chr1.qc.colors > ./callCNV.py chr1.qc.colors
# chr1:1-249238921 17585 0.99647119 0
# chr1:243193925-21484649 1139 0.999307093 0
# chr1:21502423-22485593 54 0.798851916 -
# chr1:22503367-9285531 3740 0.997231251 0
# chr1:92873305-96090725 174 0.816146473 -
# chr1:96108499-109371107 705 0.994470169 0
# chr1:109388881-115519763 326 0.816995479 -
# chr1:115537537-141901787 1390 0.994263396 0
# chr1:141919561-198008273 3016 1.146312186 +
# chr1:191131706-180899061 6818 0.998268633 0
# chr1:180899061-171097489 7849 0.999514915 0
# chr1:171097489-159125275 8860 0.995090745 0
# chr1:159125275- 8166 0.996402123 0
```

View the results with `~/QC/chr1/chr1.qc.colors > ./callCNV.py chr1.qc.colors`.

Sample	Chromosome	Start	End	Num_Probes	Segment_Mean	Segment_Call
tumor	1	1	249238921	17585	0.99647119	0
tumor	2	1	243193925	13159	0.995918783	0
tumor	3	1	21484649	1139	0.999307093	0
tumor	3	21502423	22485593	54	0.798851916	-
tumor	3	22503367	9285531	3740	0.997231251	0
tumor	3	92873305	96090725	174	0.816146473	-
tumor	3	96108499	109371107	705	0.994470169	0
tumor	3	109388881	115519763	326	0.816995479	-
tumor	3	115537537	141901787	1390	0.994263396	0
tumor	3	141919561	198008273	3016	1.146312186	+
tumor	4	1	191131706	6818	0.998268633	0
tumor	5	1	180899061	7849	0.999514915	0
tumor	6	1	171097489	8860	0.995090745	0
tumor	7	1	159125275	8166	0.996402123	0

→ Besides the last column, how is this result different from that of step 4?

