

Copy number calling and SNV classification using targeted short read sequencing

Markus Riester¹

¹Novartis Institutes for BioMedical Research, Cambridge, MA

January 4, 2019

Abstract

PureCN [1] is a purity and ploidy aware copy number caller for cancer samples inspired by the *ABSOLUTE* algorithm [2]. It was designed for hybrid capture sequencing data, especially with medium-sized targeted gene panels without matching normal samples in mind (matched whole-exome data is of course supported).

It can be used to supplement existing normalization and segmentation algorithms, i.e. the software can start from BAM files, from target-level coverage data, from copy number log-ratios or from already segmented data. After the correct purity and ploidy solution was identified, *PureCN* will accurately classify variants as germline vs. somatic or clonal vs. sub-clonal.

PureCN was further designed to integrate well with industry standard pipelines [3], but it is straightforward to generate input data from other pipelines.

Package

PureCN 1.12.2

Contents

1	Introduction	3
2	Basic input files	3
2.1	VCF	3
2.2	Target information	4
2.3	Coverage data	5
2.4	Third-party coverage tools	5
2.5	Third-party segmentation tools	5
2.6	Example data	5
3	Library-specific coverage bias	6
4	Pool of normals	7
4.1	Selection of normals for log-ratio calculation	7

Copy number calling and SNV classification using targeted short read sequencing

4.2	Artifact filtering	8
4.2.1	VCF	8
4.2.2	Coverage data	9
4.3	Artifact filtering without a pool of normals	9
5	Recommended run	10
6	Output	10
6.1	Plots	10
6.2	Data structures	15
6.2.1	Prediction of somatic status and cellular fraction	15
6.2.2	Amplifications and deletions	16
6.2.3	Find genomic regions in LOH	19
7	Curation	19
8	Cell lines	20
9	Maximizing the number of heterozygous SNPs	21
10	Advanced usage	21
10.1	Custom normalization and segmentation	21
10.1.1	Custom segmentation	21
10.1.2	Custom normalization	23
10.2	COSMIC annotation	23
10.3	ExAC and gnomAD annotation	23
10.4	Mutation burden	24
10.5	Detect cross-sample contamination	24
10.6	Power to detect somatic mutations	25
11	Limitations	28
12	Support	28
12.1	Checklist	28
12.2	FAQ	29

1 Introduction

This tutorial will demonstrate on a toy example how we recommend running *PureCN* on targeted sequencing data. To estimate tumor purity, we jointly utilize both target-level¹ coverage data and allelic fractions of single nucleotide variants (SNVs), inside - and optionally outside - the targeted regions. Knowledge of purity will in turn allow us to accurately (i) infer integer copy number and (ii) classify variants (somatic vs. germline, mono-clonal vs. sub-clonal, heterozygous vs. homozygous etc.).

¹The captured genomic regions, e.g. exons.

This requires 3 basic input files:

1. A VCF file containing germline SNPs and somatic mutations. Somatic status is not required in case the variant caller was run without matching normal sample.
2. The tumor BAM file.
3. At least one BAM file from a normal control sample, either matched or process-matched.

In addition, we need to know a little bit more about the assay. This is the annoying step since here the user needs to provide some information. Most importantly, we need to know the positions of all targets. Then we need to correct for GC-bias, for which we need GC-content for each target. Optionally, if gene-level calls are wanted, we also need for each target a gene symbol. We may also observe subtle differences in coverage of tumor compared to normal due to varying proliferation rates and we can provide replication timing data to check and correct for such a bias. To obtain best results, we can finally use a pool of normal samples to automatically learn more about the assay and its biases and common artifacts.

The next sections will show how to do all this with *PureCN* alone or with the help of *GATK* and/or existing copy number pipelines.

All the steps described in the following are available in easy to use command line scripts described in a separate vignette.

2 Basic input files

2.1 VCF

Germline SNPs and somatic mutations are expected in a single VCF file. At the bare minimum, this VCF should contain read depths of reference and alt alleles in an AD format field and a DB info flag for membership in germline databases².

Without DB flag, variant ids starting with `rs` are assumed to be in dbSNP. Population allele frequencies are expected in a POP_AF info field. These frequencies are currently only used to infer membership in germline databases when the DB flag is missing; in future versions they will be used calculate somatic prior probabilities more accurately.

If a matched normal is available, then somatic status information is currently expected in a SOMATIC info flag in the VCF. The *VariantAnnotation* package provides examples how to add info fields to a VCF in case the used variant caller does not add this flag. If the VCF contains a BQ format field containing base quality scores, *PureCN* can remove low quality calls.

²The name of the flag is customizable in *runAbsoluteCN*

Copy number calling and SNV classification using targeted short read sequencing

VCF files generated by *MuTect* [4] should work well and in general require no post-processing. *PureCN* can handle *MuTect* VCF files generated in both tumor-only and matched normal mode. Experimental support for *MuTect 2* and *FreeBayes* VCFs generated in tumor-only mode is available.

2.2 Target information

For the default segmentation function provided by *PureCN*, the algorithm first needs to calculate log-ratios of tumor vs. normal control coverage. To do this, we need to know the locations of the captured genomic regions (targets). These are provided by the manufacturer of your capture kit³. Please double check that the genome version of the target file matches the reference. Usually the manufacturer provides two files: the baits file containing the coordinates of the actual capture baits, and the target file containing the coordinates of the actual regions we want to capture. We recommend to use the baits file (and recognize the confusing nomenclature that we follow due to convention in established tools).

Default parameters assume that these targets do NOT include a "padding" to include flanking regions. *PureCN* will automatically include variants in the 50bp flanking regions if the variant caller was either run without interval file or with interval padding (See section 12.2).

PureCN will attempt to optimize the targets for copy number calling (similar to [5]):

- Large targets are split to obtain a higher resolution
- Targets in regions of low mappability are dropped
- Optionally, accessible regions in-between the target (off-target) regions are included so that available coverage information in on- and off-target reads can be used by the segmentation function. In the following, we will use *intervals* when something applies to both on-target and off-target regions and *targets* when it only applies to on-target.

It further annotates intervals by GC-content (how coverage is normalized is described later in Section 3).

PureCN provides the `preprocessIntervals` function:

```
reference.file <- system.file("extdata", "ex2_reference.fa",
  package = "PureCN", mustWork = TRUE)
bed.file <- system.file("extdata", "ex2_intervals.bed",
  package = "PureCN", mustWork = TRUE)
mappability.file <- system.file("extdata", "ex2_mappability.bigWig",
  package = "PureCN", mustWork = TRUE)

intervals <- import(bed.file)
mappability <- import(mappability.file)

preprocessIntervals(intervals, reference.file,
  mappability=mappability, output.file = "ex2_gc_file.txt")

## WARN [2019-01-04 22:59:59] No repriming scores provided.
## INFO [2019-01-04 22:59:59] Calculating GC-content...
```

A command line script described in a separate vignette provides convenient access to this function and also attempts to annotate the targets with gene symbols using the `annotateTargets` function.

³While *PureCN* can use a pool of normal samples to learn which intervals are reliable and which not, it is highly recommended to provide the correct intervals. Garbage in, garbage out.

2.3 Coverage data

The `calculateBamCoverageByInterval` function can be used to generate the required coverage data from BAM files. All we need to do is providing the desired intervals (as generated by `preprocessIntervals`):

```
bam.file <- system.file("extdata", "ex1.bam", package="PureCN",
  mustWork = TRUE)
interval.file <- system.file("extdata", "ex1_intervals.txt",
  package = "PureCN", mustWork = TRUE)

calculateBamCoverageByInterval(bam.file = bam.file,
  interval.file = interval.file, output.file = "ex1_coverage.txt")
```

2.4 Third-party coverage tools

Calculating coverage from BAM files is a common task and your pipeline might already provide this information. As alternative to `calculateBamCoverageByInterval`, *PureCN* currently supports coverage files generated by *GATK3 DepthOfCoverage*, *GATK4 CollectFragmentCounts* and by *CNVkit*. By providing files with standard file extension, *PureCN* will automatically detect the correct format and all following steps are the same for all tools. You will, however, still need the interval file generated in Section 2.2 and the third-party tool must use the exact same intervals. See also FAQ Section 12.2 for recommended settings for *GATK3 DepthOfCoverage*.

2.5 Third-party segmentation tools

PureCN integrates well with existing copy number pipelines. Instead of coverage data, the user then needs to provide either already segmented data or a wrapper function. This is described in Section 10.1.

2.6 Example data

We now load a few example files that we will use throughout this tutorial:

```
library(PureCN)

normal.coverage.file <- system.file("extdata", "example_normal.txt",
  package="PureCN")
normal2.coverage.file <- system.file("extdata", "example_normal2.txt",
  package="PureCN")
normal.coverage.files <- c(normal.coverage.file, normal2.coverage.file)
tumor.coverage.file <- system.file("extdata", "example_tumor.txt",
  package="PureCN")
seg.file <- system.file("extdata", "example_seg.txt",
  package = "PureCN")
vcf.file <- system.file("extdata", "example.vcf.gz", package="PureCN")
```

```
interval.file <- system.file("extdata", "example_intervals.txt",  
                             package="PureCN")
```

3 Library-specific coverage bias

In coverage normalization, we distinguish between assay- and library-specific biases. Assay-specific biases, for example due to probe density, probe capture efficiency and read mappability, are best removed with a pool of normal samples (Section 4.1, [5]). In other words, by examining the coverage of particular intervals in a pool of normals, we can estimate how well this assay captures these intervals and will then adjust the tumor coverage accordingly.

Other biases are library-specific, meaning a patient sample captured in different libraries may display dramatically different coverage profiles across libraries. Data from great sequencing centers usually show relatively small technical variance nowadays, but some biases are not completely avoidable. The most important library-specific bias is due to GC-content, i.e. regions of high AT- or GC-content are not always captured with exactly the same efficiency in tumor and normals.

We usually also observe that early replicating regions have a slightly higher coverage than late replicating regions [6, 7]. Since there is often a significant difference in proliferation rates of tumor and normal, the pool of normals might also not completely adjust for this small bias.

As first step, we thus correct the raw coverage of all samples, tumor and normal, for these two major sources of library-specific coverage biases (Figure 1). For GC-normalization, we use a 2-step loess normalization [8]. For the replication timing bias, a linear model of log-transformed coverage and provided replication timing score is used.

```
correctCoverageBias(normal.coverage.file, interval.file,  
                    output.file="example_normal_loess.txt", plot.bias=TRUE)
```

All the following steps in this vignette assume that the coverage data are normalized. The example coverage files are already GC-normalized. We provide a convenient command line script for generating normalized coverage data from BAM files or from *GATK* coverage files (see Quick vignette).

Copy number calling and SNV classification using targeted short read sequencing

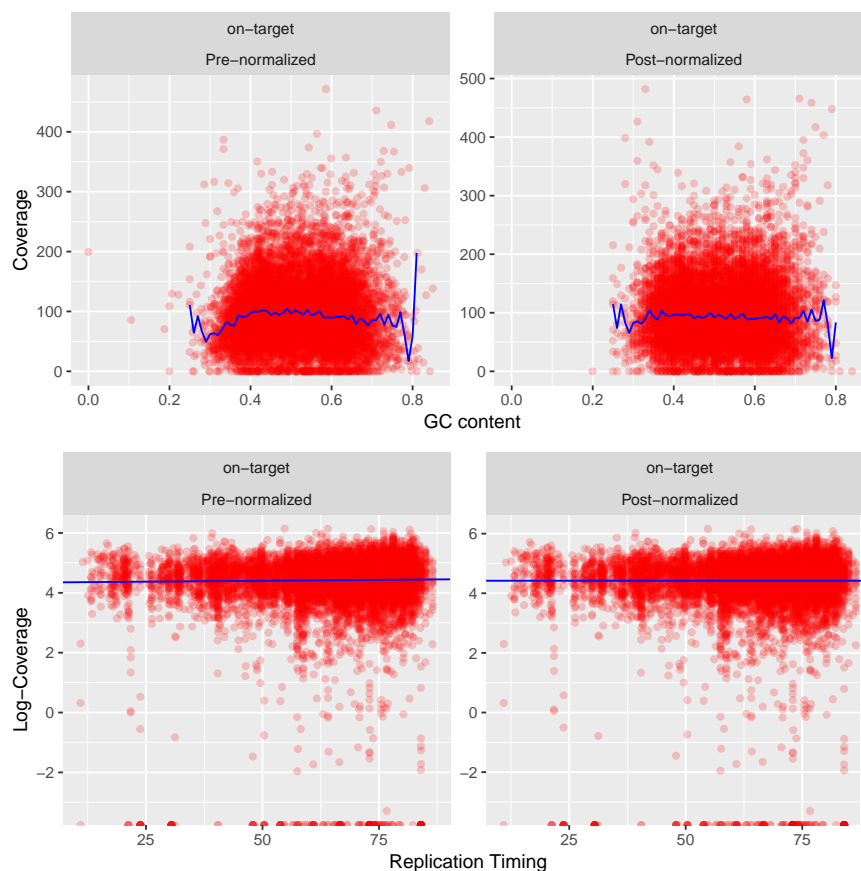


Figure 1: Coverage before and after normalization for GC-content and replication timing

This plot shows coverage as a function of on- and off-target GC-content and replication timing before and after normalization. Each dot is an interval. The example files are already GC-normalized; real data will show more dramatic differences.

4 Pool of normals

4.1 Selection of normals for log-ratio calculation

For calculating copy number log-ratios of tumor vs. normal, *PureCN* requires coverage from a process-matched normal sample. Using a normal that was sequenced using a similar, but not identical assay, rarely works. Differently covered genomic regions simply result in too many log-ratio outliers. This section describes how to identify good process-matched normals in case no matched normal is available or in case the matched normal has low or uneven coverage.

The `createNormalDatabase` function builds a database of coverage files (a command line script providing this functionality is described in a separate vignette):

```
normalDB <- createNormalDatabase(normal.coverage.files)

## INFO [2019-01-04 23:00:12] 576 on-target bins with low coverage in all samples.
## WARN [2019-01-04 23:00:12] You are likely not using the correct baits file!
```

Copy number calling and SNV classification using targeted short read sequencing

```
## WARN [2019-01-04 23:00:12] Allosome coverage missing, cannot determine sex.
## WARN [2019-01-04 23:00:12] Allosome coverage missing, cannot determine sex.
## INFO [2019-01-04 23:00:12] Processing on-target regions...
## INFO [2019-01-04 23:00:12] Removing 930 intervals with low coverage in normalDB.
## INFO [2019-01-04 23:00:12] Removing 1 intervals with zero coverage in more than 3% of normalDB.

# serialize, so that we need to do this only once for each assay
saveRDS(normalDB, file="normalDB.rds")
```

Again, please make sure that all coverage files were GC-normalized prior to building the database (Section 3). Internally, `createNormalDatabase` determines the sex of the samples and trains a PCA that is later used for denoising tumor coverage:

```
normalDB <- readRDS("normalDB.rds")
pool <- calculateTangentNormal(tumor.coverage.file, normalDB)
```

4.2 Artifact filtering

It is important to remove as many artifacts as possible, because low ploidy solutions are typically punished more by artifacts than high ploidy solutions. High ploidy solutions are complex and usually find ways of explaining artifacts reasonably well. The following steps in this section are optional, but recommended since they will reduce the number of samples requiring manual curation, especially when matching normal samples are not available.

4.2.1 VCF

We recommend running *MuTect* with a pool of normal samples to filter common sequencing errors and alignment artifacts from the VCF. *MuTect* requires a single VCF containing all normal samples, for example generated by the *GATK3 CombineVariants* tool (see Section 12.2).

It is highly recommended to provide *PureCN* this combined VCF as well; it will help the software correcting non-reference read mapping biases. This is described in the `setMappingBiasVcf` documentation. To reduce memory usage, the normal panel VCF can be reduced to contain only variants present in 4 or more samples (the VCF for *MuTect* should however contain variants present in 2-3 samples).

Because these VCFs can become huge with large pools of normals, we can optionally pre-compute the mapping bias, thus avoiding parsing these VCFs for every sample:

```
# speed-up future runtimes by pre-calculating variant mapping biases
normal.panel.vcf.file <- system.file("extdata", "normalpanel.vcf.gz",
                                     package="PureCN")

bias <- calculateMappingBiasVcf(normal.panel.vcf.file, genome = "h19")
## INFO [2019-01-04 23:00:14] Processing variants 1 to 5000...

saveRDS(bias, "mapping_bias.rds")
normal.panel.vcf.file <- "mapping_bias.rds"
```


Copy number calling and SNV classification using targeted short read sequencing

4.2.2 Coverage data

We next use coverage data of normal samples to estimate the expected variance in coverage per interval:

```
interval.weight.file <- "interval_weights.txt"
calculateIntervalWeights(normalDB$normal.coverage.files, interval.weight.file)

## INFO [2019-01-04 23:00:15] Loading coverage data...
## INFO [2019-01-04 23:00:15] Mean target coverages: 71X (tumor) 99X (normal).
## INFO [2019-01-04 23:00:15] Mean target coverages: 71X (tumor) 43X (normal).
```

This function calculates copy number log-ratios using all normal samples provided in the `normal.coverage.files` argument. Assuming that all normal samples are in general diploid, a high variance in log-ratio is indicative of an interval with either common germline alterations or frequent artifacts; high or low copy number log-ratios in these intervals are unlikely measuring somatic copy number events.

This `interval.weight.file` is automatically generated by the `NormalDB.R` script described in the Quick vignette.

4.3 Artifact filtering without a pool of normals

By default, *PureCN* will exclude targets with coverage below 15X from segmentation (with a pool of normals, targets are filtered based on the coverage and variance in normal database only). For variants in the provided VCF, the same 15X cutoff is applied. *MuTect* applies more sophisticated artifact tests and flags suspicious variants. If *MuTect* was run in matched normal mode, then both potential artifacts and germline variants are rejected, that means we cannot just filter by the PASS/REJECT *MuTect* flags. The `filterVcfMuTect` function optionally reads the *MuTect* 1.1.7 stats file and will keep germline variants, while removing potential artifacts. Without the stats file, *PureCN* will use only the filters based on read depths as defined in `filterVcfBasic`. Both functions are automatically called by *PureCN*, but can be easily modified and replaced if necessary.

We can also use a BED file to blacklist regions expected to be problematic, for example the simple repeats track from the UCSC:

```
# Instead of using a pool of normals to find low quality regions,
# we use suitable BED files, for example from the UCSC genome browser.

# We do not download these in this vignette to avoid build failures
# due to internet connectivity problems.
downloadFromUCSC <- FALSE
if (downloadFromUCSC) {
  library(rtracklayer)
  mySession <- browserSession("UCSC")
  genome(mySession) <- "hg19"
  simpleRepeats <- track(ucscTableQuery(mySession,
    track="Simple Repeats", table="simpleRepeat"))
  export(simpleRepeats, "hg19_simpleRepeats.bed")
}
```

```
snp.blacklist <- "hg19_simpleRepeats.bed"
```

5 Recommended run

Finally, we can run *PureCN* with all that information:

```
ret <- runAbsoluteCN(normal.coverage.file=pool,
# normal.coverage.file=normal.coverage.file,
  tumor.coverage.file=tumor.coverage.file, vcf.file=vcf.file,
  genome="hg19", sampleid="Sample1",
  interval.file=interval.file, normalDB=normalDB,
# args.setMappingBiasVcf=list(normal.panel.vcf.file=normal.panel.vcf.file),
# args.filterVcf=list(snp.blacklist=snp.blacklist,
# stats.file=mutect.stats.file),
  args.segmentation=list(interval.weight.file=interval.weight.file),
  post.optimize=FALSE, plot.cnv=FALSE, verbose=FALSE)

## WARN [2019-01-04 23:00:16] Allosome coverage missing, cannot determine sex.
## WARN [2019-01-04 23:00:16] Allosome coverage missing, cannot determine sex.
```

The `normal.coverage.file` argument points to a coverage file obtained from either a matched or a process-matched normal sample, but can be also a small pool of best normals (Section 4.1).

The `normalDB` argument (Section 4.1) provides a pool of normal samples and for example allows the segmentation function to skip targets with low coverage or common germline deletions in the pool of normals. If available, a VCF containing all variants from the normal samples should be provided via `args.setMappingBiasVcf` to correct read mapping biases. The files specified in `args.filterVcf` help *PureCN* filtering SNVs more efficiently for artifacts as described in Sections 4.2 and 4.3. The `snp.blacklist` is only necessary if neither a matched normal nor a large pool of normals is available.

The `post.optimize` flag will increase the runtime by about a factor of 2-5, but might return slightly more accurate purity estimates. For high quality whole-exome data, this is typically not necessary for copy number calling (but might be for variant classification, see Section 6.2.1). For smaller targeted panels, the runtime increase is typically marginal and `post.optimize` should be always set to `TRUE`.

The `plot.cnv` argument allows the segmentation function to generate additional plots if set to `TRUE`. Finally, `verbose` outputs important and helpful information about all the steps performed and is therefore set to `TRUE` by default.

6 Output

6.1 Plots

We now create a few output files:

Copy number calling and SNV classification using targeted short read sequencing

```
file.rds <- "Sample1_PureCN.rds"
saveRDS(ret, file=file.rds)
pdf("Sample1_PureCN.pdf", width=10, height=11)
plotAbs(ret, type="all")
dev.off()

## pdf
## 2
```

The RDS file now contains the serialized return object of the `runAbsoluteCN` call. The PDF contains helpful plots for all local minima, sorted by likelihood. The first plot in the generated PDF is displayed in Figure 2 and shows the purity and ploidy local optima, sorted by final likelihood score after fitting both copy number and allelic fractions.

```
plotAbs(ret, type="overview")
```

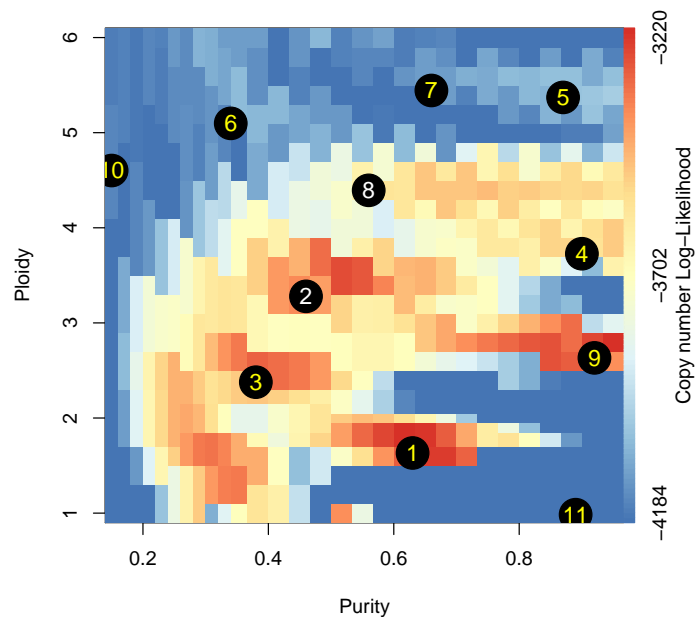


Figure 2: Overview

The colors visualize the copy number fitting score from low (blue) to high (red). The numbers indicate the ranks of the local optima. Yellow fonts indicate that the corresponding solutions were flagged, which does not necessarily mean the solutions are wrong. The correct solution (number 1) of this toy example was flagged due to large amount of LOH.

We now look at the main plots of the maximum likelihood solution in more detail.

```
plotAbs(ret, 1, type="hist")

## NULL
```

Figure 3 displays a histogram of tumor vs. normal copy number log-ratios for the maximum likelihood solution (number 1 in Figure 2). The height of a bar in this plot is proportional to the fraction of the genome falling into the particular log-ratio copy number range. The vertical dotted lines and numbers visualize the, for the given purity/ploidy combination,

Copy number calling and SNV classification using targeted short read sequencing

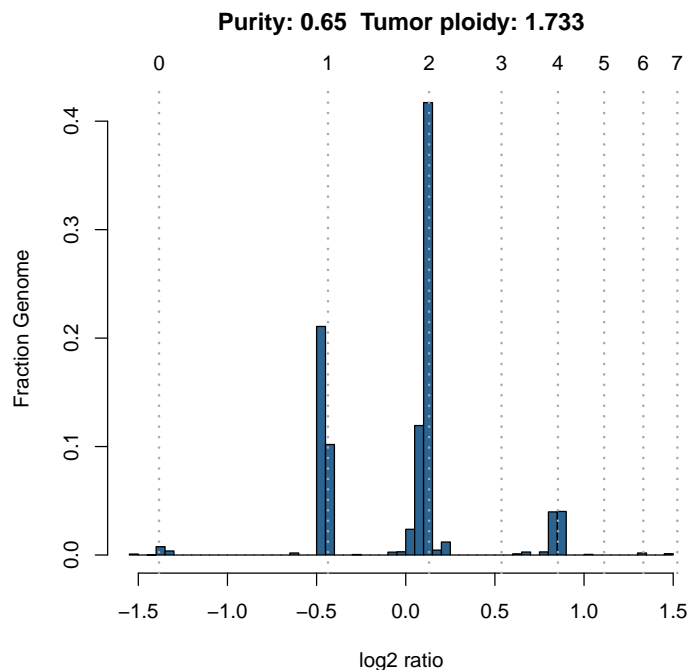


Figure 3: Log-ratio histogram

expected log-ratios for all integer copy numbers from 0 to 7. It can be seen that most of the log-ratios of the maximum likelihood solution align well to expected values for copy numbers of 0, 1, 2 and 4.

```
plotAbs(ret, 1, type="BAF")
```

Germline variant data are informative for calculating integer copy number because unbalanced maternal and paternal chromosome numbers in the tumor portion of the sample lead to unbalanced germline allelic fractions. Figure 4 shows the allelic fractions of predicted germline SNPs. The goodness of fit (GoF) is provided on an arbitrary scale in which 100% corresponds to a perfect fit and 0% to the worst possible fit. The latter is defined as a fit in which allelic fractions on average differ by 0.2 from their expected fractions. Note that this does not take purity into account and low purity samples are expected to have a better fit. In the middle panel, the corresponding copy number log-ratios are shown. The lower panel displays the calculated integer copy numbers, corrected for purity and ploidy. We can zoom into particular chromosomes (Figure 5).

```
plotAbs(ret, 1, type="BAF", chr="chr19")
```

```
plotAbs(ret, 1, type="AF")
```

Finally, Figure 6 provides more insight into how well the variants fit the expected values.

Copy number calling and SNV classification using targeted short read sequencing

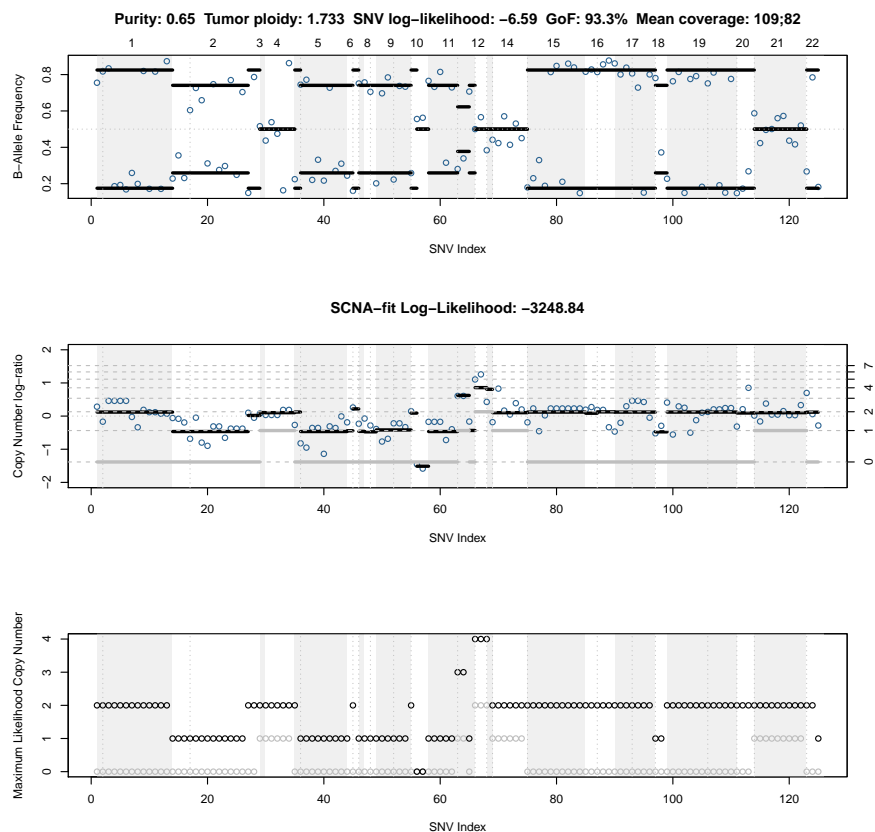


Figure 4: B-allele frequency plot

Each dot is a (predicted) germline SNP. The first panel shows the allelic fractions as provided in the VCF file. The alternating grey and white background colors visualize odd and even chromosome numbers, respectively. The black lines visualize the expected (not the average!) allelic fractions in the segment. These are calculated using the estimated purity and the total and minor segment copy numbers. These are visualized in black and grey, respectively, in the second and third panel. The second panel shows the copy number log-ratios, the third panel the integer copy numbers.

Copy number calling and SNV classification using targeted short read sequencing

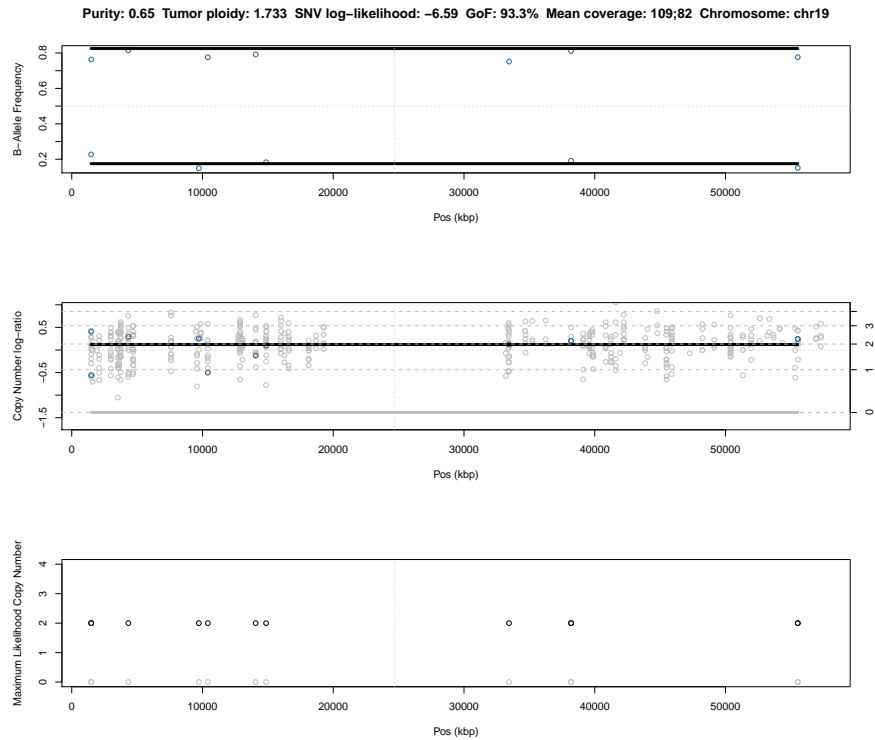


Figure 5: Chromosome plot

Similar to Figure 4, but zoomed into a particular chromosome. The grey dots in the middle panel visualize copy number log-ratios of targets without heterozygous SNPs, which are omitted from the previous genome-wide plot. The x-axis now indicates genomic coordinates in kbps.

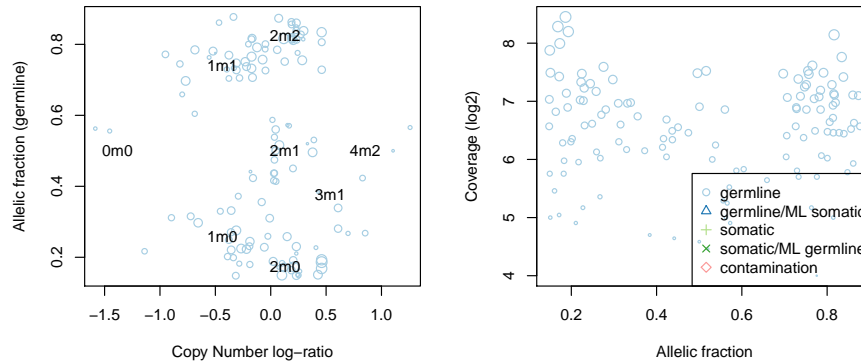


Figure 6: Allele fraction plots

Each dot is again a (predicted) germline SNP. The size of dots indicate quality, defined here as the product of mapping bias and coverage. The shapes visualize the different SNV groups based on prior and posterior probabilities. The labels show the expected values for all called states; 2m1 would be diploid, heterozygous, 2m2 diploid, homozygous. The relationship of allelic fraction and coverage is shown in the top right panel. This plot normally also shows somatic mutations in two additional panels, with the left panel showing the same plot as for germline SNPs and the bottom right panel a histogram of cellular fraction of predicted somatic mutations. This toy example contains only germline SNPs however.

6.2 Data structures

The R data file (`file.rds`) contains gene-level copy number calls, SNV status and LOH calls. The purity/ploidy combinations are sorted by likelihood and stored in `ret$results`.

```
names(ret)

## [1] "candidates" "results" "input"
```

We provide convenient functions to extract information from this data structure and show their usage in the next sections. We recommend using these functions instead of accessing the data directly since data structures might change in future versions.

6.2.1 Prediction of somatic status and cellular fraction

To understand allelic fractions of particular SNVs, we must know the (i) somatic status, the (ii) tumor purity, the (iii) local copy number, as well as the (iv) number of chromosomes harboring the mutations or SNPs. One of *PureCN* main functions is to find the most likely combination of these four values. We further assign posterior probabilities to all possible combinations or states. Availability of matched normals reduces the search space by already providing somatic status.

The `predictSomatic` function provides access to these probabilities. For predicted somatic mutations, this function also provides cellular fraction estimates, i.e. the fraction of tumor cells with mutation. Fractions significantly below 1 indicate sub-clonality⁴:

```
head(predictSomatic(ret), 3)
```

##	chr	start	end	ID	REF	ALT	SOMATIC.M0	SOMATIC.M1	SOMATIC.M2	SOMATIC.M3	SOMATIC.M4	SOMATIC.M5	SOMATIC.M6	SOMATIC.M7	GERMLINE.M0	GERMLINE.M1	GERMLINE.M2	GERMLINE.M3	GERMLINE.M4	GERMLINE.M5	GERMLINE.M6	GERMLINE.M7	GERMLINE.CONTHIGH	GERMLINE.CONTLOW	GERMLINE.HOMOZYGOUS	ML.SOMATIC	POSTERIOR.SOMATIC	ML.M	ML.C	ML.M.SEGMENT	M.SEGMENT.POSTERIOR	M.SEGMENT.FLAGGED	ML.AR	AR
## 1	chr1	114515871	114515871	chr1114515871xxx	G	A	5.969924e-101																											
## 2	chr1	150044293	150044293	chr1150044293xxx	T	G	7.667563e-91																											
## 3	chr1	158449835	158449835	chr1158449835xxx	A	G	4.162785e-146																											
## 1		4.511273e-38	4.183033e-07	3.512481e-268																														
## 2		1.887887e-38	4.276546e-10	2.828745e-270																														
## 3		1.939661e-61	1.303470e-14	3.954647e-272																														
## 1		0	1.574106e-68	2.173694e-15			0.9999996	2.459725e-266																										
## 2		0	2.307980e-63	2.728936e-18			1.0000000	1.386484e-264																										
## 3		0	2.291025e-104	2.254699e-29			1.0000000	3.766019e-264																										
## 1		0	0	0			0	6.174341e-43																										
## 2		0	0	0			0	3.072618e-21																										
## 3		0	0	0			0	5.875696e-27																										
## 1		1.263476e-286		0		FALSE	4.183033e-07	2																										
## 2		1.395456e-241		0		FALSE	4.276546e-10	2																										
## 3		0.000000e+00		0		FALSE	1.303470e-14	2																										
## 1	2	0		1		FALSE	0.825	0.755183																										
## 2	2	0		1		FALSE	0.825	0.817078																										
## 3	2	0		1		FALSE	0.825	0.834266																										

⁴This number can be above 1 when the observed allelic fraction is higher than expected for a clonal mutation. This may be due to random sampling, wrong copy number, sub-clonal copy number events, or wrong purity/ploidy estimates.

Copy number calling and SNV classification using targeted short read sequencing

```
## AR.ADJUSTED MAPPING.BIAS ML.LOH CN.SUBCLONAL CELLFRACTION FLAGGED
## 1 0.7736009 0.976192 TRUE FALSE NA FALSE
## 2 0.8370054 0.976192 TRUE FALSE NA FALSE
## 3 0.8546126 0.976192 TRUE FALSE NA FALSE
## log.ratio depth prior.somatic prior.contamination on.target seg.id
## 1 0.2842136 184 9.9e-05 0.01 1 1
## 2 -0.1686186 138 9.9e-05 0.01 1 1
## 3 0.4596841 217 9.9e-05 0.01 1 1
## gene.symbol
## 1 HIPK1
## 2 VPS45
## 3 <NA>
```

The output columns are explained in Table 1.

To annotate the input VCF file with these values:

```
vcf <- predictSomatic(ret, return.vcf=TRUE)
writeVcf(vcf, file="Sample1_PureCN.vcf")
```

For optimal classification results:

- Set `post.optimize=TRUE` since small inaccuracies in purity can decrease the classification performance significantly
- Provide `args.setMappingBias` a pool of normal VCF to obtain position-specific mapping bias information
- Exclude variants in regions of low mappability
- Use a somatic posterior probability cutoff of 0.8 and 0.2 for somatic and germline variants, respectively. This appears to be a good compromise of call rate and accuracy. If the beta-binomial model was selected in the `model` argument of `runAbsoluteCN`, these cutoffs might need to be relaxed to get acceptable call rates.
- Add a `Cosmic.CNT` info field to the VCF or provide a COSMIC VCF in `runAbsoluteCN` (see Section 10.2).

Note that the posterior probabilities assume that the purity and ploidy combination is correct. Before classifying variants, it is thus recommended to manually curate flagged samples.

6.2.2 Amplifications and deletions

To call amplifications, we recommend using a cutoff of 6 for focal amplifications and a cutoff of 7 otherwise. For homozygous deletions, a cutoff of 0.5 is useful to allow some heterogeneity in copy number.

For samples that failed `PureCN` calling we recommended using common log-ratio cutoffs to call amplifications, for example 0.9.

This strategy is implemented in the `callAlterations` function:

```
gene.calls <- callAlterations(ret)
head(gene.calls)
```


Copy number calling and SNV classification using targeted short read sequencing

Table 1: `predictSomatic` output columns

Column name	Description
<code>chr, start, end</code>	Variant coordinates
<code>ID</code>	The variant ID as provided in the VCF
<code>REF, ALT</code>	The reference and alt alleles
<code>SOMATIC.M*</code>	Posterior probabilities for all somatic states. M0 to M7 are multiplicity values, i.e. the number of chromosomes harboring the mutation (e.g. 1 heterozygous, 2 homozygous if copy number C is 2). <code>SOMATIC.M0</code> represents a sub-clonal state (somatic mutations by definition have a multiplicity larger than 0).
<code>GERMLINE.M*</code>	Posterior probabilities for all heterozygous germline states
<code>GERMLINE.CONTHIGH</code>	Posterior probability for contamination. This state corresponds to homozygous germline SNPs that were not filtered out because reference alleles from another individual were sequenced, resulting in allelic fractions smaller than 1.
<code>GERMLINE.CONTLOW</code>	Posterior probability for contamination. This state corresponds to non-reference alleles only present in the contamination.
<code>GERMLINE.HOMOZYGOUS</code>	Posterior probability that SNP is homozygous in normal. Requires the <code>model.homozygous</code> option in <code>runAbsoluteCN</code> . See Section 8.
<code>ML.SOMATIC</code>	TRUE if the maximum likelihood state is a somatic state
<code>POSTERIOR.SOMATIC</code>	The posterior probability that the variant is somatic (sum of all somatic state posterior probabilities)
<code>ML.M</code>	Maximum likelihood multiplicity
<code>ML.C</code>	Maximum likelihood integer copy number
<code>ML.M.SEGMENT</code>	Maximum likelihood minor segment copy number
<code>M.SEGMENT.POSTERIOR</code>	Posterior probability of <code>ML.M.SEGMENT</code>
<code>M.SEGMENT.FLAGGED</code>	Segment flag indicating <code>ML.M.SEGMENT</code> is unreliable, either due to low posterior probability (< 0.5) or few variants ($< \text{min.variants.segment}$). Indels are always flagged.
<code>ML.AR</code>	Expected allelic fraction of the maximum likelihood state
<code>AR</code>	Observed allelic fraction (as provided in VCF)
<code>AR.ADJUSTED</code>	Observed allelic fraction adjusted for mapping bias
<code>ML.LOH</code>	TRUE if variant is most likely in LOH
<code>CN.SUBCLONAL</code>	TRUE if copy number segment is sub-clonal
<code>CELLFRACTION</code>	Fraction of tumor cells harboring the somatic mutation
<code>FLAGGED</code>	Flag indicating that call is unreliable (currently only due to high mapping bias and high pool of normal counts)
<code>log.ratio</code>	The copy number log-ratio (tumor vs. normal) for this variant
<code>depth</code>	The total sequencing depth at this position
<code>prior.somatic</code>	Prior probability of variant being somatic
<code>prior.contamination</code>	Prior probability that variant is contamination from another individual
<code>on.target</code>	1 for variants within intervals, 2 for variants in flanking regions, 0 for off-target variants
<code>seg.id</code>	Segment id
<code>pon.count</code>	Number of hits in the <code>normal.panel.vcf.file</code>
<code>gene.symbol</code>	Gene symbol if available

Copy number calling and SNV classification using targeted short read sequencing

```
##      chr      start      end C seg.mean seg.id number.targets
## EIF2A   chr3 150270143 150301699 6   1.3351     5           13
## MIR548H2 chr3 151531951 151538241 6   1.3351     5           3
## GPNMB    chr7 23286477 23313844 7   1.4723    17          11
## NRBF2   chr10 64911910 64913979 0  -1.4318    22           2
## SH2D4B  chr10 82298088 82403838 0  -1.5118    26           8
## SLC35G1  chr10 95653791 95661248 0  -1.3002    28           3
##      gene.mean  gene.min  gene.max focal breakpoints      type
## EIF2A   1.5159387 0.7544062 2.139217 TRUE           0 AMPLIFICATION
## MIR548H2 0.8261594 0.4682628 1.200025 TRUE           0 AMPLIFICATION
## GPNMB    1.4490939 0.9397072 1.820078 TRUE           0 AMPLIFICATION
## NRBF2   -1.4318247 -1.4415146 -1.422135 FALSE          0 DELETION
## SH2D4B  -1.5141982 -1.9281007 -1.260723 FALSE          0 DELETION
## SLC35G1  -1.5988526 -1.8066816 -1.449356 FALSE          0 DELETION
##      num.snps.segment  M M.flagged loh
## EIF2A                 0 NA      NA   NA
## MIR548H2              0 NA      NA   NA
## GPNMB                 0 NA      NA   NA
## NRBF2                 0 NA      NA   NA
## SH2D4B                 2 0      TRUE TRUE
## SLC35G1              0 NA      NA   NA
```

It is also often useful to filter the list further by known biology, for example to exclude non-focal amplifications of tumor suppressor genes. The Sanger Cancer Gene Census [9] for example provides such a list.

The output columns of `callAlterations` are explained in Table 2.

Table 2: `callAlterations` output columns

Column name	Description
<code>chr, start, end</code>	Gene coordinates
<code>C</code>	Segment integer copy number
<code>seg.mean</code>	Segment mean of copy number log-ratios (not adjusted for purity/ploidy)
<code>seg.id</code>	Segment id
<code>number.targets</code>	Number of targets annotated with this symbol
<code>gene.*</code>	Gene copy number log-ratios (not adjusted for purity/ploidy)
<code>focal</code>	TRUE for focal amplification, as defined by the <code>fun.focal</code> argument in <code>runAbsoluteCN</code>
<code>breakpoints</code>	Number of breakpoints between <code>start</code> and <code>end</code>
<code>type</code>	Amplification or deletion
<code>num.snps.segment</code>	Number of SNPs in this segment informative for LOH detection (requires VCF)
<code>M</code>	minor copy number of segment (requires VCF)
<code>M.flagged</code>	flag indicating that <code>M</code> is unreliable (requires VCF)
<code>loh</code>	TRUE if segment is in LOH, FALSE if not and NA if number of SNPs is insufficient (requires VCF)

6.2.3 Find genomic regions in LOH

The `gene.calls` data.frame described above provides gene-level LOH information. To find the corresponding genomic regions in LOH, we can use the `callLOH` function:

```
loh <- callLOH(ret)
head(loh)
```

##	chr	start	end	arm	C	M	type
## 1	chr1	114515871	121535434	p 2	0	WHOLE ARM	COPY-NEUTRAL LOH
## 2	chr1	124535434	248085104	q 2	0	WHOLE ARM	COPY-NEUTRAL LOH
## 3	chr2	10262881	92326171	p 1	0	WHOLE ARM	LOH
## 4	chr2	95326171	231775678	q 1	0	WHOLE ARM	LOH
## 5	chr2	236403331	239039169	q 2	0	COPY-NEUTRAL	LOH
## 6	chr3	11888043	90504854	p 2	1		

The output columns are explained in Table 3.

Table 3: `callLOH` output columns

Column name	Description
chr, start, end	Segment coordinates
arm	Chromosome arm
C	Segment integer copy number
M	Minor integer copy number ($M + N = C, M \leq N$)
type	LOH type if $M = 0$

7 Curation

For prediction of variant status (germline vs. somatic, sub-clonal vs. clonal, homozygous vs. heterozygous), it is important that both purity and ploidy are correct. We provide functionality for curating results:

```
createCurationFile(file.rds)
```

This will generate a CSV file in which the correct purity and ploidy values can be manually entered. It also contains a column "Curated", which should be set to `TRUE`, otherwise the file will be overwritten when re-run.

Then in R, the correct solution (closest to the combination in the CSV file) can be loaded with the `readCurationFile` function:

```
ret <- readCurationFile(file.rds)
```

This function has various handy features, but most importantly it will re-order the local optima so that the curated purity and ploidy combination is ranked first. This means `plotAbs(ret, 1, type="hist")` would show the plot for the curated purity/ploidy combination, for example.

The default curation file will list the maximum likelihood solution:

Copy number calling and SNV classification using targeted short read sequencing

```
read.csv("Sample1_PureCN.csv")

## Sampleid Purity Ploidy Sex Contamination Flagged Failed Curated
## 1 Sample1 0.65 1.733021 ? 0 TRUE FALSE FALSE
##
## Comment
## 1 EXCESSIVE LOSSES;EXCESSIVE LOH
```

[PureCN](#) currently only flags samples with warnings, it does not mark any samples as failed. The `Failed` column in the curation file can be used to manually flag samples for exclusion in downstream analyses. See Table 4 for an explanation of all flags.

Table 4: `createCurationFile` flags

Flag	Description
EXCESSIVE LOH	> 50% of genome in LOH and ploidy < 2.6
EXCESSIVE LOSSES	$\geq 1\%$ of genome deleted
HIGH AT- OR GC-DROPOUT	High GC-bias exceeding cutoff in <code>max.dropout</code>
HIGH PURITY	(when <code>model.homozygous=FALSE</code>). For very high purity samples, it is recommended to set <code>model.homozygous=TRUE</code> . See Section 8.
LOW PURITY	Purity < 30%
LOW BOOTSTRAP VALUE	<code>bootstrapResults</code> identified multiple plausible solutions
NOISY LOG-RATIO	Log-ratio standard deviation > <code>max.logr.sdev</code>
NOISY SEGMENTATION	More than <code>max.segments</code>
NON-ABERRANT	$\geq 99\%$ of genome has identical copy number and $\geq 0.5\%$ has second most common state
POLYGENOMIC	$\geq 0.75 \times$ <code>max.non.clonal</code> fraction of the genome in sub-clonal state
POOR GOF	GoF < <code>min.gof</code>
POTENTIAL SAMPLE CONTAMINATION	Significant portion of dbSNP variants potentially cross-contaminated
RARE KARYOTYPE	Ploidy < 1.5 or > 4.5

8 Cell lines

Default parameters assume some normal contamination. In 100% pure samples without matching normal samples such as cell lines, we cannot distinguish homozygous SNPs from LOH by looking at single allelic fractions alone. It is thus necessary to keep homozygous variants and include this uncertainty in the likelihood model. This is done by setting the `runAbsoluteCN` argument `model.homozygous=TRUE`. If matched normals are available, then variants homozygous in normal are automatically removed since they are uninformative.

For technical reasons, the maximum purity *PureCN* currently models is 0.99. We recommend setting `test.purity=seq(0.9,0.99,by=0.01)` in `runAbsoluteCN` for cell lines.

Please note that in order to detect homozygous deletions in 100% pure samples, you will need to provide a `normalDB` in `runAbsoluteCN` to filter low quality targets efficiently (Section 5).

9 Maximizing the number of heterozygous SNPs

It is possible to use SNPs in off-target reads in the variant fitting step by running *MuTect* without interval file and then setting the `filterVcfBasic` argument `remove.off.target.snvs` to `FALSE`. We recommend a large pool of normals in this case and then generating SNP blacklists as described in Sections 4.2 and 4.3. Remember to also run all the normals in *MuTect* without interval file.

An often better alternative to including all off-target reads is only including variants in the flanking regions of targets (between 50-100bp). This will usually significantly increase the number of heterozygous SNPs (see Section 12.2). These SNPs are automatically added if the variant caller was run without interval file or with interval padding.

10 Advanced usage

10.1 Custom normalization and segmentation

Copy number normalization and segmentation are crucial for obtaining good purity and ploidy estimates. If you have a well-tested pipeline that produces clean results for your data, you might want to use *PureCN* as add-on to your pipeline. By default, we will use *DNAcopy* [10] to segment normalized target-level coverage log-ratios. It is straightforward to replace the default with other methods and the `segmentationCBS` function can serve as an example.

The next section describes how to replace the default segmentation. For the probably more uncommon case that only the coverage normalization is performed by third-party tools, see Section 10.1.2.

10.1.1 Custom segmentation

It is possible to provide already segmented data, which is especially recommended when matched SNP6 data are available or when third-party segmentation tools are not written in R. Otherwise it is usually however better to customize the default segmentation function, since the algorithm then has access to the raw log-ratio distribution⁵. The expected file format for already segmented copy number data is⁶:

ID	chrom	loc.start	loc.end	num.mark	seg.mean
Sample1	1	61723	5773942	2681	0.125406444072723
Sample1	1	5774674	5785170	10	-0.756511807441712

Since its likelihood model is exon-based, *PureCN* currently still requires an interval file to generate simulated target-level log-ratios from a segmentation file. For simplicity, this interval file is expected either via the `tumor.coverage.file` or via the `interval.file` argument (see

⁵If the third-party tool provides target-level log-ratios, then these can be provided via the `log.ratio` argument in addition to `seg.file` though. See also Section 10.1.2.

⁶This segmentation file can contain multiple samples, in which case the provided `sampleid` must match a sample in the column ID

Copy number calling and SNV classification using targeted short read sequencing

Figure 7). Note that *PureCN* will re-segment the simulated log-ratios using the default *segmentationCBS* function, in particular to identify regions of copy-number neutral LOH and to cluster segments with similar allelic imbalance and log-ratio. The provided interval file should therefore cover all significant copy number alterations⁷. Please check that the log-ratios are similar to the ones obtained by the default *PureCN* segmentation and normalization.

⁷If this behaviour is not wanted, because maybe the custom function already identifies CN/LOH reliably, *segmentationCBS* can be replaced with a minimal version.

```
retSegmented <- runAbsoluteCN(seg.file=seg.file,
  interval.file=interval.file, vcf.file=vcf.file,
  max.candidate.solutions=1, genome="hg19",
  test.purity=seq(0.3,0.7,by=0.05), verbose=FALSE,
  plot.cnv=FALSE)
```

```
## WARN [2019-01-04 23:01:44] Allosome coverage missing, cannot determine sex.
## WARN [2019-01-04 23:01:44] Allosome coverage missing, cannot determine sex.
```

The *max.candidate.solutions* and *test.purity* arguments are set to non-default values to reduce the runtime of this vignette.

```
plotAbs(retSegmented, 1, type="BAF")
```

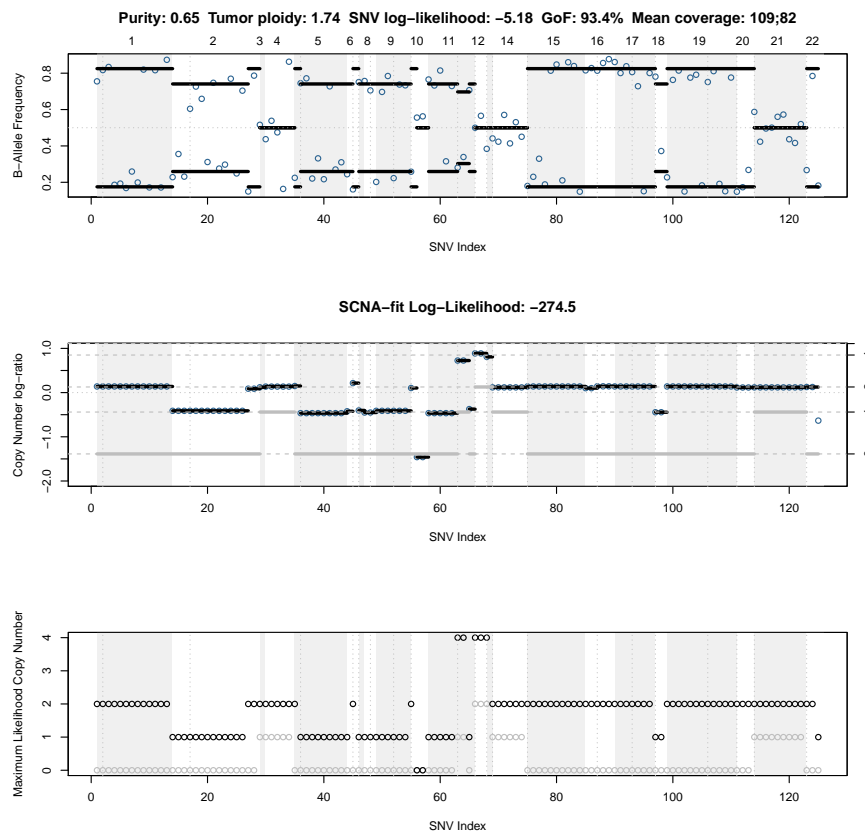


Figure 7: B-allele frequency plot for segmented data

This plot shows the maximum likelihood solution for an example where segmented data are provided instead of coverage data. Note that the middle panel shows no variance in log-ratios, since only segment-level log-ratios are available.

10.1.2 Custom normalization

If third-party tools such as *GATK4* are used to calculate target-level copy number log-ratios, and *PureCN* should be used for segmentation and purity/ploidy inference only, it is possible to provide these log-ratios:

```
# We still use the log-ratio exactly as normalized by PureCN for this
# example
log.ratio <- calculateLogRatio(readCoverageFile(normal.coverage.file),
                             readCoverageFile(tumor.coverage.file))

retLogRatio <- runAbsoluteCN(log.ratio=log.ratio,
                             interval.file=interval.file, vcf.file=vcf.file,
                             max.candidate.solutions=1, genome="hg19",
                             test.purity=seq(0.3,0.7,by=0.05), verbose=FALSE,
                             normalDB=normalDB, plot.cnv=FALSE)

## WARN [2019-01-04 23:02:05] Allosome coverage missing, cannot determine sex.
## WARN [2019-01-04 23:02:05] Allosome coverage missing, cannot determine sex.
```

Again, the `max.candidate.solutions` and `test.purity` arguments are set to non-default values to reduce the runtime of this vignette. It is highly recommended to compare the log-ratios obtained by *PureCN* and the third-party tool, since some pipelines automatically adjust log-ratios for a default purity value. Note that this example uses a pool of normals to filter low quality targets. Interval coordinates are again expected in either a `interval.file` or a `tumor.coverage.file`. If a tumor coverage file is provided, then all targets below the coverage minimum are further excluded.

10.2 COSMIC annotation

If a matched normal is not available, it is also helpful to provide `runAbsoluteCN` the COSMIC database [11] via `cosmic.vcf.file` (or via a `Cosmic.CNT` INFO field in the VCF). While this has limited effect on purity and ploidy estimation due the sparsity of hotspot mutations, it often helps in the manual curation to compare how well high confidence germline (dbSNP) vs. somatic (COSMIC) variants fit a particular purity/ploidy combination.

For variant classification (Section 6.2.1), providing COSMIC annotation also avoids that hotspot mutations with dbSNP id get assigned a very low prior probability of being somatic.

10.3 ExAC and gnomAD annotation

PureCN is not automatically annotating input VCFs with data from common germline databases such as ExAC. See Section 2.1 for ways to tell *PureCN* where to find either a summary binary flag (i.e. likely germline yes/no) or population allele frequencies.

10.4 Mutation burden

The `predictSomatic` function described in Section 6.2.1 can be used to efficiently remove private germline mutations. This in turn allows the calculation of mutation burden for unmatched tumor samples. A wrapper function for this specific task is included as `callMutationBurden`:

```
callableBed <- import(system.file("extdata", "example_callable.bed.gz",
  package = "PureCN"))

callMutationBurden(ret, callable=callableBed)

## somatic.ontarget somatic.all private.germline.ontarget
## 1 0 0 0
## private.germline.all callable.bases.ontarget callable.bases.flanking
## 1 0 1511056 2223991
## callable.bases.all somatic.rate.ontarget somatic.rate.ontarget.95.lower
## 1 3063762 0 0
## somatic.rate.ontarget.95.upper private.germline.rate.ontarget
## 1 2.441256 0
## private.germline.rate.ontarget.95.lower
## 1 0
## private.germline.rate.ontarget.95.upper
## 1 2.441256
```

The `callableBed` file should be a file parsable by `rtracklayer`. This file can specify genomic regions that are callable, for example as obtained by *GATK3 CallableLoci*. This is optional, but if provided can be used to accurately calculate mutation rates per megabase. Variants outside the callable regions are not counted. Private germline rates should be fairly constant across samples; outliers here should be manually inspected.

The output columns are explained in Table 5.

10.5 Detect cross-sample contamination

It is important to correctly handle heterozygous common SNPs that do not have an expected allelic fraction of 0.5 in normal samples. These can be SNPs in poor quality regions (as already described, see Section 4.2.1), but also SNPs from cross-sample contaminated DNA. Without matched normals, detection of those problematic SNPs is not trivial.

For cross-sample contamination, *PureCN* by default always tests for a 1% contamination and assigns common SNPs to a contamination state when allelic fractions are either close to 0 or close to 1 and when this cannot be explained by CNAs. The main purpose of these states is to provide a bin for common SNPs that for artifactual reasons do not fit any other state.

This tool applies a simple heuristic to flag samples for cross-contamination: Given the coverage and putative contamination rate based on allelic fractions of potentially contaminated SNPs, how many SNPs do we expect to detect based on our power to detect variants at that contamination rate? If the expected number is much higher than observed, then significant contamination is unlikely; observed SNPs close to 0 or 1 are more likely artifacts or the contamination rate is much lower than the minimum tested. Otherwise *PureCN* will perform a post-optimization in which contamination rate is optimized in additional variant fitting steps.

Copy number calling and SNV classification using targeted short read sequencing

Table 5: `callMutationBurden` output columns

Column name	Description
<code>somatic.ontarget</code>	Number of mutations in target regions
<code>somatic.all</code>	Total number of mutations. Depending on input VCF and <code>runAbsoluteCN</code> arguments, this might include calls in flanking regions and off-targets reads.
<code>private.germline.ontarget</code>	Number of private germline SNPs in targets. For matched tumor/normal samples, this is the number of variants annotated as somatic, but classified as germline.
<code>private.germline.all</code>	Total number of private germline SNPs
<code>callable.bases.ontarget</code>	Number of callable on-target bases
<code>callable.bases.flanking</code>	Number of callable on-target and flanking bases
<code>callable.bases.all</code>	Total number of callable bases. With default parameters includes off-target regions that were ignored by <code>runAbsoluteCN</code> .
<code>somatic.rate.ontarget</code>	Somatic mutations per megabase in target regions
<code>somatic.rate.ontarget.95.lower</code>	Lower 95% of confidence interval
<code>somatic.rate.ontarget.95.upper</code>	Upper 95% of confidence interval
<code>private.germline.rate.ontarget</code>	Private germline mutations per megabase in target regions
<code>private.germline.rate.ontarget.95.lower</code>	Lower 95% of confidence interval
<code>private.germline.rate.ontarget.95.upper</code>	Upper 95% of confidence interval

Cross-sample contamination can also result in increased observed heterozygosity on chrX for males, which in turn often results in a `PureCN` warning that sex inferred from coverage and VCF are in conflict.

By default, cross-contamination is tested in the range from 1 to 7.5%. Catastrophic failures with higher contamination might not get flagged.

10.6 Power to detect somatic mutations

As final quality control step, we can test if coverage and tumor purity are sufficient to detect mono-clonal or even sub-clonal somatic mutations. We strictly follow the power calculation by Carter et al. [2].

The following Figure 8 shows the power to detect mono-clonal somatic mutations as a function of tumor purity and sequencing coverage (reproduced from [2]):

```
purity <- c(0.1,0.15,0.2,0.25,0.4,0.6,1)
coverage <- seq(5,35,1)
power <- lapply(purity, function(p) sapply(coverage, function(cv)
  calculatePowerDetectSomatic(coverage=cv, purity=p, ploidy=2,
```

Copy number calling and SNV classification using targeted short read sequencing

```
verbose=FALSE)$power))

# Figure S7b in Carter et al.
plot(coverage, power[[1]], col=1, xlab="Sequence coverage",
     ylab="Detection power", ylim=c(0,1), type="l")

for (i in 2:length(power)) lines(coverage, power[[i]], col=i)
abline(h=0.8, lty=2, col="grey")
legend("bottomright", legend=paste("Purity", purity),
      fill=seq_along(purity))
```

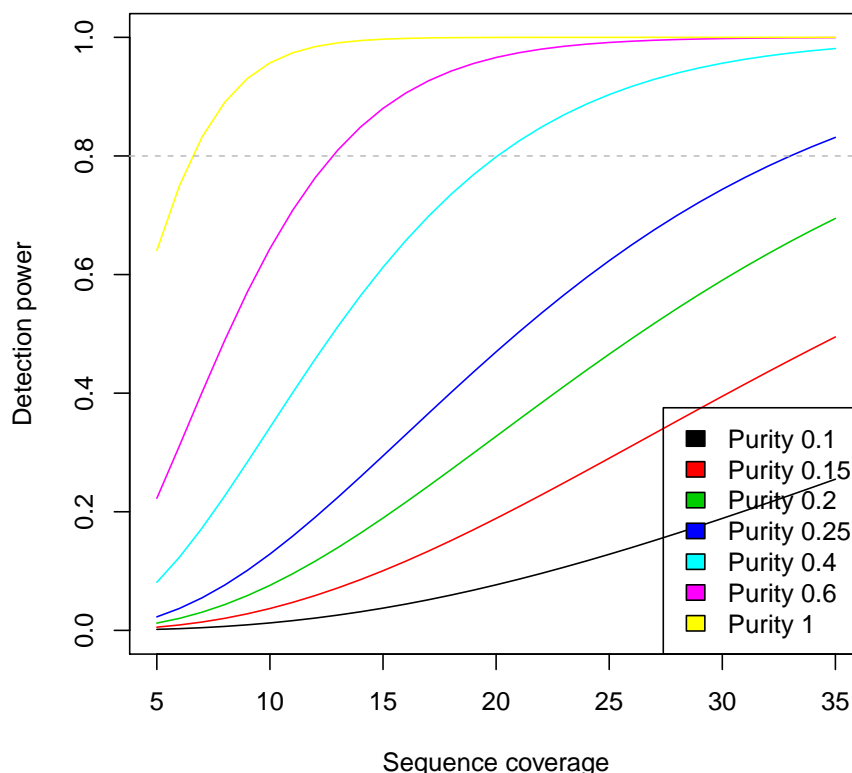


Figure 8: Power to detect mono-clonal somatic mutations

Reproduced from [2].

Figure 9 then shows the same plot for sub-clonal mutations present in 20% of all tumor cells:

```
coverage <- seq(5,350,1)
power <- lapply(purity, function(p) sapply(coverage, function(cv)
  calculatePowerDetectSomatic(coverage=cv, purity=p, ploidy=2,
    cell.fraction=0.2, verbose=FALSE)$power))
plot(coverage, power[[1]], col=1, xlab="Sequence coverage",
     ylab="Detection power", ylim=c(0,1), type="l")

for (i in 2:length(power)) lines(coverage, power[[i]], col=i)
abline(h=0.8, lty=2, col="grey")
legend("bottomright", legend=paste("Purity", purity),
```

Copy number calling and SNV classification using targeted short read sequencing

```
fill=seq_along(purity))
```

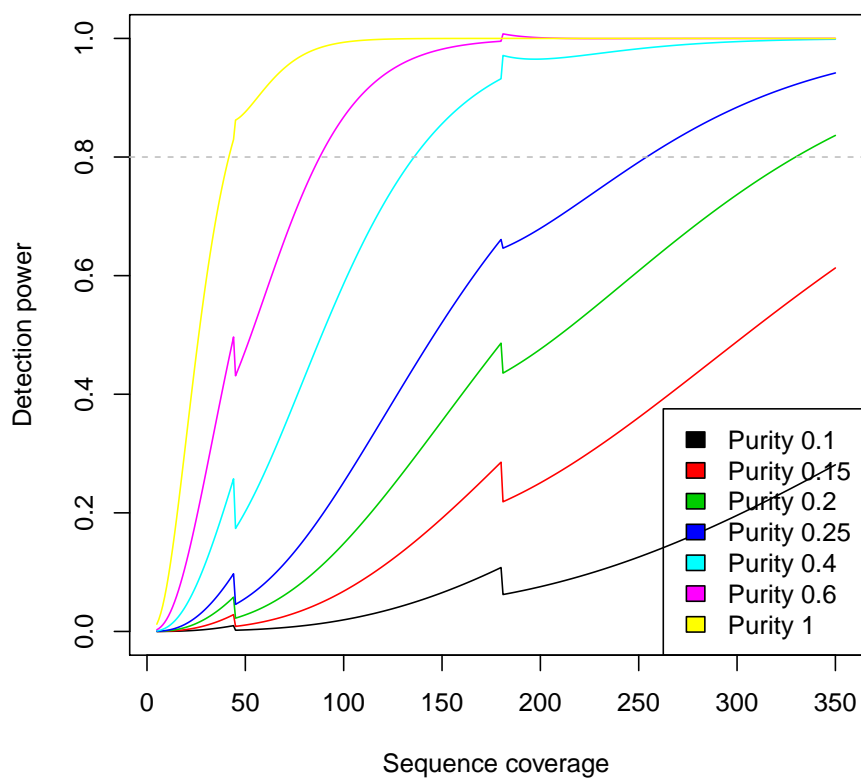


Figure 9: Power to detect sub-clonal somatic mutations present in 20% of all tumor cells
Reproduced from [2].

11 Limitations

PureCN currently assumes a completely diploid normal genome. For human samples, it tries to detect sex by calculating the coverage ratio of chromosomes X and Y and will then remove sex chromosomes in male samples⁸. For non-human samples, the user needs to manually remove all non-diploid chromosomes from the coverage data and specify `sex="diploid"` in the *PureCN* call.

⁸Loss of Y chromosome (LOY) can result in wrong female calls, especially in high purity samples or if LOY is in both tumor and contaminating normal cells. The software throws a warning in this case when germline SNPs are provided.

While *PureCN* supports and models sub-clonal somatic copy number alterations, it currently assumes that the majority of alterations are mono-clonal. For most clinical samples, this is reasonable, but very heterogeneous samples are likely not possible to call without manual curation. Poly-genomic tumors are often called as high ploidy or low purity. The former usually happens when sub-clonal losses are called as 2 copies and mono-clonal losses correctly as 1 copy. The latter when sub-clonal losses are called mono-clonal, which only happens when there are far more sub-clonal than mono-clonal losses. Please note however that unless purities are very high, algorithms that model poly-genomic tumors do not necessarily have a higher call rate, since they tend to overfit noisy samples or similarly confuse true high-ploidy with poly-genomic tumors. Due to the lack of signal, manual curation is also recommended in low purity samples or very quiet genomes.

12 Support

If you encounter bugs or have problems running *PureCN*, please post them at

- https://support.bioconductor.org/p/new/post/?tag_val=PureCN or
- <https://github.com/lima1/PureCN/issues>.

If *PureCN* throws user errors, then there is likely a problem with the input files. If the error message is not self-explanatory, feel free to seek help at the support site.

In your report, please add the outputs of the `runAbsoluteCN` call (with `verbose=TRUE`, or the `*.log` file in *PureCN.R*) and `sessionInfo()`. Please also check that your problem is not already covered in the following sections.

For general feedback such as suggestions for improvements, feature requests, complaints, etc. please do not hesitate to send us an email.

12.1 Checklist

- Used the correct interval files provided by the manufacturer of the capture kit and the genome version of the interval file matches the reference. Ideally used the baits file, not the targets file (in Agilent data, the baits files are called "covered" and the targets are "regions").
- For hybrid capture data, included off-target reads in the coverage calculation
- BAM files were generated following established best practices and tools finished successfully.
- Checked standard QC metrics such as AT/GC dropout and duplication rates.
- Tumor and normal data were obtained using the same capture kit and pipeline.

Copy number calling and SNV classification using targeted short read sequencing

- Coverage data of tumor and normal were GC-normalized.
- The VCF file contains germline variants (i.e. not only somatic calls).
- Maximized the number of high coverage heterozygous SNPs, for example by running *MuTect* with a 50-75bp interval padding (Section 9). The `runAbsoluteCN` output lists the percentage of targets with variants and this should be around 10-15%. Ultra-deep sequencing data can provide good SNP allelic fractions in the 100-200bp flanking regions.
- If a pool of normal samples is available, followed the steps in Section 4.2.
- Read the output of `runAbsoluteCN` with `verbose=TRUE`, fixed all warnings.
- If third-party segmentation tools are used, checked that normalized log-ratios are not biased, i.e. very similar compared to *PureCN* log-ratios (some pipelines already adjust for a default normal contamination).

12.2 FAQ

If the ploidy is frequently too high, please check:

- Does the log-ratio histogram (Figure 3) look noisy? If yes, then
 - Is the coverage sufficient? Tumor coverages below 80X can be difficult, especially in low purity samples. Normal coverages below 50X might result in high variance of log-ratios. See Section 4.1 for finding a good normal sample for log-ratio calculation.
 - Is the coverage data of both tumor and normal GC-normalized? If not, see `correctCoverageBias`.
 - Is the quality of both tumor and normal sufficient? A high AT or GC-dropout might result in high variance of log-ratios. Challenging FFPE samples also might need parameter tuning of the segmentation function. See `segmentationCBS`. A high expected tumor purity allows more aggressive segmentation parameters, such as `prune.hclust.h=0.2` or higher.
 - Was the correct target interval file used (genome version and capture kit, see Section 2.4)? If unsure, ask the help desk of your sequencing center.
 - Were the normal samples run with the same assay and pipeline?
 - Did you provide `runAbsoluteCN` all the recommended files as described in Section 5?
 - For whole-genome data, you will get better results using a specialized third-party segmentation method as described in section 10.1, since our default is optimized for targeted sequencing.
- Otherwise, if log-ratio peaks are clean as in Figure 3:
 - Was MuTect run without a matched normal? If yes, then make sure to provide either a pool of normal VCF or a SNP blacklist (if no pool of normal samples is available) as described in Sections 4.2 and 4.3.
 - A high fraction of sub-clonal copy-number alterations might also result in a low ranking of correct low ploidy solutions (see Section 11).

Copy number calling and SNV classification using targeted short read sequencing

If the ploidy is frequently too low:

- *PureCN* with default parameters is conservative in calling genome duplications.
- This should only affect low purity samples (< 35%), since in higher purity samples the duplication signal is usually strong enough to reliably detect it.
- In whole-exome data, it is usually safe to decrease the `max.homozygous.loss` default, since such large losses are rare.

Will *PureCN* work with my data?

- *PureCN* was designed for medium-sized (>2-3Mb) targeted panels. The more data, the better, best results are typically achieved in whole-exome data.
- The number of heterozygous SNPs is also important (>1000 per sample). Copy number baits enriched in SNPs are therefore very helpful (see Section 9).
- Some users got acceptable results with small (<1Mb) panels. Try to find a perfect off-target bin width (`average.off.target.width` in `preprocessIntervals`) and maximize the number of heterozygous SNPs by including as much padding as possible. Keep in mind that without tiling baits, you will only have a poor resolution to detect LOH.
- Coverages below 80X are difficult unless purities are high and coverages are even.
- *PureCN* also needs process-matched normal samples, again, the more the better.
- Samples with tumor purities below 15-20% usually cannot be analyzed with this algorithm and *PureCN* might return very wrong purity estimates. In high coverage samples with low duplication rates, this limit can be close to 10%.
- Whole-genome data is not officially supported and specialized tools will likely provide better results. Third-party segmentation tools designed for this data type would be again required.
- Amplicon sequencing data is also not officially supported. If the assay contains tiling probes (at least with 1Mb spacing) and uses a barcode protocol that reduces PCR bias of measured allelic fractions, then this method might produce acceptable results. Setting the `model` argument of `runAbsoluteCN` to `betabin` is recommended. Specialized segmentation tools might be again better than our default.

If you have trouble generating input data *PureCN* accepts, please check:

- For problems related to generating valid coverage data, either consult the *GATK* manual for the *DepthOfCoverage* or *CollectFragmentCounts* tools or Section 2.3 for the equivalent function in *PureCN*. If you use *DepthOfCoverage* and off-target intervals as generated by *IntervalFile.R* (See Quick vignette), make sure to run it with parameters `-omitDepthOutputAtEachBase` and `-interval_merging OVERLAPPING_ONLY`.
- Currently only VCF files generated by *MuTect 1* are officially supported and well tested. A minimal example *MuTect* call would be:

```
$JAVA7 -Xmx6g -jar $MUTECT \  
--analysis_type MuTect -R $REFERENCE \  
--dbSNP $DBSNP_VCF \  
--cosmic $COSMIC_VCF \  
-I:normal $BAM_NORMAL \  

```

Copy number calling and SNV classification using targeted short read sequencing

```
-I:tumor $BAM_TUMOR \
-o $OUT/${ID}_bwa_mutect_stats.txt \
-vcf $OUT/${ID}_bwa_mutect.vcf
```

The normal needs to be matched; without matched normal, only provide the tumor BAM file (do NOT provide a process-matched normal here). The default output file is the stats or call-stats file; this can be provided in addition to the required VCF file via `args.filterVcf` in `runAbsoluteCN`. If provided, it may help *PureCN* filter artifacts. This requires *MuTect* in version 1.1.7. This version is currently available at <https://software.broadinstitute.org/gatk/download/mutect> and requires Java 1.7 (it does not work with Java 1.8).

Note that this *MuTect* VCF will contain variants in off-target reads. By default, *PureCN* will remove variants outside the provided targets and their 50bp flanking regions. We highly recommend finding good values for each assay. A good cutoff will maximize the number of heterozygous SNPs and keep only an acceptable number of lower quality calls. This cutoff is set via `interval.padding` in `args.filterVcf`. See Section 9.

- Support for GATK4 and Mutect2 is still experimental. When matched normals are available, this will require version 4.0.3.0 and higher and specifying the `-genotype-germline-sites` flag.
- For VCFs generated by other callers, the required dbSNP annotation can be added for example with *bcftools*:

```
bcftools annotate --annotation $DBSNP_VCF \
--columns ID --output $OUT/${ID}_bwa_dbsnp.vcf.gz --output-type z \
$OUT/${ID}_bwa.vcf.gz
```

- To generate a mappability file with the *GEM* library:
 - Calculate mappability, set kmer size to length of mapped reads.

```
REFERENCE="hg38.fa"
PREF=`basename $REFERENCE .fa`
THREADS=4
KMER=100
gem-indexer -T ${THREADS} -c dna -i ${REFERENCE} -o ${PREF}_index
gem-mappability -T ${THREADS} -I ${PREF}_index.gem -l ${KMER} \
-o ${PREF}_${KMER} -m 2 -e 2
gem-2-wig -I ${PREF}_index.gem -i ${PREF}_${KMER}.mappability \
-o ${PREF}_${KMER}
```

- Convert to bigWig format, for example using the UCSC *wigToBigWig* tool:

```
cut -f1,2 ${REFERENCE}.fai > ${PREF}.sizes
wigToBigWig ${PREF}_${KMER}.wig ${PREF}.sizes ${PREF}_${KMER}.bw
```

- To generate the normal panel VCF (`normal.panel.vcf.file`) with *GATK3*:
 - Run *MuTect* on the normal with `-I:tumor $BAM_NORMAL` and optionally with the `-artifact_detection_mode` flag.
 - Remove the empty `none` sample from the VCF:

Copy number calling and SNV classification using targeted short read sequencing

```
$JAVA -Xmx6g -jar $GATK3 \
--analysis_type SelectVariants -R $REFERENCE \
--exclude_sample_expressions none \
-V $OUT/${ID}_bwa_mutect_artifact_detection_mode.vcf \
-o $OUT/${ID}_bwa_mutect_artifact_detection_mode_no_none.vcf
```

- Merge the VCFs:

```
CMD="java -Xmx12g -jar $GATK3 -T CombineVariants --minimumN 3 -R $REFERENCE"
for VCF in $OUT/*bwa_mutect_artifact_detection_mode_no_none.vcf;
do
    CMD="$CMD --variant $VCF"
done
CMD="$CMD -o $OUT/normals_merged_min3.vcf"
echo $CMD > $OUT/merge_normals_min3.sh
bgzip $OUT/normals_merged_min3.vcf
tabix $OUT/normals_merged_min3.vcf.gz
```

Questions related to pool of normals. As described in detail in Sections 4.1 and 4.2, a pool of normal samples is used in *PureCN* for coverage normalization (to adjust for target-specific capture biases) and for artifact filtering. A few recommendations:

- Matched normals are obviously perfect for identifying most alignment artifacts and mapping biases. While not critical, we still recommend generating a `normal.panel.vcf` file for *MuTect* and `setMappingBiasVcf` using the available normals.
- Without matched normals, we need process-matched normals for coverage normalization. We recommend more than 5 from 2-3 different library preparation and sequencing batches. The more, the better.
- These normals used for coverage normalization should be ideally sequenced to at least half the coverage as the tumor samples. This is completely different from matched normals for germline calling where 30-40X provide sufficient power to detect heterozygosity.
- For artifact removal, the more normals available, the more rare artifacts are removed. We recommend at least 10, 50 or more are ideal. The more artifacts are removed, the less likely *PureCN* output requires manual curation (Section 4.2).
- For position-specific mapping bias correction, the more normals are available, the more rare SNPs will have reliable mapping bias estimates. This requires again at least about 10 normals to be useful, 100 or more are ideal.
- With smaller pool of normals, we additionally recommend filtering SNPs from low quality regions (Section 4.3). Additionally, it is worth trying the beta-binomial function instead of the default in the `model` argument of `runAbsoluteCN`. This will incorporate uncertainty of observed variant allelic fractions in the variant fitting step.
- *Do I really need a pool of normals? I only have tumor samples.* Unfortunately, yes. If you used a commercial capture kit, you might be able to obtain control samples from the vendor or the public domain. This is not optimal, but usually better than nothing.

Copy number calling and SNV classification using targeted short read sequencing

- It is safe to include multiple normals from the same individual. Fewer common germline CNVs in `calculateIntervalWeights` and fewer SNPs in `calculateMappingBiasVcf` will be detected. But especially when the controls were sequenced in multiple batches, these replicates will still provide useful information for coverage normalization.

Questions related to manual curation. *PureCN*, like most other related tools, essentially finds the most simple explanation of the data. There are three major problems with this approach:

- First, hybrid capture data can be noisy and the algorithm must distinguish signal from noise; if the algorithm mistakes noise for signal, then this often results in wrong high ploidy calls (see Sections 4.2 and 4.3). If all steps in this vignette were followed, then *PureCN* should ignore common artifacts. Noisy samples thus often have outlier ploidy values and are often automatically flagged by *PureCN*. The correct solution is in most of these cases ranked second or third.
- The second problem is that signal can be sparse, i.e. when the tumor purity is very low or when there are only few somatic events. Manual curation is often easy in the latter case. For example when small losses are called as homozygous, but corresponding germline allele-frequencies are unbalanced (a complete loss would result in balanced germline allele frequencies, since only normal DNA is left). Future versions might improve calling in these cases by underweighting uninformative genomic regions.
- The third problem is that tumor evolution is fast and complex and very difficult to incorporate into general likelihood models. Sometimes multiple solutions explain the data equally well, but one solution is then often clearly more consistent with known biology, for example LOH in tumor suppressor genes such as *TP53*. A basic understanding of both the algorithm and the tumor biology of the particular cancer type are thus important for curation. Fortunately, in most cancer types, such ambiguity is rather rare. See also Section 11.

If all or most of the samples are flagged as:

Noisy segmentation: The default of 300 for `max.segments` is calibrated for high quality and high coverage whole-exome data. For whole-genome data or lower coverage data, this value needs to be re-calibrated. In case the copy number data looks indeed noisy, please see the first FAQ. Please be aware that *PureCN* will apply more aggressive segmentation parameters when the number of segments exceeds this cutoff. If the high segment count is real, this might confound downstream analyses.

High AT/GC dropout: If the data is GC-normalized, then there might be issues with either the target intervals or the provided GC content. Please double check that all files are correct and that all the coverage files are GC-normalized (Section 3).

Sex mismatch of coverage and VCF: If the panel contains baits for chromosome Y, then the interval file was probably generated without mappability file (Section 2.2). Cross-sample contamination (Section 10.5) can also cause sex mismatches.

References

- [1] Markus Riester, Angad P. Singh, A. Rose Brannon, Kun Yu, Catarina D. Campbell, Derek Y. Chiang, and Michael P. Morrissey. PureCN: copy number calling and SNV classification using targeted short read sequencing. *Source code for biology and medicine*, 11:13, Dec 2016.
- [2] Scott L. Carter, Kristian Cibulskis, Elena Helman, Aaron McKenna, Hui Shen, Travis Zack, Peter W. Laird, Robert C. Onofrio, Wendy Winckler, Barbara A. Weir, Rameen Beroukhi, David Pellman, Douglas A. Levine, Eric S. Lander, Matthew Meyerson, and Gad Getz. Absolute quantification of somatic DNA alterations in human cancer. *Nature biotechnology*, 30(5):413–421, May 2012.
- [3] Aaron McKenna, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew Kernysky, Kiran Garimella, David Altshuler, Stacey Gabriel, Mark Daly, and Mark A. DePristo. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome research*, 20(9):1297–1303, Sep 2010.
- [4] Kristian Cibulskis, Michael S. Lawrence, Scott L. Carter, Andrey Sivachenko, David Jaffe, Carrie Sougnez, Stacey Gabriel, Matthew Meyerson, Eric S. Lander, and Gad Getz. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature biotechnology*, 31(3):213–219, Mar 2013.
- [5] Eric Talevich, A. Hunter Shain, Thomas Botton, and Boris C. Bastian. CNVkit: Genome-wide copy number detection and visualization from targeted DNA sequencing. *PLoS computational biology*, 12(4):e1004873, Apr 2016.
- [6] Amnon Koren, Robert E. Handsaker, Nolan Kamitaki, Rosa Karlić, Sulagna Ghosh, Paz Polak, Kevin Eggan, and Steven A. McCarroll. Genetic variation in human dna replication timing. *Cell*, 159(5):1015–1026, Nov 2014.
- [7] Kortine Kleinheinz, Isabell Bludau, Daniel Huebschmann, Michael Heinold, Philip Kensche, Zuguang Gu, Cristina Lopez, Michael Hummel, Wolfram Klapper, Peter Moeller, Inga Vater, Rabea Wagener, Benedikt Brors, Reiner Siebert, Roland Eils, and Matthias Schlesner. ACEseq - allele specific copy number estimation from whole genome sequencing. *bioRxiv*, 2017. URL: <https://www.biorxiv.org/content/early/2017/10/29/210807>, doi:10.1101/210807.
- [8] Gavin Ha, Andrew Roth, Jaswinder Khattra, Julie Ho, Damian Yap, Leah M. Prentice, Nataliya Melnyk, Andrew McPherson, Ali Bashashati, Emma Laks, Justina Biele, Jiarui Ding, Alan Le, Jamie Rosner, Karey Shumansky, Marco A. Marra, C. Blake Gilks, David G. Huntsman, Jessica N. McAlpine, Samuel Aparicio, and Sohrab P. Shah. TITAN: inference of copy number architectures in clonal cell populations from tumor whole-genome sequence data. *Genome research*, 24(11):1881–1893, Nov 2014.
- [9] P. Andrew Futreal, Lachlan Coin, Mhairi Marshall, Thomas Down, Timothy Hubbard, Richard Wooster, Nazneen Rahman, and Michael R. Stratton. A census of human cancer genes. *Nature reviews. Cancer*, 4(3):177–183, Mar 2004.
- [10] E. S. Venkatraman and Adam B. Olshen. A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics (Oxford, England)*, 23(6):657–663, Mar 2007.

- [11] Simon A. Forbes, David Beare, Prasad Gunasekaran, Kenric Leung, Nidhi Bindal, Harry Boutselakis, Minjie Ding, Sally Bamford, Charlotte Cole, Sari Ward, Chai Yin Kok, Mingming Jia, Tisham De, Jon W. Teague, Michael R. Stratton, Ultan McDermott, and Peter J. Campbell. COSMIC: exploring the world's knowledge of somatic mutations in human cancer. *Nucleic acids research*, 43(Database issue):D805–D811, Jan 2015.

Session Info

- R version 3.5.2 (2018-12-20), x86_64-w64-mingw32
- Locale: LC_COLLATE=C, LC_CTYPE=English_United States.1252, LC_MONETARY=English_United States.1252, LC_NUMERIC=C, LC_TIME=English_United States.1252
- Running under: Windows Server 2012 R2 x64 (build 9600)
- Matrix products: default
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, stats4, utils
- Other packages: Biobase 2.42.0, BiocGenerics 0.28.0, BiocParallel 1.16.5, Biostrings 2.50.2, DNACopy 1.56.0, DelayedArray 0.8.0, GenomeInfoDb 1.18.1, GenomicRanges 1.34.0, IRanges 2.16.0, PureCN 1.12.2, Rsamtools 1.34.0, S4Vectors 0.20.1, SummarizedExperiment 1.12.0, VariantAnnotation 1.28.5, XVector 0.22.0, matrixStats 0.54.0
- Loaded via a namespace (and not attached): AnnotationDbi 1.44.0, BSgenome 1.50.0, BiocManager 1.30.4, BiocStyle 2.10.0, DBI 1.0.0, GenomeInfoDbData 1.2.0, GenomicAlignments 1.18.1, GenomicFeatures 1.34.1, Matrix 1.2-15, R6 2.3.0, RColorBrewer 1.1-2, RCurl 1.95-4.11, RSQLite 2.1.1, Rcpp 1.0.0, Rhdf5lib 1.4.2, VGAM 1.0-6, XML 3.98-1.16, assertthat 0.2.0, bindr 0.1.1, bindrcpp 0.2.2, biomaRt 2.38.0, bit 1.1-14, bit64 0.9-7, bitops 1.0-6, blob 1.1.1, colorspace 1.3-2, compiler 3.5.2, crayon 1.3.4, data.table 1.11.8, digest 0.6.18, dplyr 0.7.8, evaluate 0.12, formatR 1.5, futile.logger 1.4.3, futile.options 1.0.1, ggplot2 3.1.0, glue 1.3.0, grid 3.5.2, gridExtra 2.3, gtable 0.2.0, highr 0.7, hms 0.4.2, htmltools 0.3.6, http 1.4.0, knitr 1.21, labeling 0.3, lambda.r 1.2.3, lattice 0.20-38, lazyeval 0.2.1, magrittr 1.5, memoise 1.1.0, munsell 0.5.0, pillar 1.3.1, pkgconfig 2.0.2, plyr 1.8.4, prettyunits 1.0.2, progress 1.2.0, purrr 0.2.5, rhdf5 2.26.2, rlang 0.3.0.1, rmarkdown 1.11, rtracklayer 1.42.1, scales 1.0.0, splines 3.5.2, stringi 1.2.4, stringr 1.3.1, tibble 2.0.0, tidyselect 0.2.5, tools 3.5.2, xfun 0.4, yaml 2.2.0, zlibbioc 1.28.0