

PureCN - Quick Start

This tutorial provides a quick overview of the command line tools shipping with *PureCN*. For the R package and more detailed information, see the main vignette.

Update from previous stable versions

PureCN is backward compatible with input generated by the previous stable version. For optimal results, we recommend starting from scratch following the tutorial below. *PureCN* 1.8 introduced off-target read support and minor improvements that are not available with previously generated data.

Prepare environment and files

- Start R and enter the following to get the path to the command line scripts:

```
system.file("extdata", package="PureCN")
## [1] "/tmp/RtmpgFBWVP/Rinst3fd625bdc24a/PureCN/extdata"
```

- Exit R and store this path in an environment variable, for example in BASH:

```
$ export PURECN="/path/to/PureCN/extdata"
$ Rscript $PURECN/PureCN.R --help
Usage: /path/to/PureCN/inst/extdata/PureCN.R [options] ...
```

- Generate an interval file from a BED file containing target coordinates:

```
$ Rscript $PURECN/IntervalFile.R --infile baits_hg19.bed \
  --fasta hg19.fa --outfile baits_hg19_gcgene.txt \
  --offtarget --genome hg19 \
  --export baits_optimized_hg19.bed \
  --mappability wgEncodeCrgMapabilityAlign100mer.bigWig
```

Internally, this script uses *rtracklayer* to parse the `infile`. Make sure that the file format matches the file extension. See the *rtracklayer* documentation for problems loading the file.

The `-offtarget` flag will include off-target reads. Including them is recommended except for Amplicon data.

The `-genome` version is needed to annotate exons with gene symbols.

The `-export` argument is optional. If provided, this script will store the modified intervals as BED file for example (again every *rtracklayer* format is supported). This is useful when the coverages are calculated with third-party tools like GATK.

The `-mappability` argument should provide a *rtracklayer* parsable file with a mappability score in the first meta data column. If provided, off-target regions will be restricted to regions specified in this file. On-target regions with low mappability will be excluded. For a test run, you can skip this argument or simply download the file from the UCSC website (see the FAQ section of the main vignette for instruction how to generate such a file).

Create VCF files

PureCN does not ship with a variant caller. Use a third-party tool to generate a VCF for each sample. A few recommendations:

- Use *MuTect* 1.1.7 if possible
- Support for *MuTect* 2 and *FreeBayes* is available for tumor-only VCFs, but currently poorly tested and only very limited artifact filtering will be performed for these callers.
- Since germline SNPs are needed to infer allele-specific copy numbers, the provided VCF needs to contain both somatic and germline variants.
- Run the variant caller with a 50 base pair interval padding to increase the number of heterozygous SNPs

Run PureCN with internal segmentation

The following describes *PureCN* runs with internal copy number normalization and segmentation.

Coverage

For each sample, tumor and normal:

```
# Calculate and GC-normalize coverage from a BAM file
$ Rscript $PURECN/Coverage.R --outdir $OUT/$SAMPLEID \
    --bam ${SAMPLEID}.bam \
    --gogene baits_hg19_gogene.txt

# GC-normalize coverage from a GATK DepthOfCoverage file
Rscript $PURECN/Coverage.R --outdir $OUT/$SAMPLEID \
    --coverage ${SAMPLEID}.coverage.sample_interval_summary \
    --gogene baits_hg19_gogene.txt
```

Similar to GATK, this script also takes a text file containing a list of BAM or coverage file names (one per line). The file extension must be .list:

```
# Calculate and GC-normalize coverage from a list of BAM files
$ Rscript $PURECN/Coverage.R --outdir $OUT/$SAMPLEID \
    --bam normals.list \
    --gogene baits_hg19_gogene.txt \
    --cpu 4
```

NormalDB

To build a normal database, copy the paths to all GC-normalized normal coverage files in a single text file, line-by-line:

```
ls -a normal*loess.txt | cat > example_normal.list

# From already GC-normalized files
```

```
$ Rscript $PURECN/NormalDB.R --outdir $OUT \
  --coveragefiles example_normal.list \
  --genome hg19
```

- The resulting `normalDB_hg19.rds` file contains absolute paths to the coverage files. It is thus necessary to re-run this command when the coverage files are moved.
- Consider generating different databases when differences are significant, e.g. for samples with different read lengths or insert size distributions
- In particular, do not mix normal data obtained with different capture kits (e.g. Agilent SureSelect v4 and v6)

PureCN

Now that the assay-specific files are created and all coverages calculated, we run PureCN.R to normalize, segment and determine purity and ploidy:

```
cd $OUT/${SAMPLEID}

# Without a matched normal (minimal test run)
$ Rscript $PURECN/PureCN.R --out . --tumor ${SAMPLEID}_coverage_loess.txt \
  --normaldb ../normalDB_hg19.rds \
  --sampleid ${SAMPLEID} \
  --vcf ${SAMPLEID}_mutect.vcf \
  --pool 5 \
  --genome hg19 --gchgce baits_hg19_gchgce.txt

# Production pipeline run
$ Rscript $PURECN/PureCN.R --out . --tumor ${SAMPLEID}_coverage_loess.txt \
  --normaldb ../normalDB_hg19.rds \
  --normal_panel $NORMAL_PANEL \
  --sampleid ${SAMPLEID} \
  --vcf ${SAMPLEID}_mutect.vcf \
  --statsfile ${SAMPLEID}_mutect_stats.txt \
  --pool 5 \
  --genome hg19 --gchgce baits_hg19_gchgce.txt
  --snpblacklist hg19_simpleRepeats.bed \
  --targetweightfile ../target_weights_hg19.txt
  --force --postoptimize

# With a matched normal (test run)
$ Rscript $PURECN/PureCN.R --out . --tumor ${SAMPLEID}_coverage_loess.txt \
  --normal ${SAMPLEID_NORMAL}_coverage_loess.txt \
  --sampleid ${SAMPLEID} \
  --vcf ${SAMPLEID}_mutect.vcf \
  --genome hg19 \
  --normaldb ../normalDB_hg19.rds \
  --gchgce baits_hg19_gchgce.txt

# Recreate output after manual curation of ${SAMPLEID}.csv
$ Rscript $PURECN/PureCN.R --rds ${SAMPLEID}.rds
```

A few important recommendations:

- Even if matched normals are available, it is often safer to use the normal database for coverage normalization
- Providing the normal database in addition to a matched normal is optional, but helps ignoring low quality regions in the segmentation
- The `-pool` flag specifies how many normal samples should be used to calculate the tumor vs. normal coverage. In a large, homogeneous database, this should be set to 10. In a heterogeneous database (for example different insert sizes, coverages, library preparation protocols, etc.), this value is set to an appropriate group size between 2 and 10.
- The normal panel VCF file is useful for mapping bias correction and especially recommended without matched normals. See the FAQ of the main vignette how to generate this file. It is not essential for test runs.
- The *MuTect* 1.1.7 stats file (the main output file besides the VCF) should be provided for better artifact filtering. If the VCF was generated by a pipeline that performs good artifact filtering, this file is not needed.
- The `-postoptimize` flag defines that purity should be optimized using both variant allelic fractions and copy number instead of copy number only. This results in a significant runtime increase for whole-exome data.
- If `-out` is a directory, it will use the sample id as file prefix for all output files. Otherwise *PureCN* will use `-out` as prefix.

Run PureCN with third-party segmentation

If you already have a segmentation from third-party tools (for example CNVkit, EXCAVATOR2). For a test run:

```
Rscript $PURECN/PureCN.R --out $OUT/$SAMPLEID \
  --sampleid $SAMPLEID \
  --segfile $OUT/$SAMPLEID/${SAMPLEID}.cnvkit.seg \
  --vcf ${SAMPLEID}_mutect.vcf \
  --genome hg19 --gchgane baits_hg19_gchgane.txt
```

See the main vignette for more details and file formats.

For a production pipeline run we provide a bit more information about the assay and genome:

```
# Recommended CNVkit example
Rscript $PURECN/PureCN.R --out $OUT/$SAMPLEID \
  --sampleid $SAMPLEID \
  --tumor $OUT/$SAMPLEID/${SAMPLEID}.cnr \
  --segfile $OUT/$SAMPLEID/${SAMPLEID}_cnvkit.seg \
  --normal_panel $NORMAL_PANEL \
  --vcf ${SAMPLEID}_mutect.vcf \
  --statsfile ${SAMPLEID}_mutect_stats.txt \
  --snpblacklist hg19_simpleRepeats.bed \
  --genome hg19 \
  --funsegmentation none \
  --force --postoptimize
```

- The `-funsegmentation` argument controls if the data should be re-segmented using germline BAFs (default). Set this value to `none` if the provided segmentation should be used as is.

Dx

Dx.R extracts copy number and mutation metrics from PureCN.R output.

```
# Provide a BED file with callable regions, for examples obtained by
# GATK CallableLoci. Useful to calculate mutations per megabase and
# to exclude low quality regions.
grep CALLABLE ${SAMPLEID}_callable_status.bed > \
${SAMPLEID}_callable_status_filtered.bed

Rscript $PureCN/Dx.R --out . --rds ${SAMPLEID}.rds \
--callable ${SAMPLEID}_callable_status_filtered.bed
```

Reference

Table 1: IntervalFile

Argument name	Corresponding PureCN argument	PureCN function
<code>-fasta</code>	<code>reference.file</code>	<code>calculateGCContentByInterval</code>
<code>-infile</code>	<code>interval.file</code>	<code>calculateGCContentByInterval</code>
<code>-offtarget</code>	<code>off.target</code>	<code>calculateGCContentByInterval</code>
<code>-targetwidth</code>	<code>average.target.width</code>	<code>calculateGCContentByInterval</code>
<code>-offtargetwidth</code>	<code>average.off.target.width</code>	<code>calculateGCContentByInterval</code>
<code>-offtargetseqlevels</code>	<code>off.target.seqlevels</code>	<code>calculateGCContentByInterval</code>
<code>-mappability</code>	<code>mappability</code>	<code>calculateGCContentByInterval</code>
<code>-genome</code>	<code>txdb, org</code>	<code>annotateTargets</code>
<code>-outfile</code>		
<code>-export</code>		<code>rtracklayer::export</code>
<code>-version -v</code>		
<code>-force -f</code>		
<code>-help -h</code>		

Table 2: Coverage

Argument name	Corresponding PureCN argument	PureCN function
<code>-bam</code>	<code>bam.file</code>	<code>calculateBamCoverageByInterval</code>
<code>-bai</code>	<code>index.file</code>	<code>calculateBamCoverageByInterval</code>
<code>-coverage</code>	<code>coverage.file</code>	<code>correctCoverageBias</code>
<code>-gchgene</code>	<code>gc.gene.file</code>	<code>correctCoverageBias</code>
<code>-method</code>	<code>method</code>	<code>correctCoverageBias</code>
<code>-keepduplicates</code>	<code>keep.duplicates</code>	<code>calculateBamCoverageByInterval</code>
<code>-outdir</code>		
<code>-cpu</code>		Number of CPUs to use
<code>-seed</code>		
<code>-version -v</code>		
<code>-force -f</code>		
<code>-help -h</code>		

Table 3: NormalDB

Argument name	Corresponding PureCN argument	PureCN function
<code>-coveragefiles</code>	<code>normal.coverage.files</code>	<code>createNormalDatabase</code>
<code>-maxmeancoverage</code>	<code>max.mean.coverage</code>	<code>createNormalDatabase</code>
<code>-assay -a</code>	Optional assay name	Used in output file names.
<code>-genome -g</code>	Optional genome version	Used in output file names.
<code>-outdir -o</code>		
<code>-version -v</code>		
<code>-force -f</code>		
<code>-help -h</code>		

Table 4: PureCN

Argument name	Corresponding PureCN argument	PureCN function
<code>-sampleid -i</code>	<code>sampleid</code>	<code>runAbsoluteCN</code>
<code>-normal</code>	<code>normal.coverage.file</code>	<code>runAbsoluteCN</code>
<code>-tumor</code>	<code>tumor.coverage.file</code>	<code>runAbsoluteCN</code>
<code>-vcf</code>	<code>vcf.file</code>	<code>runAbsoluteCN</code>
<code>-rds</code>	<code>file.rds</code>	<code>readCurationFile</code>
<code>-normal_panel</code>	<code>normal.panel.vcf.file</code>	<code>setMappingBiasVcf</code>
<code>-normaldb</code>	<code>normalDB</code> (serialized with <code>saveRDS</code>)	<code>findBestNormal, filterTargets</code>
<code>-segfile</code>	<code>seg.file</code>	<code>runAbsoluteCN</code>
<code>-sex</code>	<code>sex</code>	<code>runAbsoluteCN</code>
<code>-pool</code>	<code>pool, num.normals</code>	<code>findBestNormal</code>
<code>-genome</code>	<code>genome</code>	<code>runAbsoluteCN</code>
<code>-gchgene</code>	<code>gc.gene.file</code>	<code>runAbsoluteCN</code>
<code>-statsfile</code>	<code>stats.file</code>	<code>filterVcfMuTect</code>
<code>-minaf</code>	<code>af.range</code>	<code>filterVcfBasic</code>
<code>-snpbblacklist</code>	<code>snp.blacklist</code>	<code>filterVcfBasic</code>
<code>-error</code>	<code>error</code>	<code>runAbsoluteCN</code>
<code>-funsegmentation</code>	<code>fun.segmentation</code>	<code>runAbsoluteCN</code>
<code>-alpha</code>	<code>alpha</code>	<code>segmentationCBS</code>
<code>-maxsegments</code>	<code>max.segments</code>	<code>runAbsoluteCN</code>
<code>-targetweightfile</code>	<code>target.weight.file</code>	<code>segmentationCBS</code>
<code>-minpurity</code>	<code>test.purity</code>	<code>runAbsoluteCN</code>
<code>-maxpurity</code>	<code>test.purity</code>	<code>runAbsoluteCN</code>
<code>-minploidy</code>	<code>min.ploidy</code>	<code>runAbsoluteCN</code>
<code>-maxploidy</code>	<code>max.ploidy</code>	<code>runAbsoluteCN</code>
<code>-postoptimize</code>	<code>post.optimize</code>	<code>runAbsoluteCN</code>
<code>-modelhomozygous</code>	<code>model.homozygous</code>	<code>runAbsoluteCN</code>
<code>-model</code>	<code>model</code>	<code>runAbsoluteCN</code>
<code>-logratioalibration</code>	<code>log.ratio.calibration</code>	<code>runAbsoluteCN</code>
<code>-outvcf</code>	<code>return.vcf</code>	<code>predictSomatic</code>
<code>-out -o</code>		
<code>-seed</code>		
<code>-version -v</code>		
<code>-force -f</code>		
<code>-help -h</code>		

Table 5: Dx

Argument name	Corresponding PureCN argument	PureCN function
<code>-rds</code>	<code>file.rds</code>	<code>readCurationFile</code>
<code>-callable</code>	<code>callable</code>	<code>callMutationBurden</code>
<code>-exclude</code>	<code>exclude</code>	<code>callMutationBurden</code>
<code>-out</code>		
<code>-version -v</code>		
<code>-force -f</code>		
<code>-help -h</code>		

Session Info

- R version 3.4.3 (2017-11-30), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_US.UTF-8, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Running under: Ubuntu 16.04.3 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.6-bioc/R/lib/libRblas.so
- LAPACK: /home/biocbuild/bbs-3.6-bioc/R/lib/libRlapack.so
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, stats4, utils
- Other packages: Biobase 2.38.0, BiocGenerics 0.24.0, Biostrings 2.46.0, DNAcopy 1.52.0, DelayedArray 0.4.1, GenomeInfoDb 1.14.0, GenomicRanges 1.30.0, IRanges 2.12.0, PureCN 1.8.1, Rsamtools 1.30.0, S4Vectors 0.16.0, SummarizedExperiment 1.8.0, VariantAnnotation 1.24.2, XVector 0.18.0, matrixStats 0.52.2
- Loaded via a namespace (and not attached): AnnotationDbi 1.40.0, BSgenome 1.46.0, BiocParallel 1.12.0, BiocStyle 2.6.1, DBI 0.7, GenomeInfoDbData 0.99.1, GenomicAlignments 1.14.1, GenomicFeatures 1.30.0, Matrix 1.2-12, R6 2.2.2, RColorBrewer 1.1-2, RCurl 1.95-4.8, RMySQL 0.10.13, RSQLite 2.0, Rcpp 0.12.14, VGAM 1.0-4, XML 3.98-1.9, assertthat 0.2.0, backports 1.1.1, biomaRt 2.34.0, bit 1.1-12, bit64 0.9-7, bitops 1.0-6, blob 1.1.0, colorspace 1.3-2, compiler 3.4.3, data.table 1.10.4-3, digest 0.6.12, edgeR 3.20.1, evaluate 0.10.1, futile.logger 1.4.3, futile.options 1.0.0, ggplot2 2.2.1, grid 3.4.3, gtable 0.2.0, highr 0.6, htmltools 0.3.6, knitr 1.17, labeling 0.3, lambda.r 1.2, lattice 0.20-35, lazyeval 0.2.1, limma 3.34.3, locfit 1.5-9.1, magrittr 1.5, memoise 1.1.0, munsell 0.4.3, plyr 1.8.4, prettyunits 1.0.2, progress 1.1.2, rlang 0.1.4, rmarkdown 1.8, rprojroot 1.2, rtracklayer 1.38.2, scales 0.5.0, splines 3.4.3, stringi 1.1.6, stringr 1.2.0, tibble 1.3.4, tools 3.4.3, yaml 2.1.15, zlibbioc 1.24.0