# HANDYCNV: AN R PACKAGE FOR STANDARDIZED SUMMARY, ANNOTATION, COMPARISON, AND VISUALIZATION OF CNV AND CNVR

## J. Zhou[1,2], L. Liu[1,2], T.J. Lopdell [3], D. Garrick[2*] and Y. Shi[1*]

[1]School of Agriculture, Ningxia University, Yinchuan, China
[2]AL Rae Centre of Genetics and Breeding, Massey University, Hamilton, New Zealand
[3] Research and Development, Livestock Improvement Corporation, Ruakura Road, Hamilton, New Zealand

## SUMMARY

There is a need for a pipeline to provide standard, reproducible and timesaving post-analysis of CNV (copy number variants) from SNP (Single nucleotide polymorphisms) chip genotyping. We present a package built with a dozen functions that can convert the coordinates of SNP map files, compare the positions of SNPs between the given maps, summarize the CNVs, call CNVRs (Copy number variation regions), provide gene annotation, compare CNV, CNVR and the annotated gene lists, and visualize CNVs at both individual and population level.

## INTRODUCTION

Copy number variants are a type of structural variation of a DNA fragment which comprise the deletion or duplication type depending upon how many copies an individual has compared with the two copies in the diploid reference genome. These structure variants could change the structure or dosage of genes that might further affect the phenotypes. Studies on CNVs have become common in livestock research in recent years. The fluorescent signal intensity of SNPs chip provides the general source to detect CNV, and thanks to the wide application of genome-wide association studies and genomic selection in animal breeding, there is now a lot of SNP data suitable for analysis of structural variants. Software such as PennCNV (Wang *et al*. 2007), CNVPartition (Illumina) and SVS Golden Helix (Bozeman) are designed to detect CNV from SNP data, but each method has its own advantages and shortcomings, so it is recommended to use more than one method to infer CNVs (Winchester *et al*. 2009), therefore comparison of multiple CNV results is a normal task in characterizing structural variation.

When doing CNV analysis we are curious about all the information related to any CNV region, not only at the individual level but also at the population level. For instance, we want to know how many individuals have a CNV in a common region? What kind of type of CNVs are there? Are these individuals from the same farm or are they progenies of the same sire? Are there any genes in the CNV region? What are the gene frequencies? How about the signal intensities, call rate, minor allele frequency and linkage disequilibrium conditions of these CNVs? Besides, the common post-analysis of CNV studies includes provision of summary CNVs, generation of CNVR, comparison of CNVs from different software, finding consensus CNVR by comparing results to the gold standard CNV database, gene annotation in the CNV region and CNV-based regression analysis. To accomplish all these tasks various tools are typically required. Therefore, we integrated these functions into a package to make post-CNV analysis easy and reproducible. The use of some developed R package like the Tidyverse family (Wickham *et al*. 2019) made our function development much easier. We believe this package will be convenient to others who are doing similar work. The source code can be found in the Github repository (https://github.com/JH-Zhou/HandyCNV).

## MATERIALS AND METHODS

The pipeline and results of this package are shown in Figure 1. The Demo Data are the CNV results from the GeneSeek GGP Bovine 150k BeadChip detected by PennCNV (Wang *et al*. 2007)

and CNVPartition (Illumina). All the input files and demo code can be found in Github (https://github.com/JH-Zhou/HandyCNV). Run through all functions in HandyCNV need prepare CNV Results, SNP maps, Reference Gene List, Pedigree, Plink files (Bim, Bed and Fam) and SNP Signal Intensity in total. All the input files have a fixed format, and the file requirements depends on which functions the users are using. Here we only introduce the input data format and some noticeable methods we used in some functions, therefore, we will not cover the data structures or interpret how to use the results in this article, more details can be found by browsing our Demo Data.
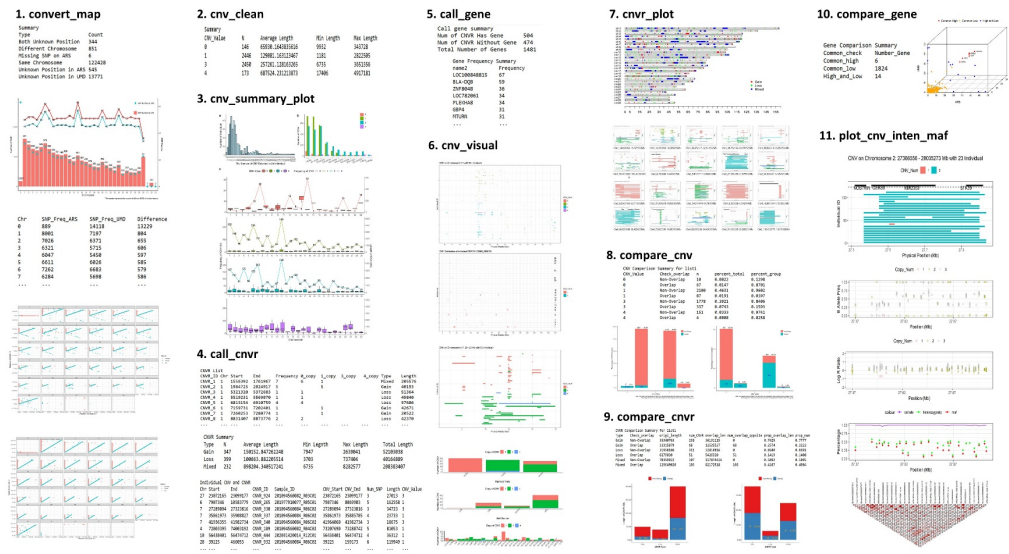


**Figure 1. Pipeline and results of HandyCNV for the post-analysis of CNV**

The first function is convert_map which is used to convert map files from the original to an objective map file provided by the user. There are differences between genome assemblies, for example, in which some SNP might locate on a different chromosome or on the same chromosome but in a different order between different assemblies. Most Bovine SNP chips have been using the UMD3.1 (Shamimuzzaman *et al*. 2019) as the default reference genome assembly, but with the release of new reference genome ARS-UCD1.2 with high continuity, accuracy, and completeness (Rosen *et al*. 2020),it may be of interest to convert the coordinates to the latest assembly to help further research. Four columns are required with no header in the input map files whose columns are Chromosome, SNP ID, Morgan Position (UMD) or Physical Position (ARS) and Physical Position (unit: bp) (Table 1, Table 2).

**Table 1. Original map format (UMD 3.1)**

| 14 | ARS-BFGL-BAC-10172 | 0 | 6371334 |
|----|--------------------|-------|----------|
| 14 | ARS-BFGL-BAC-1020 | 0 | 7928189 |
| 14 | ARS-BFGL-BAC-10245 | 28.23 | 31819743 |
| 14 | ARS-BFGL-BAC-10345 | 0 | 6133529 |

**Table 2. Objective map format (ARS)**

| 14 | ARS-BFGL-BAC-10172 | 5.34266 | 5342658 |
|----|--------------------|---------|---------|
| 14 | ARS-BFGL-BAC-1020  | 6.88966 | 6889656 |
| 14 | ARS-BFGL-BAC-10245 | 30.1241 | 30124134 |
| 14 | ARS-BFGL-BAC-10345 | 5.10573 | 5105727 |

The cnv_clean function is designed to convert the CNV results to a standard format, the output clean CNV file is used as input data in many of the other functions. It supports PennCNV and CNVPartition default output results, the length of CNVs are calculated as one plus the end position minus the start position. The CNV results from other software can be prepared as the template format to use in the remaining functions (Table 3). The function cnv_summary_plot will generate several plots to show the number, length group, type, and frequency details of CNVs on individuals and on chromosomes.

**Table 3 Template Format of Clean CNV**

| Sample_ID | Chr | Start | End | CNV_Value | Length |
|-----------|-----|-------|-----|-----------|--------|
| 201094560060_R02C01 | 11 | 106224443 | 106359588 | 4 | 135146 |
| 201094560060_R02C01 | 12 | 58073538 | 58417437 | 1 | 343900 |
| 201094560060_R02C01 | 19 | 27576066 | 27643677 | 4 | 67612 |
| 201094560060_R02C02 | 1 | 88638760 | 88904687 | 3 | 265928 |

The call_cnvr function will merge the CNVs which have at least one bp overlapping length to a CNVR. The results are the non-redundant CNVRs,but this method could cause misleading information while reporting the genes and comparing the overlapping length on CNVR. This is because it may appear all CNVs in a CNVR are the same length but in reality there are often lots of short disparate CNVs. To solve this problem, combine call_gene and cnv_visual function will plot all genes located on CNVs of every individual in a CNVR. The call_gene function needs the user to provide the reference genes which can be downloaded from the UCSC website (http://hgdownload.soe.ucsc.edu/downloads.html).

The compare_cnv and compare_cnvr functions with the similar strategies, when the results have the same version coordinates they will compare directly, but when the coordinates are from different versions, it will convert the position for each file at first then make comparison between the coordinates of the latest version. The overlapped region between two interval results may be slightly different, when reporting and plotting the number and length of overlapping regions correspond to each input files, respectively.

**RESULTS AND DISCUSSION**

**When do you need to convert the coordinates of SNP or CNV?** The first scenario is when a new reference genome is released. Take the Bovine reference genome as example, the lasted version (ARS-UCD1.2) has higher coverage and accuracy of its genome assembly than the previous commonly used UMD3.1, so it may help to improve the accuracy of SNP-based CNVs detection by using the latest reference genome. The second scenario is to make comparison between results from different reference genomes. There are lots of studies that have reported CNVs using previous reference genomes, and we may want to compare their results to our assembly.

**Why do we need to visualize CNVR?** CNV is of interest relative to the comparison of individuals but the CNVR are mostly of interest at the population level. The common method to generate CNVR is to merge all overlapping CNVs from every individual into a common region, then make gene annotation and comparison on CNVRs on the population level.

The main shortcoming of SNP-based CNV detection is that it cannot report the exact start or end position because of the limited marker density, so when we merge these CNV intervals to a common CNVR the actual situation is that not all the CNVs with the same break points as the CNVR, therefore, not all the genes within a CNVR has the same frequency with the CNVs (Figure 1. 11). Sometimes we might find an interested candidate gene within a high frequency CNVR but if only a few individuals have CNV of that gene, the better way to avoid this mistake is to report the gene frequency by counting how many CNVs with this gene in a CNVR, but this not enough, because of some genes may have CNV in just a partial fragment rather than the entire gene, in this circumstance plotting all CNVs and annotated genes in a CNVR by the start and end position can make it much clearer to understand what is happening. We are often curious about all the information in a CNVR in a population, such as the relationship between SNPs in that region, so visualizing a CNVR by plotting all related information in one figure is a good solution.

**What are the limitations of this study?** First, we have only used it in bovine studies, so some functions may need to be revised to be used in other species. Second, the linkage disequilibrium (LD) plots are based on the Gaston package which was drawing the base plot only, for some CNVRs with fewer number of SNPs the plot size was not well controlled while merging it to other plots, and this could lead some CNVR plots to be unsuitable without further modifications. Third, plots of CNVs on the population level are suitable for small populations but could be too busy for large populations. Fourth, the functions for regression analysis between CNVs or CNVRs and phenotype are still being developed.

## CONCLUSIONS

Here we present an R package called HandyCNV in the initial version which includes several functions for tasks such as converting SNP maps, generating CNVR, genome annotation, comparing and visualizing of CNV and CNVR and reporting summary results on each step. This tool provides a standard, reproducible and timesaving post-analysis of copy number variants.

## ACKNOWLEDGEMENTS

## REFERENCES

Wang K., Li M., Hadley D., et al. Penn C.N.V. (2007) *Genome Res*. **17**(11): 1665.

Illumina. GenomeStudio. https://www.illumina.com/techniques/microarrays/array-data-analysis-experimental-design/genomestudio.html

Bozeman, MT: Golden Helix I. SNP & Variation Suite ™ (Version 8.x). http://www.goldenhelix.com

Winchester L., Yau C., and Ragoussis J. (2009) *Briefings Funct Genomics Proteomics*. **8**(5): 353.

Wickham H, Averick M, Bryan J, et al. (2019) *J Open Source Softw*.

Shamimuzzaman M., Le Tourneau J.J., Unni D.R., et al. (2019) *Nucleic Acids Res*.

Rosen BD, Bickhart DM, Schnabel RD, et al. (2020) *Gigascience*.