

ggoncplot: an R package for interactive visualisation of somatic mutation data from cancer patient cohorts

Sam El-Kamand¹, Julian M. W. Quinn¹, and Mark J. Cowley^{1,2}

¹ Children's Cancer Institute, Australia ² School of Clinical Medicine, UNSW Medicine & Health, Australia ¶ Corresponding author

DOI: 10.xxxxxx/draft

Software

- Review
- Repository
- Archive

Editor: Open Journals

Reviewers:

- @openjournals

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License (CC BY 4.0).

Summary

The ggoncplot R package generates interactive oncoplots to visualize mutational patterns across patient cancer cohorts (Figure 1). Oncoplots, also called oncoprints, reveal patterns of gene co-mutation and include marginal plots that indicate co-occurrence of gene mutations with tumour and clinical features. It is useful to relate gene mutation patterns seen in an oncoplot to patterns in other plot types, including gene expression t-SNE plots or methylation UMAPs. The simplest and most intuitive approach to examining such relations is to link plots dynamically such that samples selected in an oncoplot can be highlighted in other plots, and vice versa. There are, however, no existing oncoplot-generating R packages that support dynamic data linkage between different plots. To address this gap and enable rapid exploration of a variety of data types we constructed the ggoncplot package for the production of oncoplots that are easily integrated with custom visualisations and that support synchronised data-selections across plots (Figure 2). ggoncplot is available on GitHub at <https://github.com/selkamand/ggoncplot>.

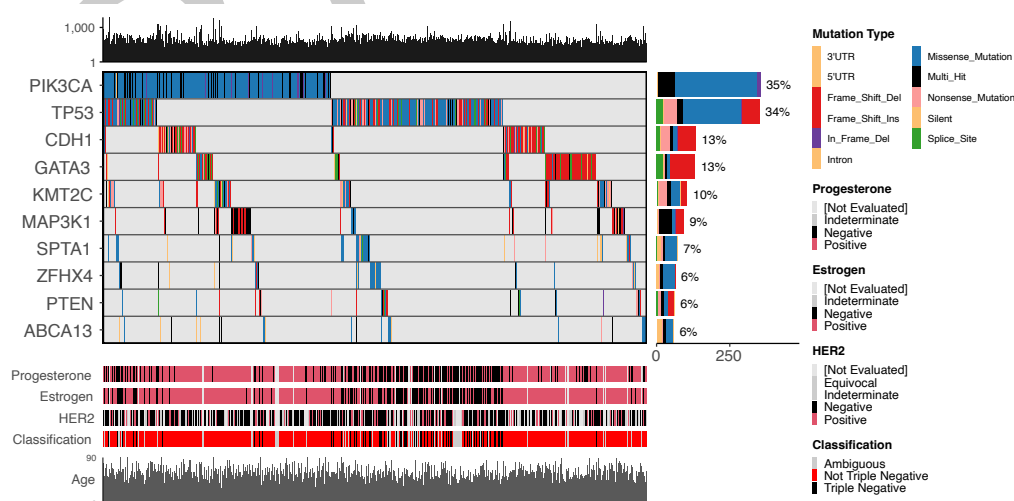


Figure 1: ggoncplot output visualising mutational trends in the TCGA breast carcinoma cohort. Individual patient samples are plotted on the x-axis, hierarchically sorted so that samples with the most frequent gene mutations appear on the leftmost side. The plot indicates that PIK3CA is the most frequently mutated gene, followed by TP53. Marginal plots indicate the total number of mutations per sample (top), and the number of samples showing mutations in each gene, coloured by mutation type (right). A range of clinical features, including progesterone and estrogen receptor status are shown on the marginal plot at the bottom. A detailed description of the ggoncplot sorting algorithm is available [here](#).

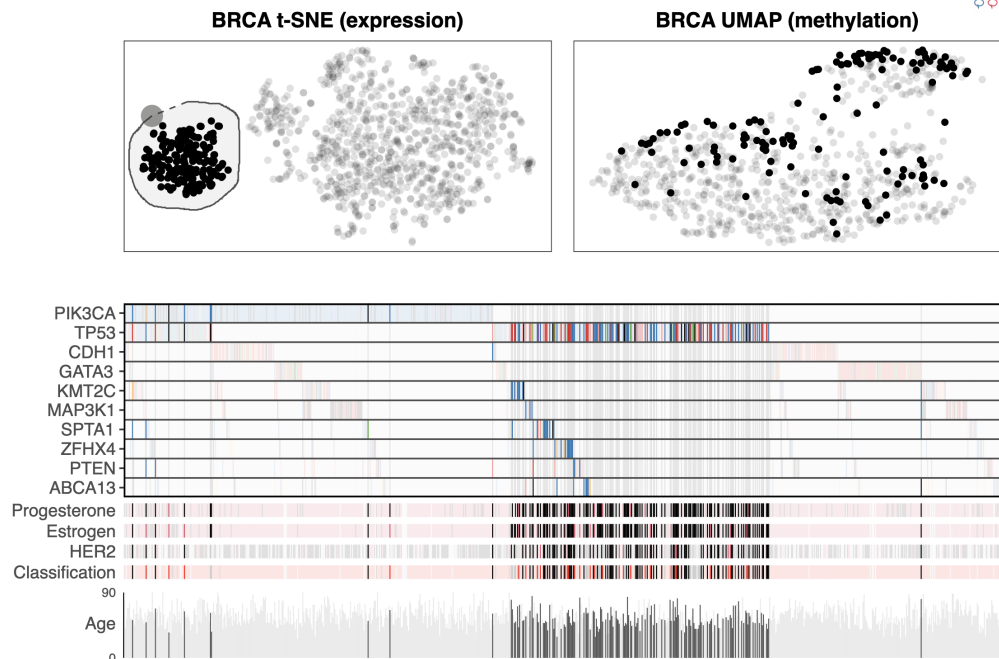


Figure 2: Example of the `ggoncplot` shown in Figure 1, where the oncplot has been dynamically cross-linked to a gene expression t-SNE plot (top left) and a methylation UMAP (top right). Here, the lasso tool was used to select a cluster of gene expression data points (i.e., individual samples) in the t-SNE plot. Selected samples were automatically highlighted on the UMAP and oncplot. This reveals that samples which cluster on the left of the t-SNE plot also cluster in the oncplot, chiefly containing mutations in TP53 and wild type PIK3CA. The plots of progesterone, estrogen, HER2 status and triple negative classification show that the samples selected in the t-SNE are enriched for triple negative breast cancers. In contrast to the oncplot, the methylation UMAP shows no strong clustering, consistent with knowledge of methylation patterns in triple negative breast cancer. Expression and methylation plots were produced using the `express` package.

Statement of Need

Oncoplots are highly effective for visualising mutation data in cancer cohorts but are challenging to generate with the major R plotting systems (base, lattice, or ggplot2) due to their algorithmic and graphical complexity. Simplifying the process of generating oncoplots would make them more accessible to researchers. Existing packages including ComplexHeatmap (Gu, 2022), maftools (Mayakonda et al., 2018), and genVisR (Skidmore et al., 2016) all make static oncoplots easier to create, but there is still a significant unmet need for a user-friendly method of creating oncoplots with the following features:

- **Interactive plots:** Customizable tooltips, cross-selection of samples across different plots, and auto-copying of sample identifiers on click. This enables exploration of multiomic datasets as shown in [Figure 2](#).
- **Support for tidy datasets:** Compatibility with tidy, tabular mutation-level formats (MAF files or relational databases), typical of cancer cohort datasets. This greatly improves the range of datasets that can be quickly and easily visualised in an oncoplot.
- **Auto-colouring:** Automatic selection of accessible colour palettes for datasets where the consequence annotations are aligned with standard variant effect dictionaries (PAVE, SO, or MAF).
- **Versatility:** The ability to visualize entities other than gene mutations, such as noncoding

features (e.g., promoter or enhancer mutations) and non-genomic entities (e.g., microbial presence in microbiome datasets).
We developed ggoncplot as the first R package to address all these challenges together (Figure 3). Examples of all key features are available in the ggoncplot manual.

Table with 5 columns: Property, complexheatmap, maftools, GenVisR, ggoncplot. Rows compare features like Sample sorting algorithm, Plotting framework, Automatic rendering of clinical annotations, etc.

Figure 3: Comparison of R packages for creating oncoplots. 1Requires the shiny and interactiveComplex-Heatmap packages. 2Exclusively colours tiles based on mutation impact which must be described using valid MAF variant classification terms. 3Requires the user to first summarise mutations at the gene level and format as a sample by gene matrix with mutations separated by semicolons (wide format). 4For MAF inputs the most severe consequence is chosen, however for non-MAF datasets users must manually define the mutation impact hierarchy. 5Non-unique mutation types are treated as one observation, however if different mutation types affect one gene, the individual mutations can be plotted with different shapes or sizes in a user-configured manner.

Acknowledgements

The results shown here are in whole or part based upon data generated by the TCGA Research Network: https://www.cancer.gov/tcga. Methylation, expression, and genome datasets were obtained from the Xena TCGA Pan-Cancer Atlas Hub (Goldman et al., 2020).

We thank the developers of the packages integral to ggoncplot, especially David Gohel for ggiraph (Gohel & Skintzos, 2024), which enables its interactivity, and Thomas Lin Pedersen for patchwork (Pedersen, 2024) and ggplot2 maintenance. We also acknowledge Hadley Wickham and all contributors to ggplot2 (Wickham, 2016). Additionally, we thank Dr. Marion Mateos for her insightful feedback during the development of ggoncplot.

References

Gohel, D., & Skintzos, P. (2024). Ggiraph: Make 'ggplot2' graphics interactive. https://davidgohel.github.io/ggiraph/

- 53 Goldman, M. J., Craft, B., Hastie, M., Repčeka, K., McDade, F., Kamath, A., Banerjee,
54 A., Luo, Y., Rogers, D., Brooks, A. N., Zhu, J., & Haussler, D. (2020). Visualizing and
55 interpreting cancer genomics data via the xena platform. *Nature Biotechnology*, 38(6),
56 675–678. <https://doi.org/10.1038/s41587-020-0546-8>
- 57 Gu, Z. (2022). Complex heatmap visualization. *iMeta*, 1(3), e43. [https://doi.org/10.1002/](https://doi.org/10.1002/imt2.43)
58 [imt2.43](https://doi.org/10.1002/imt2.43)
- 59 Mayakonda, A., Lin, D.-C., Assenov, Y., Plass, C., & Koeffler, H. P. (2018). Maftools:
60 Efficient and comprehensive analysis of somatic variants in cancer. *Genome Research*, 28,
61 1747–1756. <https://doi.org/10.1101/gr.239244.118>
- 62 Pedersen, T. L. (2024). *Patchwork: The composer of plots*. [https://patchwork.data-imaginist.](https://patchwork.data-imaginist.com)
63 [com](https://patchwork.data-imaginist.com)
- 64 Skidmore, Z. L., Wagner, A. H., Lesurf, R., Campbell, K. M., Kunisaki, J., Griffith, O. L., &
65 Griffith, M. (2016). GenVisR: Genomic visualizations in r. *Bioinformatics*, 32, 3012–3014.
66 <https://doi.org/10.1093/bioinformatics/btw325>
- 67 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
68 ISBN: 978-3-319-24277-4

DRAFT