

Data visualization

escheR: unified multi-dimensional visualizations with Gestalt principles

Boyi Guo ¹, Louise A. Huuki-Myers ², Melissa Grant-Peters ^{3,4},
Leonardo Collado-Torres ^{1,2}, Stephanie C. Hicks ^{1,5,6,7,*}

¹Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, 21205, United States

²Lieber Institute for Brain Development, Johns Hopkins Medical Campus, Baltimore, MD, 21205, United States

³Genetics and Genomic Medicine, Great Ormond Street Institute of Child Health, University College London, London, WC1N 1EH, United Kingdom

⁴Aligning Science Across Parkinson's (ASAP) Collaborative Research Network, Chevy Chase, MD, 20815, United States

⁵Department of Biomedical Engineering, Johns Hopkins School of Medicine, Baltimore, MD, 21218, United States

⁶Center for Computational Biology, Johns Hopkins University, Baltimore, MD, 21205, United States

⁷Malone Center for Engineering in Healthcare, Johns Hopkins University, Baltimore, MD, 21218, United States

*Corresponding author. Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, 615 N Wolfe St, Baltimore, MD, 21205, United States.
E-mail: shicks19@jhu.edu

Associate Editor: Sofia Forslund

Abstract

Summary: The creation of effective visualizations is a fundamental component of data analysis. In biomedical research, new challenges are emerging to visualize multi-dimensional data in a 2D space, but current data visualization tools have limited capabilities. To address this problem, we leverage Gestalt principles to improve the design and interpretability of multi-dimensional data in 2D data visualizations, layering aesthetics to display multiple variables. The proposed visualization can be applied to spatially-resolved transcriptomics data, but also broadly to data visualized in 2D space, such as embedding visualizations. We provide an open source R package *escheR*, which is built off of the state-of-the-art *ggplot2* visualization framework and can be seamlessly integrated into genomics toolboxes and workflows.

Availability and implementation: The open source R package *escheR* is freely available on Bioconductor (<https://bioconductor.org/packages/escheR>).

1 Introduction

Visualization is an indispensable component of data analysis, providing clarity that connects quantitative evidence to key conclusions (McGowan *et al.* 2023). In biomedical research, visualization receives growing recognition as essential: many scientists rely on visualization to complete their cognitive process from analysis to insight, including analytic validation of automated pipelines and scientific communication (O'Donoghue 2021). However, an important challenge in biomedical research is how to visualize increasingly complex, multi-dimensional data (O'Donoghue *et al.* 2010).

Here, we focus on two types of visualizations in biomedical research, but note that the proposed ideas could be extended beyond these applications: (i) embedding visualizations, which project data into some low-dimensional embedding or mathematical space [e.g. Principal Components Analysis, *t*-distributed Stochastic Neighbor Embedding (*t*-SNE), and Uniform Manifold Approximation and Projection (UMAP)] and (ii) *in situ* visualizations (Dries *et al.* 2021), which aim to visualize molecules captured from *in situ* imaging or sequencing technologies where *in situ* refers to 'in its original place'.

Both of these visualizations represent data in a 2D space and are motivated by recent advances in experimental technologies that profile molecules, including DNA, RNA, and proteins, at a single-cell or spatial resolution. Some most popular technologies include single-cell/nucleus RNA-sequencing (sc/snRNA-seq) and *in situ* spatially-resolved transcriptomics. (Vandereyken *et al.* 2023).

A common and fundamental challenge with both of these visualizations is how to visualize multi-dimensional information in a 2D space. For example, in *in situ* visualizations, we often want to create a spatial map to visualize a continuous (e.g. gene expression) or discrete (e.g. cell type or spatial domain) variable representing molecular information in the original spatial location. However, it is challenging to simultaneously visualize multi-dimensional data, such as information from disparate data domains (such as expression domain and spatial domain) or disparate data modalities (such as transcriptomics and proteomics) in the same plot. Currently, best practices for this include making two different plots displayed side-by-side (Fig. 1A and B), one for gene expression and one for spatial domains. This creates cognitive gaps on

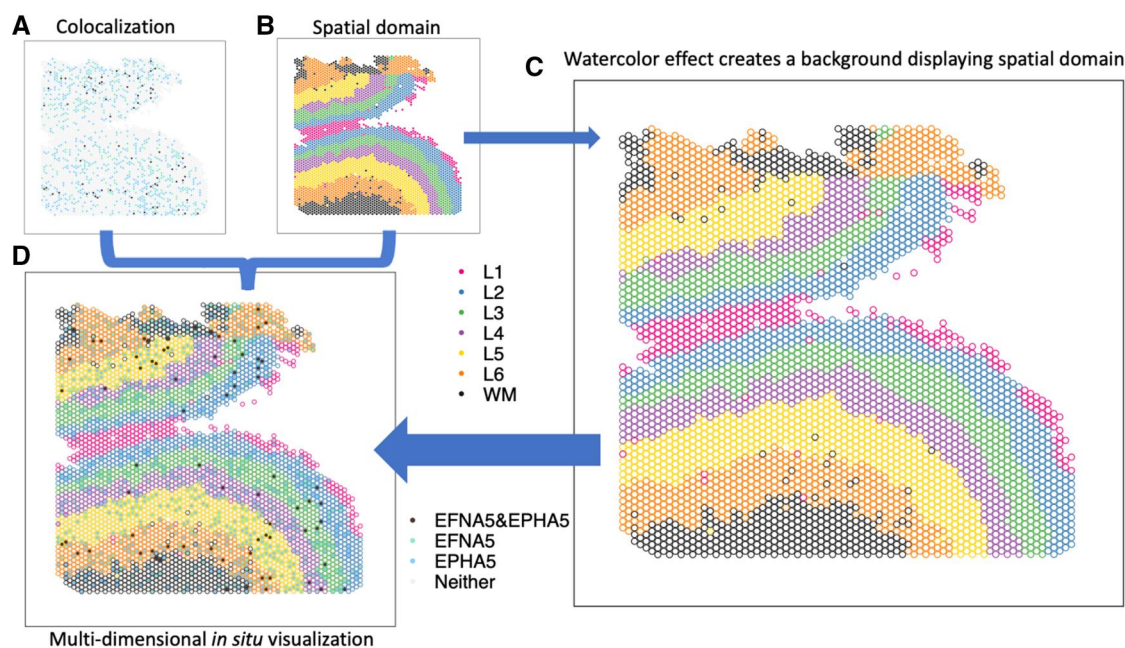


Figure 1. *escheR* enables multidimensional spatial visualizations following the Gestalt principles of design. (A and B) The traditional visualization displays the colocalization plot of the expression of two genes *EFNA5* and *EPHA5* (A) and the spatial domains from the dorsolateral prefrontal cortex in postmortem human brain (Huuki-Myers *et al.* 2023) (B) side-by-side, creating challenges to cognitively connecting colocalization status to spatial domains. (C) The watercolor effect enables displaying spatial domains by color-coding only outlines of circles. (D) *escheR* enables the multidimensional *in situ* visualization that simultaneously displays the cortex layers and the colocalization status, substantially improving interpretability.

how to associate the disparate information or how to interpret the biological findings of this multi-dimensional information regarding their microenvironment or colocalization. While interactive visualizations (Pardo *et al.* 2022, Sriworarat *et al.* 2023) have the potential to mitigate this challenge, they are infeasible for scientific communications in static media, such as printed work. Developing a static and unified visualization that enables the simultaneous display of multiple dimensions of information is crucial for biomedical research.

2 Results

To address these challenges, here we leverage the Gestalt (German for ‘unified whole’) principles (Todorovic 2008) as a path to visualize multi-dimensional data in 2D visualizations. We focus on the two types of data visualizations previously introduced that are widely used in biomedical research: (i) embedding visualizations and (ii) *in situ* visualizations. We provide an R package, *escheR*, implementing these ideas, which is built on the state-of-the-art data visualization framework *ggplot2* in the R programming language. Finally, we comment on how these ideas could be extended to other types of visualization in biomedical research.

2.1 Multi-dimensional 2D visualizations with *ggplot2* and gestalt principles

Gestalt principles (Todorovic 2008) refer to a set of rules describing how humans perceive and interpret visual information and are commonly applied in art and designs. Developed in the 1920s by German psychologists Max Wertheimer, Kurt Koffka and Wolfgang Kohler, these principles help humans perceive a set of individual elements as a whole.

Here, we leverage the principles to be able to visualize multi-dimensional data in a unified 2D plot. Our approach is

to use the state-of-art data visualization framework *ggplot2* (Wickham 2016) following the Grammar of Graphics (Wilkinson 2005) and map individual variables to different aesthetics to simultaneously display disparate variables. Specifically, we apply the figure-ground segmentation (Todorovic 2008) in displaying two variables: one variable (e.g. expression) can be plotted as color-filled circles, serving as the *figure*; one variable (e.g. spatial domains) can be plotted as the backgrounds of the circles, creating a *ground* for the figure. In practice, we use the combination of `color` and `fill="transparent"` to create the background layer and `fill` to create the figure layer. When necessary to display an additional layer for a third variable, `shape` can be used to add symbols such as cross (+) and asterisk (*) to highlight in the spatial map.

For adjacent circles with limited space between them to display the background color, we use an economic implementation, colored outlines for these circles (Fig. 1C), inspired by watercolor effect (Pinna *et al.* 2001). Watercolor effect describes the phenomenon in visual perception that surface color arises from thin boundaries and hence is applied here to perceive the background color in tight space. Overall, the figure-ground segmentation creates two isolated layers in visual perception to display the two variables while maintaining the relative spatial relationship serving as a reference between the two. In addition, other fundamental principles (Todorovic 2008), such as proximity, similarity, continuity, and closure, incentivize the brain to group elements and dimensions in the visualization, guaranteeing an integrative perception of the complex multi-dimensional spatial map.

Here, we provide an open-source package called *escheR* (named after the graphic artist M.C. Escher) in the R programming language (R Core Team 2022), leading to a simplified interface to navigate the implementation of the multi-dimensional visualization in 2D space. By adapting *ggplot2*

standard, *escheR* can be seamlessly integrated with popular data objects, including `SingleCellExperiment` (Amezquita *et al.* 2020) and `SpatialExperiment` (Righelli *et al.* 2022), and easily work with popular pipelines, such as `Seurat` (Hao *et al.* 2021), `Giotto` (Dries *et al.* 2021) to name a few, and allow further theme customization with ease.

Next, we give two use cases to exemplify some utility of the proposed spatial visualization: (i) the spatially differential gene colocalization in the human dorsolateral prefrontal cortex using spatial transcriptomics data (Huuki-Myers *et al.* 2023); and (ii) multi-dimensional UMAP highlighting differential gene expression in data-driven cell clusters (Freytag and Lister 2020).

2.2 Multi-dimensional *in situ* visualization

In a recent study investigating the molecular organization of human dorsolateral prefrontal cortex (Huuki-Myers *et al.* 2023), two schizophrenia risk genes, membrane-bound ligand ephrin A5 (*EFNA5*) and ephrin type-A receptor 5 (*EPHA5*), were identified to colocalize via cell-cell communication analysis. In addition, data suggested Layer 6 was the most highly co-localized layer compared to other cortex layers. To visually examine the inference, we applied *escheR* to create a multi-dimensional *in situ* spatial map that simultaneously exhibits the cortex layers (displayed with color-coded spot outlines) and the categorized colocalization status of genes *EFNA5* and *EPHA5* (displayed with color-coded spot fill). Compared to the traditional visualization where the cortex layers and the colocalization status are visualized in two side-by-side figures (Fig. 1A and B), our proposed visualization (Fig. 1D) enables directly mapping colocalization

status to the spatial domain, simplifying the perception of two sources of information and allowing cognitive comparison across cortex layers.

2.3 Multi-dimensional embedding visualizations

The application of the proposed framework is not limited to *in situ* visualizations of spatially-resolved transcriptomics data. It is broadly applicable to data mapped to any 2D coordinate system to simultaneously display multiple variables. Such systems include euclidean space (including spatial coordinate as a special case) and data-driven embedding space, for example, UMAP and *t*-SNE. To demonstrate, we applied the proposed visualization to address the challenge of simultaneously displaying cluster membership and gene expression in a single-cell UMAP plot. To address the overplotting problem, previous work proposed to apply hexagonal binning strategy to display the gene expression (Freytag and Lister 2020). Here, the color-coded convex hulls are used to annotate different clusters of hexagons, where the cluster label of each hexagon is determined by the most popular label within each hexagon bin (Fig. 2A). However, the convex hulls create substantial overlapping areas, creating confusion when interpreting cluster memberships of hexagons in the overlapping areas. To improve the interpretability of the visualization, we replace the convex hulls with color-coded hexagons boundaries (Fig. 2B) to avoid possible membership confusion. We note that our contribution to improving the visualization is easily implemented without any modification of *scheX* as both are built upon the Grammar of Graphics (Wilkinson 2005) standard.

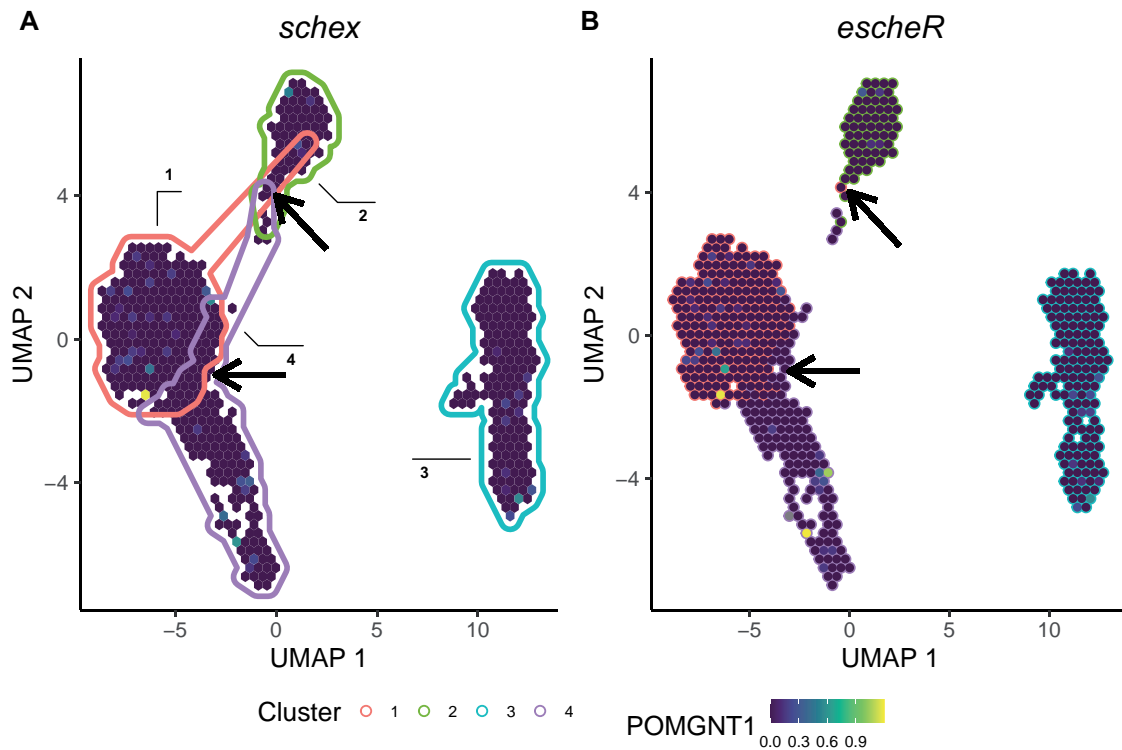


Figure 2. *escheR* improves multidimensional embedding visualizations. The gene expression of *POMGNT1* among peripheral blood mononuclear cells (Hansen *et al.*, 2022) under the UMAP representation. (A) The *schex* R/Bioconductor package uses color-coded convex hulls to annotate data-driven cell types, creating confusion when interpreting hexagons in overlapping hulls. (B) *escheR* plots hexagon-specific membership to avoid substantial overlapping of cluster membership. Due to the binning strategy, cautions in interpretation are needed due to possible cluster intermixing within each hexagon.

3 Discussion

Here, we propose an innovative multi-dimensional spatial visualization that simultaneously displays multiple variables in a 2D coordinate system. Specifically, our design leverages Gestalt principles from visual perception to create multiple visual dimensions in a spatial map by iterative layering aesthetics. Built upon `ggplot2`, we provide an open-source R package `escheR` that is seamlessly compatible with popular spatially-resolved transcriptomics and single-cell data analysis toolboxes.

Adding a third dimension to 2D plots has been a long-standing challenge in visualization (O'Donoghue *et al.* 2010). Our proposal addresses this fundamental challenge by introducing simple but effective design principles. These principles lead to visually easier-to-interpret graphics compared to 2D-color gradients and geometry annotations (Fig. 2). Unlike computer-based interactive visualizations, the proposed visualization is free from any platform and technology restriction, creating an accessible and economical solution. In addition, the proposed visualization is easily scalable and hence can be applied to all types of spatially resolved data.

Beyond the scope of biomedical research, the proposed visualization can be broadly translated to any visual analytic highlighting differentiation with respect to other measurements. To name a few, such visual analytics include examining differential tests, explaining clustering, and visualizing subgroups. However, one of the most rewarding fields to apply the proposed visualization is the rapidly expanding field of biomedical multi-omics research (Hasin *et al.* 2017), where connecting different omics (data modalities) is the fundamental goal and hence greatly appreciating innovative multi-dimensional visualization.

While the proposed design principle set-up a solid foundation, the proposed package `escheR` can be further optimized for image-based spatially resolved data. For example, to address possible overplotting of densely arranged single cells in tight space, strategies such as hatching plots (Patrick *et al.* 2023) can be incorporated. In addition, the flexibility to represent data points as polygons will be provided to allow visualization of cell morphology. In summary, we propose a novel multi-dimensional visualization, implemented in an R package `escheR`, to address the simultaneous exhibition of multiple variables in 2D plots. The proposed visualization can be broadly applicable to the visual analytics of growingly complex biomedical data and beyond.

Acknowledgements

The authors thank the anonymous reviewers for their valuable suggestions. We also would like to acknowledge Nicholas J. Eagles, Kristen R. Maynard, Mina Ryten, Leon Di Stefano, and Lukas M. Weber (appearing in alphabetic order of last name) for their helpful comments, feedback, and suggestions on `escheR` functionality. Nicholas J. Eagles and Kristen R. Maynard are employed by the Lieber Institute for Brain Development; Mina Ryten is employed by University College London; Leon Di Stefano is from the Johns Hopkins Bloomberg School of Public Health, Department of Biostatistics; Lukas M. Weber is in the Department of Biostatistics at Boston University School of Public Health.

Author contributions

B.G.: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data Curation, Writing, Visualization; L.A.H.-M.: Conceptualization, Software; M. G.-P.: Conceptualization; L.C.-T.: Conceptualization, Software; S.C.H.: Conceptualization, Resources, Writing—Review & Editing, Visualization, Supervision, Project administration, Funding acquisition.

Conflict of interest

None declared.

Funding

This project was supported by the National Institute of Mental Health [R01MH126393 to B.G. and S.C.H., U01MH122849 to L.A.H.-M. and L.C.-T.]; the Chan Zuckerberg Initiative DAF, an advised fund of Silicon Valley Community Foundation [CZF2019-002443 to S.C.H.]; the Lieber Institute for Brain Development to L.A.H.-M. and L. C.-T.; and Aligning Science Across Parkinson's [ASAP-000478, ASAP-000509 to M.G.-P.] through the Michael J. Fox Foundation for Parkinson's Research. All funding bodies had no role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

Data availability

The spatial transcriptomics dataset was obtained from `spatialLIBD` (<http://research.libd.org/spatialLIBD>). The UMAP example follows the 'using_schex' vignette in the `schex` package (<http://bioconductor.org/packages/schex>). The code that generates these figures is deposited at http://github.com/boyigu01/Manuscript_escheR (Zenodo DOI: 10.5281/zenodo.7915970). The open source software package `escheR` available in the R programming language is freely available on GitHub (<http://github.com/boyigu01/escheR>) and Bioconductor (<http://bioconductor.org/packages/escheR>).

References

- Amezquita RA, Lun AT, Becht E *et al.* Orchestrating single-cell analysis with bioconductor. *Nat Methods* 2020;17:137–45.
- Dries R, Chen J, Del Rossi N *et al.* Advances in spatial transcriptomic data analysis. *Genome Res* 2021;31:1706–18.
- Freytag S, Lister R. `schex` avoids overplotting for large single-cell RNA-sequencing datasets. *Bioinformatics* 2020;36:2291–2.
- Hansen KD, Risso D, Hicks SC. `TENxPBMCDATA`: PBMCDATA from 10X genomics. R Package Version 1.20.0. 2022. <https://bioconductor.org/packages/TENxPBMCDATA> (15 October 2023, date last accessed).
- Hao Y, Hao S, Andersen-Nissen E *et al.* Integrated analysis of multi-modal single-cell data. *Cell* 2021;184:3573–87.e29.
- Hasin Y, Seldin M, Lusic A. Multi-omics approaches to disease. *Genome Biol* 2017;18:83.
- Huuki-Myers L, Spangler A, Eagles N *et al.* A data-driven single cell and spatial transcriptomic map of the human prefrontal cortex. *Science*, 2023. In press.
- McGowan LD, Peng RD, Hicks SC. Design principles for data analysis. *J Comput Graph Stat* 2023;32:754–61.
- O'Donoghue SI. Grand challenges in bioinformatics data visualization. *Front Bioinform* 2021;1:669186. <https://doi.org/10.3389/fbinf.2021.669186>

- O'Donoghue SI, Gavin A-C, Gehlenborg N *et al.* Visualizing biological data-now and in the future. *Nat Methods* 2010;7:S2–4.
- Pardo B, Spangler A, Weber LM *et al.* spatialLIBD: an R/bioconductor package to visualize spatially-resolved transcriptomics data. *BMC Genomics* 2022;23:434.
- Patrick E, Canete NP, Iyengar SS *et al.* Spatial analysis for highly multiplexed imaging data to identify tissue microenvironments. *Cytometry A* 2023;103:593–9.
- Pinna B, Brelstaff G, Spillmann L. Surface color from boundaries: a new 'watercolor' illusion. *Vision Res* 2001;41:2669–76.
- R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing, 2022.
- Righelli D, Weber LM, Crowell HL *et al.* SpatialExperiment: infrastructure for spatially-resolved transcriptomics data in R using bioconductor. *Bioinformatics* 2022;38:3128–31.
- Sriworarat C, Nguyen A, Eagles NJ *et al.* Performant web-based interactive visualization tool for spatially-resolved transcriptomics experiments. *Biol Imaging* 2023;3:E15. <https://doi.org/10.1017/S2633903X2300017X>
- Todorovic D. Gestalt principles. *Scholarpedia* 2008;3:5345.
- Vandereyken K, Sifrim A, Thienpont B *et al.* Methods and applications for single-cell and spatial multi-omics. *Nat Rev Genet* 2023;24:494–515.
- Wickham H. *ggplot2: Elegant Graphics for Data Analysis*. New York, NY: Springer, 2016.
- Wilkinson L. *The grammar of graphics*. New York, NY, USA: Springer, 2005.