

# Annual Review of Biomedical Data Science Visualization of Biomedical Data

Seán I. O'Donoghue,<sup>1,2,3</sup> Benedetta Frida Baldi,<sup>2</sup>  
Susan J. Clark,<sup>2</sup> Aaron E. Darling,<sup>4</sup> James M. Hogan,<sup>5</sup>  
Sandeep Kaur,<sup>6</sup> Lena Maier-Hein,<sup>7</sup>  
Davis J. McCarthy,<sup>8,9</sup> William J. Moore,<sup>10</sup>  
Esther Stenau,<sup>7</sup> Jason R. Swedlow,<sup>10</sup> Jenny Vuong,<sup>1</sup>  
and James B. Procter<sup>10</sup>

<sup>1</sup>Data61, Commonwealth Scientific and Industrial Research Organisation (CSIRO), Eveleigh NSW 2015, Australia; email: sean@odonoghuelab.org

<sup>2</sup>Genomics and Epigenetics Division, Garvan Institute of Medical Research, Sydney NSW 2010, Australia

<sup>3</sup>School of Biotechnology and Biomolecular Sciences, University of New South Wales (UNSW), Kensington NSW 2033, Australia

<sup>4</sup>The ithree Institute, University of Technology Sydney, Ultimo NSW 2007, Australia

<sup>5</sup>School of Electrical Engineering and Computer Science, Queensland University of Technology, Brisbane QLD, 4000, Australia

<sup>6</sup>School of Computer Science and Engineering, University of New South Wales (UNSW), Kensington NSW 2033, Australia

<sup>7</sup>Division of Computer Assisted Medical Interventions (CAMI), German Cancer Research Center (DKFZ), 69120 Heidelberg, Germany

<sup>8</sup>European Bioinformatics Institute (EBI), European Molecular Biology Laboratory (EMBL), Wellcome Genome Campus, Hinxton CB10 1SD, United Kingdom

<sup>9</sup>St. Vincent's Institute of Medical Research, Fitzroy VIC 3065, Australia

<sup>10</sup>School of Life Sciences, University of Dundee, Dundee DD1 5EH, United Kingdom

## ANNUAL REVIEWS CONNECT

[www.annualreviews.org](http://www.annualreviews.org)

- Download figures
- Navigate cited references
- Keyword search
- Explore related articles
- Share via email or social media

Annu. Rev. Biomed. Data Sci. 2018. 1:275–304

First published as a Review in Advance on  
May 16, 2018

The *Annual Review of Biomedical Data Science* is  
online at [biodatasci.annualreviews.org](http://biodatasci.annualreviews.org)

<https://doi.org/10.1146/annurev-biodatasci-080917-013424>

Copyright © 2018 Seán I. O'Donoghue et al. This work is licensed under a Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited. See credit lines of images or other third-party material in this article for license information



## Keywords

data visualization, multivariate data, molecular biology, cell biology, tissue imaging, metagenomics

## Abstract

The rapid increase in volume and complexity of biomedical data requires changes in research, communication, and clinical practices. This includes learning how to effectively integrate automated analysis with high-data density visualizations that clearly express complex phenomena. In this review, we summarize key principles and resources from data visualization research that help address this difficult challenge. We then survey how visualization is being used in a selection of emerging biomedical research areas, including three-dimensional genomics, single-cell RNA sequencing (RNA-seq), the protein structure universe, phosphoproteomics, augmented reality-assisted surgery, and metagenomics. While specific research areas need highly

tailored visualizations, there are common challenges that can be addressed with general methods and strategies. Also common, however, are poor visualization practices. We outline ongoing initiatives aimed at improving visualization practices in biomedical research via better tools, peer-to-peer learning, and interdisciplinary collaboration with computer scientists, science communicators, and graphic designers. These changes are revolutionizing how we see and think about our data.

The eye of a Master, will do more Work than his Hand.

—Benjamin Franklin, *Poor Richard's Almanack* (1744)

## INTRODUCTION

The launch of this Annual Reviews journal is driven by the rapid increase in volume and complexity of biomedical data, requiring changes in research, communication, and clinical practices (1). Without these changes, many biomedical discoveries will remain buried in data already collected, and many misdiagnoses (now estimated at ~10–30% of all diagnoses) will remain unrecognized (2, 3), contributing to a major cause of death (4).

These changes in practice will include the development and adoption of new, automated analysis methods (e.g., clustering, modeling, machine learning). However, while necessary (**Figure 1a**), automated analysis is not sufficient: As demonstrated by Anscombe's quartet (**Figure 1b**), to find the truth, we need to visually inspect all relevant data and analyses together. Driven by this realization, data visualization has been a major research focus in computer science for decades (see the sidebar titled Data Visualization, SciVis, and InfoVis), yielding many resources that could accelerate discovery in biomedical research (5–9).

Unfortunately, relatively few scientists currently use these resources. This is evident, for example, from the widespread use of rainbow color maps (found in >50% of papers in a survey of research publications containing scientific visualizations; 10). While seemingly inconsequential, visualization research has shown that rainbow maps can obfuscate true data patterns and introduce visual artifacts (10). Sadly, many biomedical data sets (often difficult or expensive to acquire) are inspected using poor visualization methods, even though better alternatives are known.

Similarly, visualization methods are underutilized by clinicians, contributing to misdiagnoses. About half of all diagnostic errors arise from faulty cognitive processing of data (2); many of these errors can be addressed by improving how data are visualized (11). This is especially true in fields of medicine, such as radiology, where the core data are intrinsically visual, and diagnosis depends largely on visual perception (3).

We believe that the underuse of visualization methods has arisen largely because of the following misconceptions we often encounter.

Misconception 1: “The goal of data visualization is to impress.” We sometimes think of data visualization as purely aesthetic, adding an optional wow factor not present in the data itself. This can be true when creating artwork (e.g., a cover figure), but the role of data visualization in research is almost exactly the opposite: It is a necessary step, aimed at clearly revealing patterns in data.

Misconception 2: “Data visualization is easy.” Well-designed visualizations can be so easy to understand and use that we are misled into thinking they must have been easy to create. However:

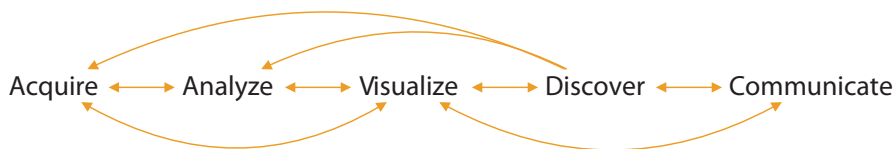
---

### Data visualization:

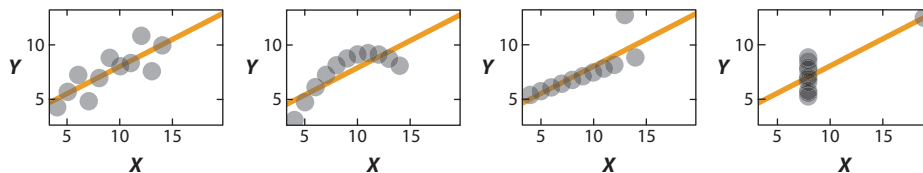
use of computer-aided, interactive visual representations of data to amplify cognition and accelerate discovery and communication

---

## a Research workflow



## b Anscombe's quartet



**Figure 1**

Role of data visualization in research. (a) As shown in this simplified model of the research workflow, data visualization is often a necessary and rate-limiting step in both discovery and communication. (b) Anscombe's quartet (117, 118) is a set of four two-dimensional data sets in which  $X$  and  $Y$  values have identical mean, variance, and correlation coefficient. They also fit an identical linear regression line (orange) with an identical coefficient of determination. Based on these statistics alone, we might expect all plots to be similar to the first, but visualization reveals surprisingly distinct patterns in each data set. This demonstrates that we cannot skip from analysis to discovery: It is almost always necessary to confirm insights from automated analysis by manually visualizing data. Panel a adapted from Reference 119. Panel b adapted with permission from Reference 117.

“Most graphs are simple, but their invention was neither simple nor obvious—the idea did not occur to the Greeks or Romans, nor even to the great 17th century mathematician-experimenters such as Newton and Leibniz” (12, p. 201).

Misconception 3: “Studying data visualization is unnecessary.” Underestimating the difficulty of data visualization can lead us to overestimate our current skills and conclude that we would gain little benefit from investing time, effort, or money in training or study.

Misconception 4: “Visualization is just a synonym for imaging.” In the life sciences, “visualization” is often used as a synonym for “imaging experiments.” In computer science, “data visualization” has a much broader meaning: In addition to imaging data, it encompasses abstract data, interactive analysis, and design, as well as visual and cognitive abilities. Furthermore, its purpose is insight, not pictures (9).

## DATA VISUALIZATION, SciVis, AND InfoVis

Computer scientists have long used the term “scientific visualization” (or “SciVis”) to describe visualization of data that directly map into two or three spatial dimensions (e.g., cartography, computed tomography scans). In contrast, the term “information visualization” (or “InfoVis”) is used to describe visualization of abstract data (e.g., classic two-dimensional data plots, network graphs). Since around the year 2000, “data visualization” has emerged as a unifying term that encompasses both of these historically separated research fields.

## AVOIDING COMMON VISUALIZATION MISTAKES

Biomedical data sets are often difficult and expensive to acquire and analyze. Ironically, when these data are visualized, we often use techniques that obfuscate true patterns in our data and introduce visual artifacts. To avoid the most common mistakes, we recommend the following strategies. (a) Avoid rainbow color maps (10, 110). (b) Use color minimally (111); color used poorly is worse than no color at all (112). (c) Avoid creating confusing, overcrowded visualizations (e.g., hairball graphs); reduce information via filtering or clustering or use a different layout (Figure 7). (d) Use three-dimensional (3D) visualizations only for spatial data that is intrinsically 3D, e.g., macromolecular structures (113), and avoid using 3D visualizations for abstract data. (e) Avoid conflating research and art. Many of the commonly used tailored tools provide powerful features that make it easy to create visualizations that are dramatic or aesthetically appealing but where the underlying scientific meaning becomes obscured. This can be useful when creating impactful artwork (e.g., a cover figure), but it undermines the goals of data visualization in research, which are always clarity and insight.

Below, we outline some key data visualization principles, many of which are straightforward, helping to create better visualizations and avoid common mistakes (see the sidebar titled Avoiding Common Visualization Mistakes). We then survey a selection of emerging biomedical research areas where visualization is playing a key role. Finally, we discuss how visualization can enhance biomedical communication, and in conclusion, we consider prospects for improving the global standard of data visualization in biomedical research.

## DATA VISUALIZATION PRINCIPLES

### Data Volume

How can data visualization help us deal with the increasingly large volume of biomedical data sets?

One straightforward answer is to get more pixels. Our visual system has extraordinary capacity and can manage much more information than is presented in many of the scientific visualizations we currently create. Information transfer speed from the eye to the brain is about 10 Mb/s, similar to wired internet (13), and well-encoded visual patterns can be recognized within 250 ms (preattentively) (14). Large, higher-resolution displays can help use more of our visual capacity. For example, connecting a laptop to a 4K display can be a cost-effective way to see more detail and improve navigation and work efficiency (15). However, larger displays become increasingly less cost effective and often have impractical user interface controls. In addition, scaling up a visualization can make global patterns (e.g., correlation) harder to perceive (16)—clearly, there is an optimal size range for visualizations.

Therefore, a second answer is to create visualizations with greater data density. Here again, our visual system has greater capacity than we typically use: Aided by redundancy and context, the eye can resolve features to 0.1 mm (5). Thus, many visualization researchers advocate creating visualizations with high data density, often a practical requirement in high-impact journals due to space limitations (5). Creating compact visualizations requires carefully selecting visual channels that encode data with high visual effectiveness. Fortunately, visualization research provides clear guidelines for this (Figure 2). Also required is time and effort to learn specialist tools (Table 1) that provide precise controls necessary for visualizing large data volumes with high data density. However, there are limits: As data density increases, visualizations can require specific, targeted strategies to remain effective (17).

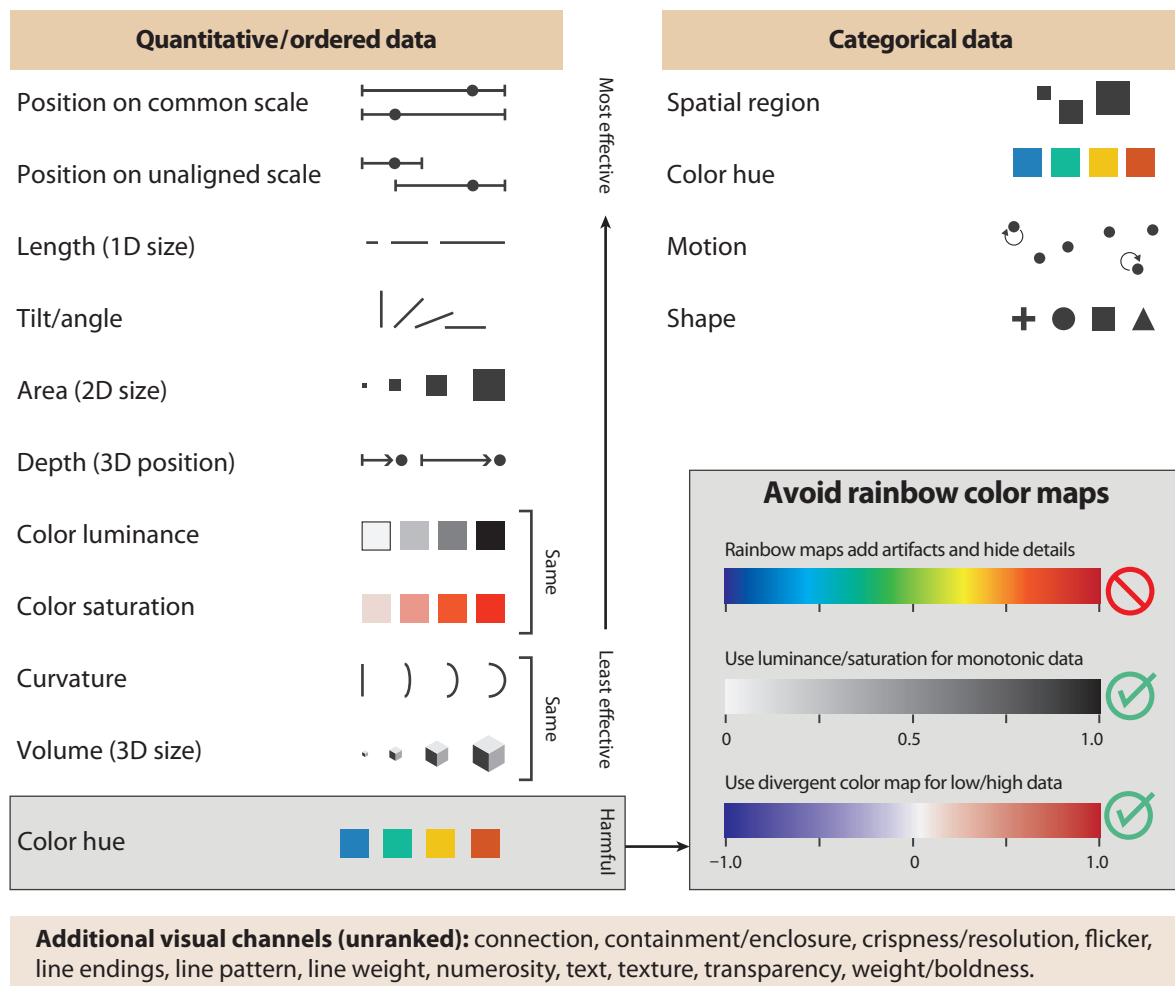
---

**Data density:** the total number of data entries shown in a visualization divided by the display area

**Visual channel:** an elementary graphical strategy used for visually encoding data (e.g., using color hue to show data categories)

**Visual effectiveness:** the accuracy and clarity with which a given visual encoding of data is conveyed to the reader

---



**Figure 2**

Visual channel ranking, showing the most- to least-effective visual channels (*top to bottom*, respectively) for encoding quantitative/ordered or categorical data. These rankings can be a useful guide in designing new visualizations and in studying exemplary visualizations (e.g., **Figures 3** and **7a**). Unfortunately, the use of color hue to encode quantitative data is still widespread in science even though visualization research clearly demonstrates that this can be not only ineffective but also harmful, introducing visual artifacts and hiding details (10). Instead (*bottom right inset*), it is more effective to use color maps tailored for specific data ranges (120). Figure adapted with permission from Reference 8, based on the approach pioneered by Mackinlay (121) and extended by others (107).

## Data Complexity

In many ways, big data volume is a small problem in the life sciences. Far more challenging is the complexity of our data, which are often multivariate, multiscale, highly interconnected, and dependent on very specific conditions.

Here, a common strategy is to use analytical methods to reduce dimensionality (e.g., clustering, principal component analysis). However, Anscombe's quartet (**Figure 1b**) reminds us that we need to visually inspect all relevant data before we draw conclusions from a simplified subset, and very often we require many more than two variables to express the complex phenomena studied in biomedical data sets.

**Table 1** Data visualization resources recommended for biomedical scientists in any field<sup>a</sup>

Resource	Description	URL
<b>Discovery<sup>b</sup></b>		
Excel <sup>c</sup>	Everyday tool for generic visualization of smaller data sets	<a href="http://microsoft.com/excel">http://microsoft.com/excel</a>
Plotly	Online tool for fast data visualization	<a href="https://plot.ly/create/">https://plot.ly/create/</a>
Tableau <sup>c</sup>	For interactive visualizations, including web based	<a href="http://tableau.com">http://tableau.com</a>
Spotfire <sup>c</sup>	For visual analysis of larger data sets and tool generation	<a href="https://spotfire.tibco.com/">https://spotfire.tibco.com/</a>
Origin <sup>c,d</sup>	For visual analysis of larger data sets	<a href="http://originlab.com">http://originlab.com</a>
Mathematica <sup>c</sup>	For visual analysis of data sets and mathematical functions	<a href="http://wolfram.com">http://wolfram.com</a>
MATLAB <sup>c</sup>	For visual analysis of data sets and mathematical functions	<a href="http://mathworks.com">http://mathworks.com</a>
Matplotlib	For tailored visualizations of data sets in Python (115)	<a href="http://matplotlib.org">http://matplotlib.org</a>
ggplot2	For tailored visualizations of large, complex data sets in R (116)	<a href="http://ggplot2.org">http://ggplot2.org</a>
D3.js	For tailored, interactive web-based visualizations	<a href="https://d3js.org">https://d3js.org</a>
<b>Communication</b>		
Photoshop <sup>c</sup>	For editing imaging data	<a href="http://adobe.com/photoshop">http://adobe.com/photoshop</a>
GIMP	Free, open-source alternative to Photoshop	<a href="http://www.gimp.org">http://www.gimp.org</a>
Illustrator <sup>c</sup>	For creating and editing vector graphics	<a href="http://adobe.com/illustrator">http://adobe.com/illustrator</a>
Inkscape	Free, open-source alternative to Illustrator	<a href="http://inkscape.org">http://inkscape.org</a>
MolecularMaya	Molecular structure plug-in for Autodesk Maya <sup>c</sup> animation suite	<a href="http://bit.ly/molmaya">http://bit.ly/molmaya</a>
BioBlender	Molecular structure plug-in for Blender animation suit	<a href="http://bioblender.org">http://bioblender.org</a>
<b>Utilities</b>		
Color Brewer	Web tool for selecting contrasting color maps	<a href="http://colorbrewer2.org">http://colorbrewer2.org</a>
Adobe Color	Web tool for designing sets of colors	<a href="http://color.adobe.com">http://color.adobe.com</a>
Paletton	Web tool for designing sets of colors	<a href="http://paletton.com">http://paletton.com</a>
<b>General Resources</b>		
BioVis	Computer science publications on biological visualizations	<a href="http://biovis.net">http://biovis.net</a>
Clarafi <sup>c</sup>	Training guides for biomedical visualization tools	<a href="http://clarafi.com">http://clarafi.com</a>
Information is Beautiful	Showcase of charts and infographics for a wide variety of data	<a href="http://bit.ly/Info_Beauty">http://bit.ly/Info_Beauty</a>
Visual Complexity	Catalog of tailored visualizations for complex data	<a href="http://visualcomplexity.com">http://visualcomplexity.com</a>
VIZBI	Collected videos and posters on tailored biological visualizations	<a href="http://vizbi.org">http://vizbi.org</a>
<b>Exemplars</b>		
PDB101	Outstanding visual explanations of protein function and structure	<a href="https://pdb101.rcsb.org">https://pdb101.rcsb.org</a>
Roche pathway	Tailored visualization showing ~3,000 metabolic reactions (72)	<a href="http://bit.ly/RochePathway">http://bit.ly/RochePathway</a>
WEHI.tv	Collection of inspiring, informative biomedical animations	<a href="http://wehi.tv">http://wehi.tv</a>

<sup>a</sup>This table covers only tools and online resources; published articles and books describing generally useful visualization methods are highlighted as annotated references in the Literature Cited. Visualization methods tailored for specific fields of biomedical research are given in the section titled Visualization for Discovery.

<sup>b</sup>Arranged in approximate order, starting at the top with easy-to-use, everyday tools for generic, one-off visualization tasks and progressing to tools that require more time and effort to use but that can manage large, complex data or reoccurring tasks.

<sup>c</sup>These tools cost money to use; the rest are free.

<sup>d</sup>Requires Microsoft Windows.

Remarkably, multivariate data of any dimensionality can be visualized in two dimensions without loss of information, using a range of generic methods (Figure 3a–d) and tools (Table 1). However, multidimensional data patterns are often scrambled and hard to recognize or interpret when encoded into two dimensions (18).

## Data Integration and Tailored Visualizations

Clearly revealing these multidimensional patterns usually requires carefully tailored visualizations that use very specific data integration strategies. Figure 3e reproduces an exemplary tailored visualization created by Charles Minard in 1869, showing Napoleon’s disastrous Russian campaign (5). To learn how to make better visualizations, readers should study in detail (guided by Figure 2) how this complex, multivariate data story has been communicated so clearly and concisely. For example, Minard’s graph shows that tailored visualizations sometimes require using less effective channels: In Figure 3e, army size would be more effectively encoded using a bar chart, but this would not be as visually expressive. Similarly, it can sometimes be required to break visual conventions: In Figure 3e, the bottom plot implicitly shows time flowing right to left, with irregularly spaced dates. Despite these departures from recommended guidelines, Minard’s graph has been described as perhaps “the best statistical graph ever drawn” (5, p. 40).

While generic visualization methods (Figure 3a–d) and tools are frequently used in biomedical research, tailored visualizations are the mainstay. Creating them requires three inter-related steps. (a) First, identify the necessary complexity. This is the subset of data that can visually express all, and only, information of most relevance to the phenomena studied. This usually means excluding data (e.g., in Figure 3e, temperature is shown only when it caused significant losses) or showing only derived features (e.g., principal components). (b) Second, identify necessary supporting context. This often means adding related information not part of the initial data set (e.g., in Figure 3e, geographic features help interpret the data set). (c) Third, invent a visual strategy that coherently integrates these data, using effective visual encoding and conventions familiar to peers.

A successful tailored visualization arranges all relevant data into a compact, immediately accessible two-dimensional (2D) view. This facilitates spatial reasoning, which, in turn, reduces the cognitive load needed to read a visualization and gain insight from data (19). By using familiar or intuitive visual conventions, successful visualization strategies also reduce the cognitive load needed when first learning to how to read them.

Fortunately, for many of the data challenges faced in biomedical research (20–24), tailored visualization methods have already been invented and implemented into working tools. These tools often use interactivity (25, 26) to facilitate combined exploration and integration of raw data, data derived via analysis, and additional supporting evidence.

However, cutting-edge research often requires us to invent novel, tailored visualizations. This difficult task can be aided by drawing ideas and inspiration from visualization resources (see Table 1 and the Annotated References 1, 5–9, 27, and 28) and by learning from peers facing similar challenges. Thus, the next section surveys tailored visualizations being used to reveal new insights across a broad range of biomedical research areas.

## VISUALIZATION FOR DISCOVERY

### Genomics and Epigenetics

We begin this survey in the field of genomics and epigenetics, where rapid advances in DNA sequencing technologies are generating a flood of data. These data are not just limited to raw DNA sequences but include an increasing spectrum of additional information that can be obtained

---

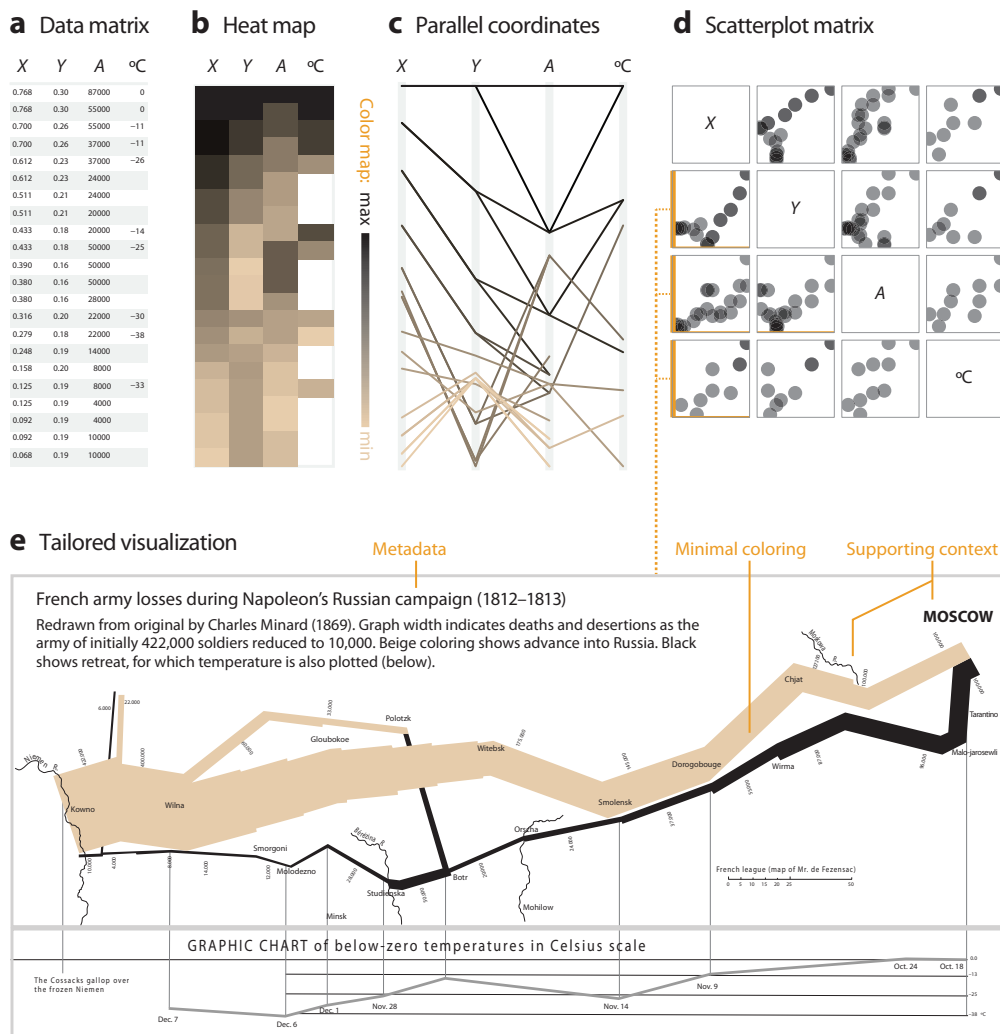
**Multivariate data:** data comprising multiple variables of any type, including quantitative, categorical (“A is cytoplasmic”), or relational (“A binds B”)

**Tailored visualization:** a strategy designed for integrating specific types of data sets, with supporting context, in a manner understood by peers

**Visual expressiveness:** how well a visualization expresses all—and only—information most relevant to the phenomena studied

**Spatial reasoning:** use of visual perception to enhance cognition; aided by organizing relevant data on a graphical display

---



**Figure 3**

Different two-dimensional (2D) views of a 4D data set. (a) Showing a data matrix as text reveals all information in two dimensions, but patterns can be difficult to detect. (b) In a heat map (122), the data matrix is visually encoded using color. This can give high data density; however, revealing data patterns often requires reordering rows and columns. Unfortunately, optical illusions can mask true patterns and introduce visual artifacts (47). (c) In parallel coordinate plots (18, 28), columns of the data matrix are represented as parallel axes, and each row becomes a series of line segments joining the axes. Correlations between adjacent axes can be easily seen, but not between nonadjacent axes. Hence, revealing data patterns often requires reordering axes and careful choice of a color map (here used to show  $X$  values). Parallel coordinates can reveal patterns not easily seen in a heat map, but they are usually not as compact. (d) A scatterplot matrix shows every pairwise combination of columns in the data matrix, thus visualizing all 2D correlations. (e) An exemplary tailored visualization of multivariate data, demonstrating many best practices. The return journey, shown here in black, uses the same data in panels *a–d*. Only the most visually expressive parts of the scatterplot matrix in panel *d* (dotted line) have been included. Note the use of very effective visual encodings (Figure 2), minimal use of color, and the addition of supporting context via geographic features. Metadata establish credibility by identifying the author, evidence sources, and methods. Panel *b* was made with Microsoft Excel, panels *c* and *d* were made with Matplotlib (115), and panel *e* was redrawn from Charles Minard's original using Illustrator.



genome wide. For example, single-molecule bisulphite sequencing can determine the methylation state of every cytosine base, chromatin immunoprecipitation sequencing can map protein–DNA interactions, and transposase-accessibility can be used to identify open regions of chromatin (29). Together, this flood of data contains unprecedented detail into the molecular structure, regulation, and function of whole organisms, but in a condensed, fragmented, and encoded form.

Visualization is widely used to help unravel this information: A very common task, and often rate-limiting step, is manual browsing of features to gain insight into function (20). The linear organization of chromosomes (**Figure 4a**) provides a natural visual layout, allowing many genomic features to be positioned on a common horizontal scale (**Figure 4b**), thus using the most effective visual channel (**Figure 2**).

**Visualization challenges.** A core challenge is multiscale navigation, both horizontally (across hundreds of millions of base pairs) and vertically (since regions can contain thousands of overlapping genomic features). Current genome browsers (20, 30) address this fairly well, using the general strategy of overview first, details upon demand (overview/details; 31). This strategy is implemented using feature clustering methods (e.g., ChromHMM; 32) and user interface controls to help users find and explore specific genomic regions and features of interest while maintaining awareness of overall chromosomal location and context (**Figure 4b**).

As new genomic technologies (e.g., single-cell or single-molecule DNA sequencing) continue to produce data of rapidly increasing volume and complexity, further innovations are needed in visualization methods. This includes improvements in multiscale navigation, error and uncertainty visualization (e.g., arising from base calling, assembly, and finishing; 20), variant analysis (33), and managing de novo assemblies for organisms where reference genomes are not available. In addition, novel tailored visualizations need to be developed to address a wide range of important, yet very specific biomedical topics, such as genomic rearrangements in cancer (34).

**Emerging frontiers.** An exciting frontier in genomics is the study of the 3D spatial organization of chromosomes. An accurate, atomic-scale model of the genome is a grand challenge that may someday be achieved thanks to recently developed experimental techniques (primarily Hi-C; 35) that can determine spatial chromatin contacts between pairs of genomic regions. These methods have low resolution and high false positive rates so cannot yet determine accurate 3D models for chromosomes (36). Nonetheless, Hi-C data can give new insights, but interpreting these data sets is difficult (37). Thus, tailored visualization methods are being developed, currently based around three alternative views. In one, Hi-C data are shown as a contact matrix (**Figure 4c**), allowing for high data density and a clear overview (38) but making it difficult to overlay other genomic features. A pyramidal layout (**Figure 4d**) addresses this issue (39) but makes it harder to see contacting regions. A circular layout (**Figure 4e**) is more compact (40), and using arcs to show contacts is a more effective visual encoding (41), but this does not allow the same data density; hence, only major contacts (calculated via clustering) are shown. Such tradeoffs between alternative views of multidimensional data are common, with each viewpoint providing different insights.

These and other advances in genomics are gradually unraveling the remaining mysteries of genomic function. Currently, however, we still lack an understanding of many core processes, such as exactly what gets transcribed and when and how this is controlled in different cell types.

## RNA Biology

RNA molecules play a leading role in biological systems, acting as messengers and sensors and forming the ribosome, one of the most ancient molecular machines. Many researchers now focus

---

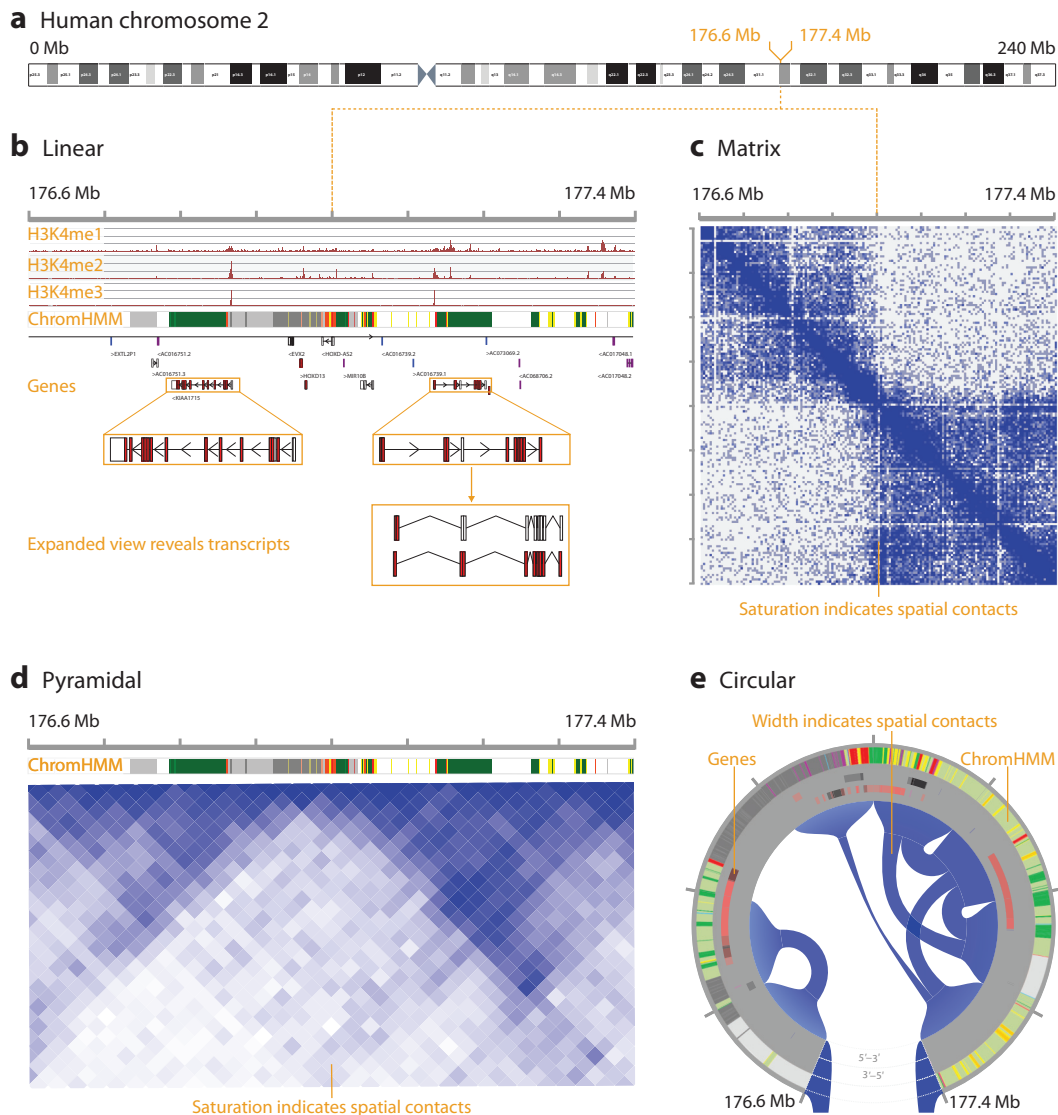
### Overview/details:

“overview first, zoom and filter, details on demand” is the visual information-seeking mantra for large data sets

### Alternative views:

different ways of visualizing the same multidimensional data set, each of which can provide different insights

---



**Figure 4**

Genomic features of human chromosome 2. (a) The linear organization of chromosomes provides a natural visual layout for mapping genomic features such as cytobands, shown in this overview of a 240 megabase (Mb) chromosome. (b) Genome browsers enable navigation to small, specific regions that often contain vast numbers of features, including epigenetic marks (H3K4me1, etc.), genes, and regulatory elements. Graphical overviews for features are created with clustering methods such as ChromHMM (32), which condense many features into a single track, using color to indicate regions with similar features. Genes are also used as a graphical overview for the often large number of transcripts they encode, which can be revealed upon demand. (c) Some features do not fit a linear layout; for example, Hi-C data (35), shown here, indicate three-dimensional spatial contacts between genomic regions and can be encoded with color saturation and a contact matrix layout. (d) Rotating the matrix and removing redundant contact data allows easier comparison with other features. (e) Connecting arcs are a more effective visual encoding for spatial contacts, and a circular layout is generally more compact. Panels *a–e* were made using Ensembl (123), Biodalliance (<https://www.biodalliance.org>; 124), JuiceBox (38), WashU EpiGenome Browser (39), and Rondo (<http://rondo.ws>; 41), respectively, and modified in Illustrator. Data in panels *c–e* are from Reference 125.

## OPTICAL ILLUSIONS CAUSED BY GROUND SUBTRACTION

Visualizations that rely on color to encode quantitative values are subject to an optical illusion known as ground subtraction. In a heat map, for example, strongly contrasting colors in a cell's neighbors can make the cell appear much higher or lower in luminance than it should. This illusion can be very strong; as demonstrated in the checker shadow image by Edward H. Adelson (<http://persci.mit.edu/gallery/checkershadow>), the human visual system can be surprisingly inaccurate at reading values encoded with color. In heat maps, this illusion can mask true patterns and introduce visual artifacts (47). As the numbers of rows and columns increase or as cell size is reduced, the effect can become worse, making it impractical to display all significant results as one very large heat map. Thus, for the display of quantities where absolute variation between observations is important, it is recommended to encode values with position or size rather than lightness, saturation, or color hue (Figure 2).

on unlocking the secrets of the RNA world; meanwhile, the measurement of RNA transcript abundances has become the workhorse of modern biology. First accomplished with microarray experiments (42), accurate measurements of the abundance of transcripts in biological samples and single cells (43) are now taken with RNA sequencing (RNA-seq) (44).

**Visualization challenges.** Interpreting the high-dimensional data sets from RNA-seq experiments remains challenging. After careful experimental design and statistical analysis (45), gene expression values judged to be significant are commonly presented as clustered heat maps (46), a technique that has dominated since the first microarray experiments (42). However, optical illusions in these visualizations make it difficult to judge the magnitude of individual values or fold changes between pairs of values (see the sidebar titled Optical Illusions Caused by Ground Subtraction). As the numbers of rows and columns increase or cell size is reduced, these effects become worse, making it impractical to display all significant results as one very large heat map. Further problems arise because the rows and columns of clustered heat maps are usually ordered to group associated genes and conditions and so highlight regulatory effects (47). Inevitably, values for genes and conditions without significant association will be placed next to one another, which exacerbates perceptual problems. Separating unrelated rows and columns (Figure 5a) can help, but this does not fundamentally address these difficulties (48), particularly for genes that cluster poorly. In such cases, there may be insufficient data to resolve those genes' regulatory networks as a 1D ordering, so it is important that the degree and support for relationships inferred from clustering are also shown. The addition of tree graphs, however, further constrains the size of the heat map that can be displayed without issue,<sup>1</sup> so we suggest that only the most informative subset of genes and conditions should be presented in this way.

**Emerging frontiers.** Single-cell RNA-seq experiments (scRNA-seq) are a revolutionary new technology that can reveal key events in differentiation normally masked in bulk RNA-seq experiments, thus providing deep insight into the behavior of cells and tissues. These data are typically visualized using dimensionality reduction methods that allow gene expression vectors to be projected onto 2D scatterplots (Figure 5b–d). scRNA-seq data allow the sequence of these events to be reconstructed—commonly termed cell pseudotime (49). Clustering and dimensionality

---

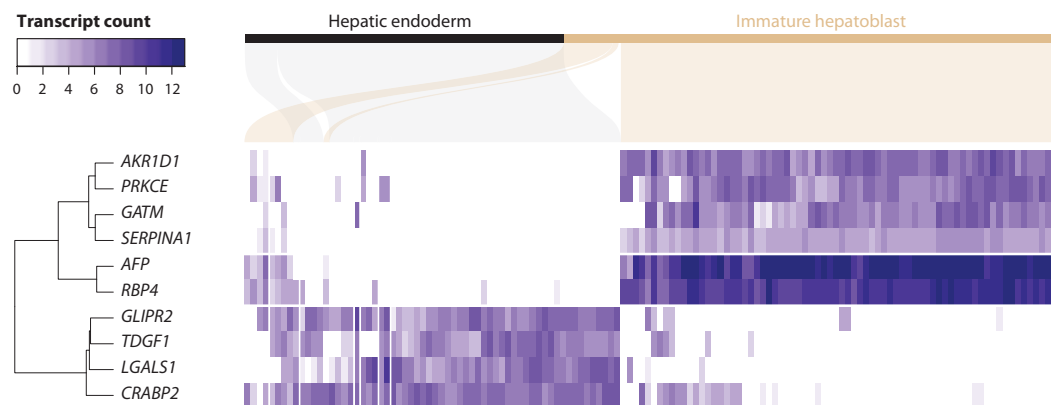
**Heat map:** a graphical representation of a matrix of data where individual values are encoded using color

**Tree graph:** a graph where all lines connect without forming loops; used for hierarchical data

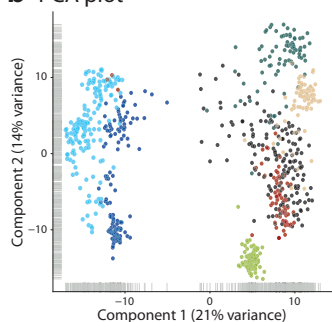
---

<sup>1</sup>We recommend rectangular heat map cells of no less than 6 mm, separated by 1.5 mm, and overlaid on a neutral background (white, black, or a color that does not contrast with those employed in the heat map).

## a Sankey diagram, tree graph, and heat map

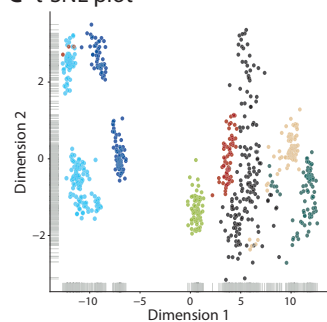


## b PCA plot

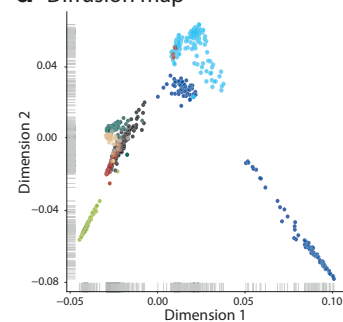


• Definitive endoderm • Endothelial • Hepatic endoderm • Immature hepatoblast • IPSC • Mature hepatocyte • Mesenchymal stem cell

## c t-SNE plot



## d Diffusion map



## Figure 5

Visualizations of single-cell RNA sequencing observations of liver bud development. (a) Clustered heat map for the top 10 differentially expressed genes in two cell types, indicated in the Sankey diagram with black and beige coloring. Saturation encodes absolute expression, and row and column positions encode genes and experimental conditions, respectively. Genes and cells with similar expression patterns are clustered to optimally order rows and columns. The cluster tree graph shows three distinct groups of expression behavior, and vertical space has been inserted to separate these sets of rows. The Sankey diagram highlights imperfect separation of the two cell types, and spaces have also been inserted to separate sets of differently behaving cells. In panels *b–d*, scatterplots show alternative views created by applying dimensionality reduction methods, each revealing different aspects of the full data set. (b) PCA groups most cell types but does not resolve cells forming the definitive endoderm and the hepatic endoderm. (c) t-SNE (126) provides more insight, revealing local similarities as well as overall variation in the data set. However, t-SNE can be more difficult to apply, as it requires setting a manually adjustable parameter (perplexity) (127). (d) Diffusion maps (128, 129) model relationships between points in the data set as a diffusion process that is then reduced to a lower-dimensional map. Here, successive developmental relationships between cells are revealed. Panel *a* was made using R and D3.js (Sankey diagram). Panels *b–d* were made using *scatter* (130). All panels were modified using Illustrator. Data in panels *a–d* are from Reference 131 and reanalyzed in R with read counts processed, as described by Hemberg et al. (<https://github.com/hemberg-lab/scRNA.seq.course>). Abbreviations: IPSC, induced pluripotent stem cell; PCA, principal component analysis; t-SNE, t-distributed stochastic neighborhood embedding.

reduction heuristics (49) allow pseudotime to be inferred, visualized, and quantitatively analyzed (Figure 5d).

Single-cell transcriptomics measurements will soon become possible at the whole-organism level (50), and we will undoubtedly require more effective methods for interpreting these data. However, the measures of abundance obtained from these experiments are only markers that

## AN EXEMPLARY BIOMEDICAL DATABANK

Compared with many areas of biomedical science, visualization methods for macromolecular structures are more advanced, largely because they build upon a solid bedrock of exceptionally well managed data. Created in 1972, the Protein Data Bank (PDB; 51) has exemplary practices and stability that facilitate reproducibility and substantially simplify the difficult task of creating and maintaining tailored visualization tools. In particular, (a) each entry is a deposition related to one specific scientific publication, not to an abstract concept (e.g., a gene or pathway) whose definition may change over time. (b) Entry identifiers are short and designed to be easy to remember. Depositors can propose an identifier, ensuring that many are meaningful. For example, the first crystal structure of the protein actin has identifier “1ATN” (114). (c) Entries include rich metadata describing how they were generated (and by whom), as well as cross-linking to related databases. (d) Raw and processed data are stored, enabling later reanalyses. (e) An international network of organizations maintains and curates the database. (f) PDB deposition is required when publishing in major journals.

Unfortunately, many biomedical databanks created since have not learned from these practices, thus requiring tool developers to contend with data formats and sources that are many, varied, and often unstable.

indicate which parts of an organism’s genome are active. To understand the biological role each gene plays, we must look beyond sequencing data; in fact, much of our current understanding has come from studying the molecular structure of RNA transcripts and the proteins they encode.

### Protein Structures

Protein structural biology aims to provide a detailed understanding of life’s molecular machinery. Thanks to decades of research worldwide, we now have 3D molecular structures (at or near atomic resolution) for ~40,000 proteins. By viewing these structures, researchers can gain insight into precise molecular mechanisms underlying many of the biochemical processes occurring within living cells. Remarkably, almost all these structures are collected in a single, exemplary database (<https://www.rcsb.org/>; 51) (see the sidebar titled An Exemplary Biomedical Databank). This has helped drive innovation in molecular graphics, which has outpaced visualization advances in many other areas of biomedical science (22).

**Visualization challenges.** Visualization is integral to structure determination and validation (22), as well as for gaining insight into protein function (e.g., with tools such as Chimera and others; 22, 52, 53). A core challenge is conveying the many different features of these large, complex data sets; this requires careful use of visualization principles (e.g., overview/details), judicious and minimal use of color, and visually expressive representations to highlight specific aspects of the data (**Figure 6a–b**).

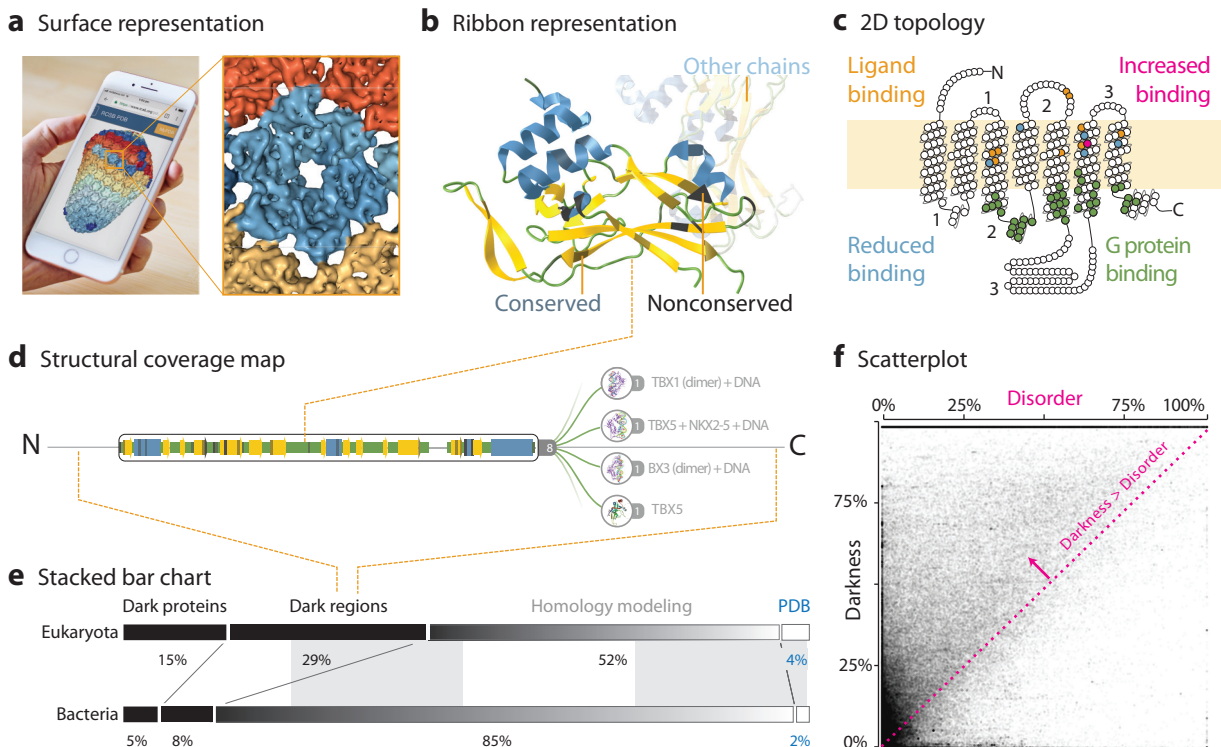
Another challenge is conveying the complex 3D shape of proteins. In special cases, shape can be communicated through specifically tailored 2D visualizations (**Figure 6c**), but ultimately, protein structures need to be viewed in three dimensions. For this reason, structural biology has been an early adopter of new visual techniques, starting with physical models (used in solving the first protein structures; 54), stereoscopic imaging (55), and virtual reality (VR) (56).<sup>2</sup> This has continued, with techniques such as low-cost VR (e.g., Visual Molecular Dynamics supports Oculus Rift; 57), very-low-cost VR (e.g., Autodesk Molecule Viewer supports Google Cardboard;

---

**Virtual reality (VR):** blocking a person’s view of their surroundings via head-mounted displays (e.g., Oculus), allowing immersion in virtually generated content

---

<sup>2</sup>VR can focus undivided attention on a data set. Although powerful, usage is limited by inconvenience, discomfort, motion sickness, and other drawbacks. By contrast, augmented reality has fewer drawbacks and looks likely to become widespread in biomedical research and in normal life.



**Figure 6**

Visualizations of protein structure data. (a) Advances in web molecular graphics now allow structures with millions of atoms to be interactively explored on a smartphone. Surface representations (shown here) and space-fill representations are useful for overviewing the arrangement of individual proteins in larger assemblies. (b) For a detailed view of a single protein, ribbon representation is useful (132), revealing how polypeptide chains fold in three dimensions; this is helped by linking to a sequence view (d). Using semitransparency for other chains in the structure can provide supporting context without clutter. Highlighting conserved or nonconserved amino acid differences to the wild-type sequence of interest gives a visual indication of model reliability (133). Typically, many such differences occur in structures inferred via homology modeling; but they are also common in Protein Data Bank (PDB) structures due to experimental limitations. (c) In special cases (e.g., GPCR transmembrane proteins, shown here), simplified two-dimensional (2D) schematics can be used to show overall topology as well as details, such as loop regions or residues where mutations have large functional effects (encoded using numbering and coloring, respectively). (d) A schematic representation of a full-length, wild-type protein sequence, with coloring indicating regions with significant sequence similarity to structures in the PDB. Details on these matching structures can be revealed upon demand using a tree graph. On average, each protein sequence matches to ~200 PDB structures (63) but contains several dark regions, with no detectable similarity to any known 3D structure (66). (e) A stacked bar chart showing the total fraction of protein residues that map to any PDB structure (either directly or via homology modelling); the remaining fraction (dark proteome) is divided into dark regions (d) and dark proteins (where a single dark region spans an entire sequence) (66). Stacked bar charts have been used to achieve moderate data density, and axes replaced by two shaded regions (indicating 25%, 50%, and 75%). Connecting lines facilitate comparison. (f) A scatterplot of darkness versus disorder (134) for ~180,000 eukaryotic proteins. Point size, color, and transparency have been adjusted to reveal an unexpected overall pattern (darkness exceeds disorder for most proteins, indicating that much of the dark proteome is not explained by disorder). Due to high data density, subtle patterns are also revealed (e.g., horizontal streaks arising from related sequence families). Panel a was made using NGL Viewer (<https://www.rcsb.org/pdb/ngl/ngl.do?pdbid=3J3Q>; 135) with PDB 3J3Q (136) and a stock photo from Unsplash, panels b and d were made using Aquaria (<http://aquaria.ws/O75333/4a04/>, 63) and Photoshop with PDB 4A04 (137), panel c was made using data from GPCRdb (138), and panels e and f were made using ggplot2 (116) with data from Reference 139. All images were modified using Illustrator.

<https://www.molviewer.com/>), 3D printing (58), commodity interaction devices (e.g., Leap Motion and Kinect; 59), augmented reality (AR) (60), crowdsourced evaluations (61), concepts from computer gaming (62), and emerging web technologies (e.g., WebGL; see **Figure 6a**).

**Emerging frontiers.** Protein structural biology is still far from complete, as many proteins still have little or no experimentally determined structural information. To address this, researchers are using high-throughput homology modeling to systematically compare all known protein sequences against all experimentally determined structures, resulting in over 100 million model structures (63). Allowing researchers to effectively explore and benefit from such large data sets requires carefully tailored visualization tools (63) that use the overview/details strategy, as well as alternative views connected via brushing and linking (**Figures 6d,e**) (64). Homology modeling currently provides structural models for about half of the eukaryotic proteome (**Figure 6e**). Interestingly, much of the remaining dark proteome currently cannot be explained (e.g., **Figure 6f**). Exploring this dark protein structure universe is an important data science challenge (65) in which visualization is playing a key role (66).

High-throughput approaches are also being applied to molecular dynamics (67), generating increasingly large, complex trajectories; these data can give insights into key events (e.g., binding with ligands or other proteins). However, unearthing those insights is a still major challenge, requiring further innovations to create very specific, tailored visualization tools (e.g., 22, 52, 57).

In addition, rapid advances in cryo-electron microscopy are making accessible much larger structures and molecular assemblies than ever before (68); this has promoted improvements in methods for visual exploration of multiscale molecular data (<http://ncbr.muni.cz/LiteMol>). Finally, high-throughput computing is also being used to integrate structural data in the construction of atomic-scale models of viruses, subcellular compartments, or even whole cells (69). The scale and complexity of these models requires the development of radically new visualization methods, bridging structural and systems biology (70).

## Systems Biology

We have long speculated about how biomolecules coordinate to perform cellular function (71), and graph-based visualizations have been key to organizing our thoughts (23). An exemplar is the Roche metabolic pathway (**Figure 7a**), initiated by Gerhard Michal in 1965 (<http://biochemical-pathways.com/>; 72); this manually tailored visualization shows thousands of metabolic reactions in a single, comprehensive view. Such pathway graphs have endured because they are visually expressive, showing causal flow and providing insight into molecular events underlying health and disease.

**Visualization challenges.** Over 4 billion biochemical reactions are currently known, and this number is rising rapidly (73). These data are typically visualized with specialized tools [e.g., Cytoscape (74) or Gephi (75)] that provide a range of automated layout methods, many based on force-directed algorithms (76), resulting in network graphs (**Figure 7b**). A force-directed layout (also known as spring embedding) can be useful for overviewing a data set; however, even small biological networks are often so interconnected that these graphs become overly cluttered (**Figure 7b**). Force-directed layout is so common that it has become something of a limiting paradigm, often used even when better strategies are available.<sup>3</sup> For example, when integrating connectivity with other data (e.g., time, subcellular location), the go-to strategy has been to overlay

---

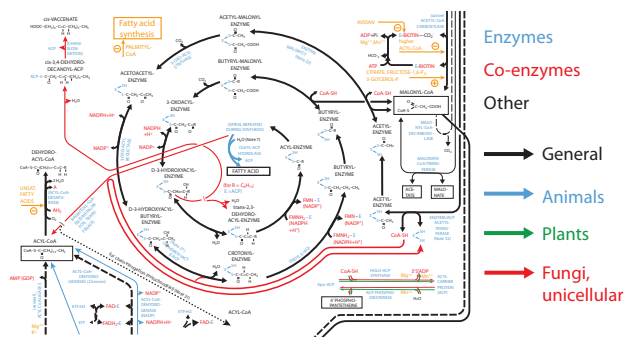
**Augmented reality (AR):** augmenting a person's normal view of their surroundings by adding computer-generated images or data (e.g., HoloLens)

**Brushing and linking:** linked alternative views, where interactive changes made in one view are automatically reflected in the other

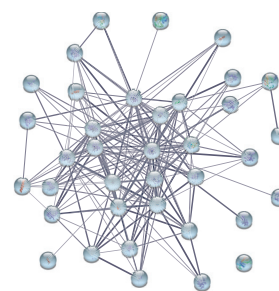
---

<sup>3</sup>Overreliance on familiar tools can lead to cognitive bias: Paraphrasing a popular adage, if your only tool is spring embedding, every data set looks like a network graph.

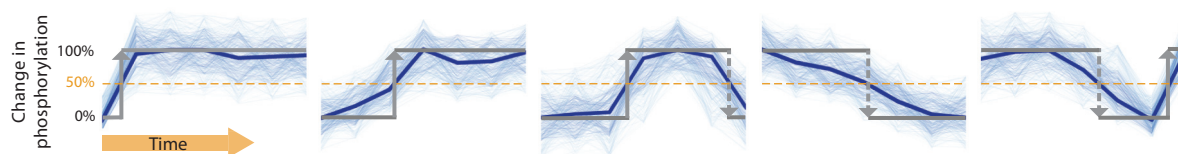
### a Roche metabolic pathway graph



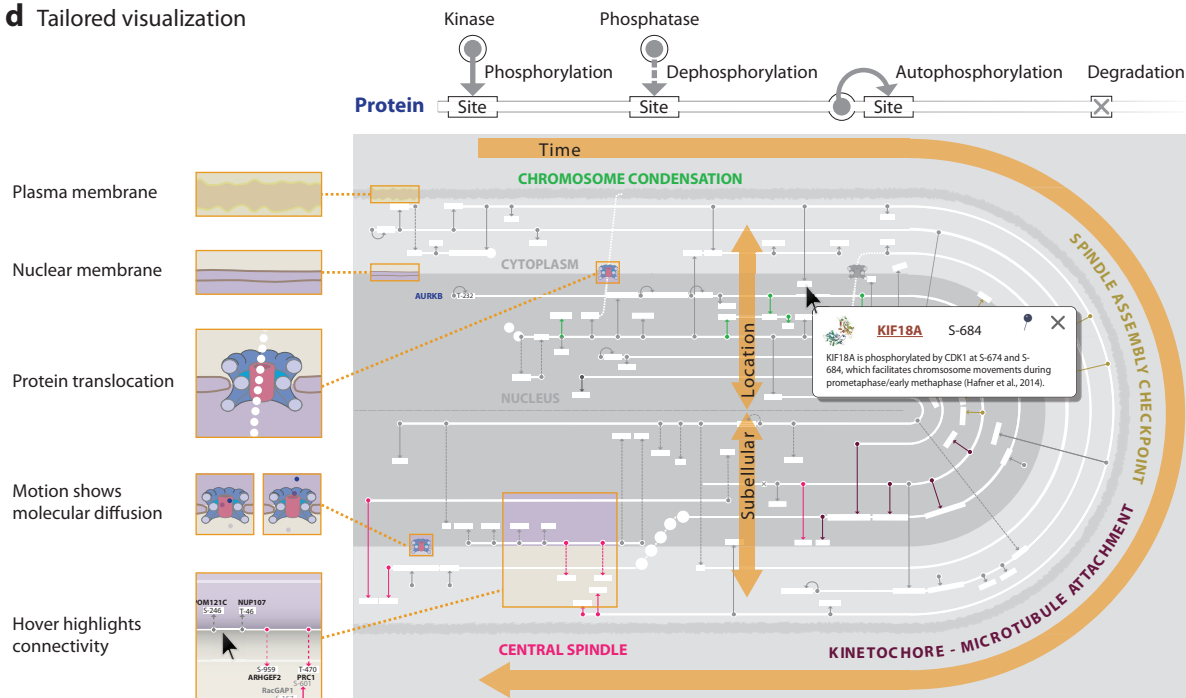
### b Network graph



### c Profile plots



### d Tailored visualization



(Caption appears on following page)



**Figure 7** (Figure appears on preceding page)

Pathway and network graphs of molecular systems. (a) Part of an exemplary tailored visualization; the full pathway shows causal flow involving ~3,000 reactions, plus supporting context (e.g., molecular structures). Note the effective use of visual channels (Figure 2). Position and shape show reaction categories, while minimal coloring is used to show—without clutter—different versions of the pathway for four categories of organisms. (b) Network graphs created via spring embedding are common (here, edge width encodes interaction confidence scores) but often too cluttered. This can be partly mitigated via, e.g., edge bundling (78). (c) The first step in visualizing phosphoproteomics data is to identify clusters of phosphosites with similar time profiles (*thin light-blue lines*). Clusters are modeled as a series of phosphorylation (*solid arrows*) and dephosphorylation (*dashed arrows*) events, each arising from a specific kinase or phosphatase (140) and occurring when the average phosphorylation (*thick blue line*) passes 50% (*yellow dotted line*). (d) A tailored visualization for phosphoproteomics data, where selected proteins are drawn as tracks in a “circular” cellular landscape. Position encodes subcellular location (with translocations shown as excursions from the track layout) and the temporal ordering of phosphoevents, each indicated via an arrow connecting a kinase or phosphatase to its substrate site. Color hue shows events that perform coordinated functions. Texture is used to show context, such as membranes. An online version has further interactive features (e.g., informative popups, motion, highlighting upon hover) that help researchers use these complex data sets to gain insight into cellular processes, such as insulin response (85) or mitosis (86). Panel a was redrawn with permission from <http://biochemical-pathways.com/> (72); panels b–d were made using STRING (141), Matplotlib (115), and Minardo (<https://minardo.org>; 86), respectively, and modified using Illustrator.

these data onto existing network layouts (23), thus reducing visual effectiveness due to clutter. Often, it is better to change the layout entirely, using position (the most effective visual channel for quantitative data; see Figure 2) to encode not just connectivity but also biological context that helps in interpreting data. This approach is taken in several tailored visualization tools, such as Cerebral (77), which uses position to encode subcellular location, edge-bundling to reduce clutter (78), and small multiples for different conditions (79).

Visually expressive layouts (e.g., Cerebral) need to be tailored for each specific scenario, and there are vastly many distinct scenarios in systems biology, since reactions vary greatly with cell type, timing, and molecular microenvironment (80). Thus, many tailored visualizations have already been developed (23, 30), and unfortunately, systems biology data are fragmented across ~700 resources (<http://pathguide.org>; 73), which partly impedes progress in the field. To help manage the many complexities of these data, researchers use a wide range of visual techniques; for example, the overview/detail strategy is used to allow subgraphs to be collapsed or expanded upon demand, thus helping users more effectively explore large graphs. However, there remains considerable scope for using visualization principles to design new, automated layouts (81, 82) and to improve the computer-aided design of manual layouts. There is also considerable scope for improving how systems biology data are organized (see the sidebar titled An Exemplary Biomedical Databank).

**Emerging frontiers.** Mass spectrometry–based proteomics (83, 84) is a rapidly emerging technology that enables systematic measurement of proteome-wide posttranslational modifications in response to stimuli, such as phosphorylation (termed phosphoproteomics). These technologies are providing new insights into fundamental cellular processes and diseases, which in turn may lead to new therapeutic interventions. However, there is a price: increased complexity. Each phosphoproteomics experiment can track highly dynamic changes in over 10,000 different phosphosites in over 5,000 proteins (84). In analyzing these data sets, it is common to first identify coregulated phosphosites, based on clustering of time profiles (Figure 7c). Several tailored visualization strategies are being developed for exploring these clusters. In one, inspired by the cyclic journey in Minard’s exemplary chart (Figure 3e), the cascade of phosphorylation events is laid out as a cyclic journey through a cellular landscape (Figure 7d). Proteins are represented as tracks and phosphoevents are positioned by time and subcellular location—two key variables in these experiments. This layout facilitates spatial reasoning about causal relationships, helping researchers use these complex data sets to gain insight into cellular processes, such as insulin response (85) or mitosis (86).

These and other advances in molecular systems biology promise to revolutionize medicine. However, realizing this promise will require software platforms capable of bridging scales, from molecules to cells, tissues, and whole organisms—a formidable challenge in which visualization plays a central role (70).

## Cellular and Tissue Imaging

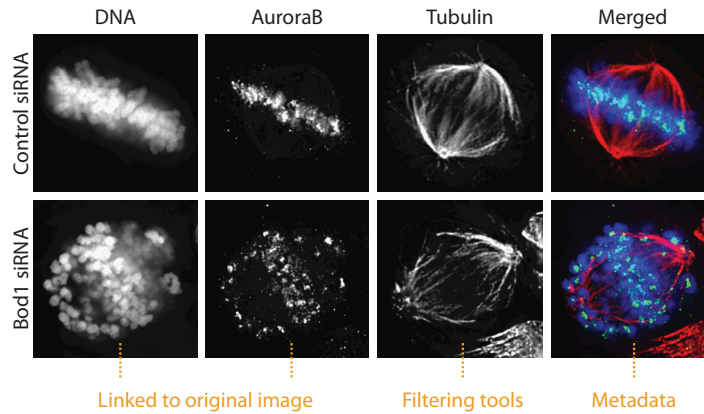
Imaging remains the primary way that we observe biological systems. Quantitative imaging data are employed throughout the biomedical sciences as a basis for research, diagnosis, and therapy. Since van Leuwenhoek created the first microscope, biologists have observed the structure and behavior of cells and how they form tissues and develop into organisms. Most recently, technological advances in labeling, sample handling, and imaging have expanded microscopy's capabilities (87), and we are now able to image complex cellular assemblies, such as neurons, in three dimensions or capture in real time the cellular processes that drive development (88). Imaging has also advanced in medicine, allowing detection and diagnosis of pathologies and providing essential guidance before and during surgery.

**Visualization challenges.** Imaging data can give quantitative and qualitative insights into morphology and function; however, this often requires extensive and sophisticated processing pipelines (89). An increasing array of tools are being developed to address this need, some of which now allow interactive visualization of raw images together with image-derived data, e.g., the Medical Imaging Interactive Toolkit (<http://mitk.org/>) and Slicer (<http://www.slicer.org/>). The primary way that we interpret these image-derived data is by encoding them as colors and annotations that are then overlaid onto the original images. This presents a challenge because biological images are usually already very complex; thus, detail often needs to be removed from the original image (for example, by transforming the image's dynamic range or color space) before derived data can be overlaid for either data exploration or publication. While many image processing and figure generation tools allow these operations, such manipulations may obscure critical details and therefore need to be documented in a reproducible manner. There is broad consensus that the transformation of image data needs to be better reported in scientific publications (90). Recent advances (91) are beginning to address these issues, allowing interactive creation of figures for online publication (**Figure 8a**) that link to original data (e.g., high-content screening, time-lapse or histological whole-slide imaging data) and to metadata (related to experimental design, image acquisition, analysis, and interpretation), thus allowing subsequent reanalysis.

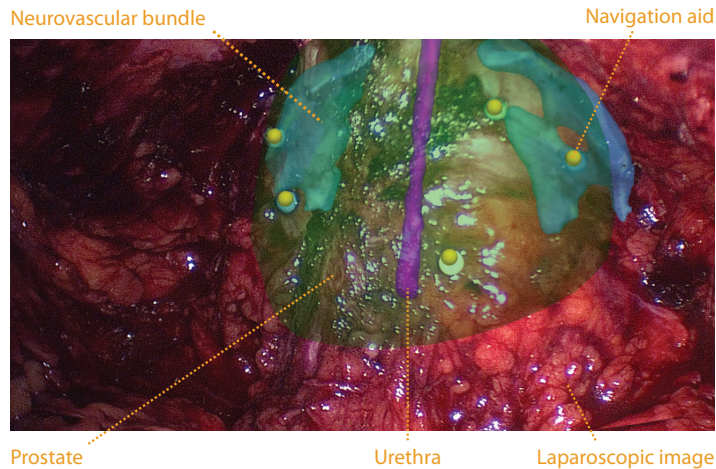
**Emerging frontiers.** An exciting frontier in medicine is the use of AR (92) to enhance, for example, live video feeds used to guide surgery with preoperative, diagnostic imaging data (**Figure 8b**). This approach promises to reduce error and increase precision during surgery by providing real-time guidance about the location and physiological status of diseased tissue. Once again, visual complexity is a major challenge; however, by integrating machine learning and semantic modeling approaches (93), surgeons today have already begun to use AR to see beyond the capabilities of normal visual perception.

There remains tremendous opportunity for improving the integration of advanced analytic tools with interactive visualization, creating new platforms, such as the Allen Cell Explorer (<http://www.allencell.org/>). Soon such improvements are set to greatly advance our understanding of the composition, structure, and dynamics of normal and pathological cells and tissues, as well as the effectiveness and precision of medical interventions.

### a Automated publication of online figures



### b Augmented reality imaging



**Figure 8**

Cellular and tissue imaging. (a) Multipart, annotated image created from, and linked to, raw data sets. By automating many routine manual tasks involved in creating well-formatted, publication-ready figures, tools such as OMERO.figure enable scientists to scale-up, easily creating figures with higher data density, and thus to address more complex questions. Figure panels can be rendered dynamically from the original image data and automatically overlaid with timestamps and scale bars, avoiding potential human error. Such tools can document all steps from the original image files to the final figure, improving data integrity, organization, and provenance. (b) Augmented reality imaging in minimally invasive surgery. Before the intervention, target and critical structures are segmented in a three-dimensional (3D) planning image. At the beginning of surgery, artificial navigation aids (fiducials) are inserted into the target organ (here, prostate), and their 3D configuration is determined from a 3D intraoperative medical image (e.g., a 3D transrectal ultrasound image). The latter is fused with the preoperative planning image using a 3D/3D registration algorithm. During surgery, the fiducials are continuously tracked in the 2D video images, and a 2D/3D registration algorithm (142) is used to find a transformation relating the endoscopic camera coordinate system with the image coordinate system of the 3D intraoperative modality. This enables the laparoscopic video image acquired during prostatectomy to be overlaid in real time with the prostate capsule and critical structures. Panel *a* was made using OMERO.figure (<http://figure.openmicroscopy.org/demo/#file/1>; 91) with data from Reference 143. Panel *b* was made using the Medical Imaging Interaction Tool Kit (144), with data from Reference 145.

---

#### Parallel coordinate plot:

a profile plot of multidimensional points, each shown as a series of line segments connecting parallel axes

#### Stacked bar chart:

a visualization in which bars representing related data are stacked on top of (or beside) each other

#### Flame graph:

a visualization of hierarchical data where width encodes branch quantity, and sub-branches are stacked on parent branches

#### Linear diagram:

shows the size of overlaps and differences among multiple sets of data; an alternative to Venn diagrams

---

## Populations and Ecosystems

Some of the most compelling questions in biology center on how our genome affects how we live and interact with other organisms and our environment and how these interactions change over time. Tree graphs (**Figure 9a**) are the primary way we visualize ancestral relationships between organisms. With sufficient data, trees can convey not only evolutionary distance but also the order in which different lineages may have evolved. Genomic sequence comparisons allow us to infer how species or individuals within a population differ from one another, but these contain localized features (e.g., SNPs, indels), copy number variants, and rearrangements that extend over millions of bases. Parallel coordinate plots (**Figure 9b**) provide one way of viewing multiple alignments of closely related bacterial genomes. These visualizations highlight the changes in genomic regions or even individual genes as they undergo mutation, rearrangements, or horizontal gene transfer.

**Visualization challenges.** Phylogenetic tree and comparative genomic visualizations are relatively mature (20, 30). For metagenomics and population sequencing, however, neither representation will suffice. Here, interactive plots or static visualizations at multiple scales are needed to capture the breadth and depth of these data. Branch structures in phylogenetic trees quickly become illegible in the presence of large sets of taxa (**Figure 9a**), so important differences in lineage must be manually highlighted. In biome analysis, the main objective is to determine the composition of samples (**Figure 9c**) and how they change and evolve over time. Stacked bar charts (**Figure 9c**) allow broad differences in composition to be shown, but it is important to also show lineage. This is often done with sunburst plots (**Figure 9d**), although using flame graphs instead makes it easier to compare multiple plots (**Figure 9e**). Similarly, the Venn diagrams (**Figure 9f**) typically used to visualize species co-occurrence can often be replaced by linear diagrams (**Figure 9g**). Overall, there remains considerable scope for improvements—for example, with many of the current tools in this research area, considerable effort is required when using them to create figures for publication that are uncluttered, effective, and visually expressive.

**Emerging frontiers.** Advances in genomic sequencing allow us to examine differences within populations and across ecosystems in unprecedented detail. A plethora of microbiome sequence data promises to revolutionize our understanding of evolution and human health. But we are struggling to develop effective visual analysis strategies because no single visual metaphor captures the richness of population-level data. Pan-genome visualizations, designed to show core and accessory genes in a species' genome, are dashboards that combine existing methods (94). Phylogeography uses maps to show phylogenetic relationships across geocoded samples (95). Pathogen surveillance, a crucial challenge that spans these fields, requires integrated, alternative views that capture population dynamics and highlight emerging resistance (96).

Our ability to observe biology at the molecular, cellular, anatomical, and physiological levels has never been greater, but making sense of these emerging data will require overcoming formidable analytical challenges, as well as the invention of fundamentally new visual metaphors (97), changing how we see and think about our data.

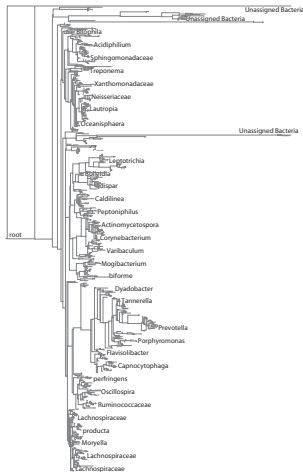
## VISUALIZATION FOR COMMUNICATION

Science is not complete until it is communicated (**Figure 1a**); however, this is often challenging due to the inherent complexity of biomedical research. Fortunately, visualization can help here as well.

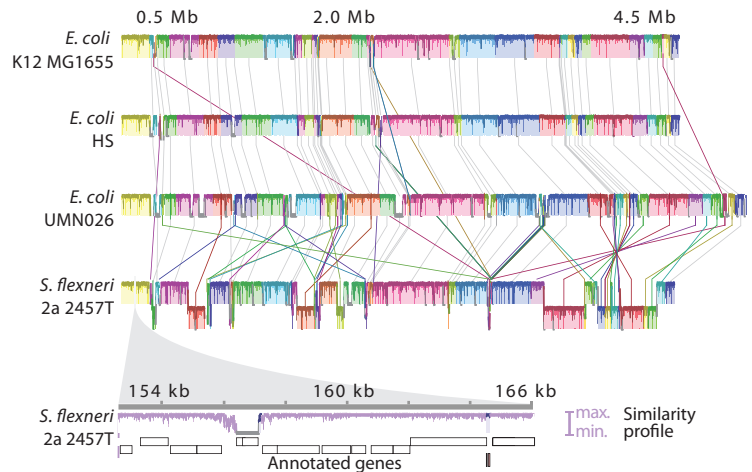
## Figures and Illustrations

In preparing a publication, visualization tools used for discovery are typically used to select static views that best express the insights found. A very small number of journals allow interactive figures; unfortunately, interactive figures in publications are often complex to produce and difficult to maintain – and, ironically, are used by very few readers. Interactive figures can augment static figures, but not replace them: just as scientific writing commits us to a particular way of describing our work, a static view commits us to a particular viewpoint that we believe best expresses the

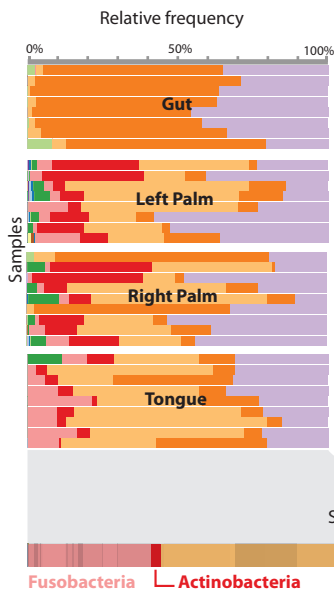
### a Phylogenetic tree



### b Parallel coordinate plot



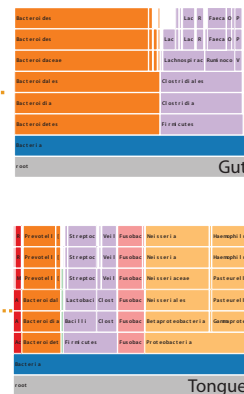
### c Stacked bar charts



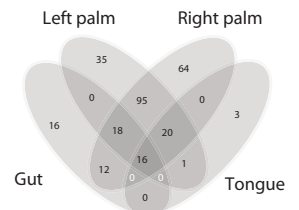
### d Sunburst plots



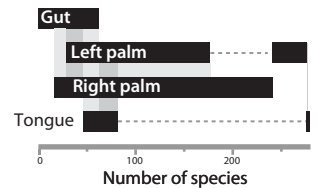
### e Flame graphs



### f Venn diagram



### g Linear diagram



(Caption appears on following page)

**Figure 9** (Figure appears on preceding page)

Phylogenetic, comparative genomic, and metagenomic data visualizations. (a) A phylogenetic tree showing evolutionary relationships and inferred operational taxonomic units (OTUs) for 16S amplicon sequencing data. While commonly used, phylogenetic trees have drawbacks: Closely related taxa become hard to resolve, and as their numbers increase, topological relationships quickly become obscured, even with the application of semantic zooming (used here to omit overlapping labels on adjacent branches). The tree contains a total of 761 leaves of which ~50 have a named OTU. (b) A multiple-genome alignment viewed as a parallel coordinate plot, with a zoomed region on a *Shigella flexneri* genome. Each genome is represented as a linear axis. Connecting lines between genomes indicate conserved regions. Lines that do not intersect with some genomes indicate horizontal gene transfer, and convergent and divergent lines correspond to gene duplication and inversion events between *S. flexneri* and the three *Escherichia coli* strains. The zoomed region of the *S. flexneri* genome reveals regions of divergence in an otherwise conserved part of the alignment. Colors are assigned to aligned sections of each genome, and a similarity profile is overlaid as a line graph in a more saturated color. When genes are inverted in some organisms, the area is shown below the genome axis. (c) Species abundance (or beta diversity) visualized as a stacked bar chart for a metagenomics analysis of microbiome samples taken from different parts of the body. Colors encode phyla of identified OTUs in samples. These charts are useful for comparing abundance across broad taxonomic levels (as shown here) but become too complex when used to show the ~280 species in each sample (zoomed region). (d) Sunburst plots showing beta diversity for pooled samples from two body sites. While these accurately portray lineage relationships, it can be difficult to compare multiple plots. (e) Flame graphs encode taxonomic rank as height and abundance as width, making it easier to compare plots and to see where taxonomic assignment is incomplete. (f) Species co-occurrence among samples from the four sites shown as a Venn diagram. Many tools offer advanced layout and shading models for Venn diagrams but can result in plots that are not visually effective. (g) A linear diagram of the same data, using the *x*-axis to show the number of co-occurring species between tissues and gray vertical boxes to highlight intersections. Panel *a* was made using Archaeopteryx (146), *b* using Mauve (147), *c* using QIIME2 (148), *d* and *e* using QuanTiTree (<http://metasystems.riken.jp/visualization/quantitree/index.htm>), *f* using the R Venn library (<https://cran.r-project.org/web/packages/venn/index.html>), and *g* using the Linear Diagram Generator (149). All panels were modified using Illustrator. Data in panels *a* and *c-g* are from the moving pictures data set (150).

phenomena revealed by our data. Although difficult, selecting static views is often an essential step in research – and can lead to new insights.

In many cases, these static views need to be postprocessed, using tools such as Illustrator or Photoshop, to improve clarity and ensure that marks and labels are consistent and readable at publication scale. In addition, since ~5% of readers and reviewers are colorblind (98), it is good practice to use color blindness proofing tools<sup>4</sup> and, where needed, modify figures to avoid relying on red-green contrast – this can usually be achieved by adjusting saturation and lightness values to increase contrast (98).

However, it sometimes can be unclear where the boundary lies between necessary improvements versus scientific fraud; thus, it is important to follow established guidelines on image and figure manipulation (99).

## Animations and Videos

Animations and videos can dramatically enhance scientific communication and are becoming easier and cheaper to produce, leading to a marked increase in scientific video content (100). Done well, scientific videos improve peer-to-peer communication and inspire public engagement and enthusiasm (27). Unfortunately, ensuring scientific accuracy typically involves considerable time and effort, as does achieving a high standard in video production, which requires learning cinematography principles, practices, and tools (e.g., Autodesk Maya, Blender, Adobe After Effects).

To help overcome these barriers, researchers are developing several specialist tools to streamline and simplify the production of scientific animations (101). Unfortunately, until these efforts become more advanced, accurate and compelling scientific videos will likely remain relatively rare.

<sup>4</sup>In Illustrator and Photoshop, choose the menu items View > Proof Setup > Color Blindness to preview how a figure will appear to people with common forms of color blindness.

## PERSPECTIVES

This review has highlighted a few specific cases where data visualization is being used to accelerate discovery, but biomedical science has thousands more. Thus, while many visualization tools are already available (1, 20–24, 30), they are often inadequate for cutting-edge data sets. Addressing this challenge requires the invention of novel, tailored visualization strategies, each adapted to specific scenarios; this can be very difficult and, in many cases, is a rate-limiting step in discovery. The resources outlined in this review can help (particularly **Table 1**). Especially noteworthy resources include the 2010 *Nature Methods* special issue on visualizing biological data (1, 20–24) and the ongoing *Nature Methods* Points of View article series focused on specific visualization issues for life scientists (6). It can also be useful to exchange experiences with peers facing similar challenges. The annual Visualizing Biological Data (VIZBI) conference provides a forum for this exchange and provides a free online collection of videos and posters from previous meetings (<http://vizbi.org/>). The VIZBI forum is also designed to help bioinformaticians connect with graphic designers, graphic artists, and biomedical communicators using illustration or animation. It can also help to connect with computer scientists researching data visualization; as well as advising on good principles and practice, they can be valuable collaborators. A forum for such engagement is provided by the annual BioVis Symposium (<http://biovis.net/>).

Tailored visualization tools play a critical role in research, some becoming widely used and highly cited (102). However, limitations in popular tools (e.g., cluttered, overly complicated user interfaces or poorly chosen defaults, such as rainbow color maps; 10) can have very negative impacts, contributing to dead-end research and incorrect diagnoses. Creating tailored tools with good visualization and design practices (103) typically requires years of sustained focus. It is often not clear how tool development and maintenance can be funded; however, it is clear that this needs to be a central issue in research funding policies.

This review has highlighted that specific research areas need highly tailored visualization tools; nonetheless, there are common generic methods (e.g., tree graphs, parallel coordinates, stacked bar charts) and strategies (e.g., clustering, alternative views, overview/details, linked views, minimal coloring) being used across many research areas. There are also common outstanding challenges, such as uncertainty visualization (104) and multiscale navigation. Unfortunately, some poor visualization practices are also common (e.g., overly cluttered visualizations). However, there are good indicators that this situation is improving and that there is increased awareness of the importance of data visualization in the life sciences (105). This is evidenced by the increased focus on visualization in mainstream conferences, as well as the emergence of more specialist meetings, such as VIZBI and BioVis.

As we seek to improve our tailored visualizations, another common challenge is how to objectively assess the quality of a particular visualization method or tool (106). An obvious and important quality measurement is rate of adoption by the community; however, the popularity of a visualization strategy often has more to do with the cognitive load required to first learn how to read it, rather than how effectively or expressively it allows data to be understood and new insights to be generated. Here again, data visualization research can help. Methods are being developed for quantitatively evaluating the effectiveness of visualizations (107) and for assessing visual information processing—for example, via eye-tracking (108) or brain activity measurements to assess cognitive load (109). Hopefully, these evaluation methods will soon provide objective measures of the value of a visualization (106) that are recognized and agreed upon by the research community. Together with other advances in data visualization and user experience design (103), this may soon provide a new generation of tools that are much more powerful yet also easier to learn and use. Such tools would significantly reduce many of the current frustrations of scientists and clinicians and revolutionize how we see and think about our data.

To understand and gain insight from the large, complex data sets generated in biomedical research, we need tailored visualization methods and tools that present the right data and analysis to the right researcher or clinician at the right time, providing a clear view of the inherent complexity in our data, not the complication of oversimplification (paraphrased from Reference 5, p. 191). The development and adoption of such methods and tools will require fundamental changes to current research, communication, training, and clinical practices. Without these changes, many biomedical insights will remain undiscovered and misdiagnoses will remain unrecognized, buried in data already collected.

## DISCLOSURE STATEMENT

A.E.D. is the founder of a company developing long-read sequencing technology for genomes and microbiomes, the Vice President of the Australian Bioinformatics and Computational Biology Society (ABACBS), and an editor at *PLOS Computational Biology*.

## ACKNOWLEDGMENTS

Much of this review emerged from discussions at VIZBI 2017 in Sydney, Australia. We are grateful for the many colleagues who contributed, helping articulate the problems of specific areas in the biosciences, and to Bang Wong, Tamara Munzner, and Jeffery Heer, who have developed groundbreaking best practices. We are grateful for funding organizations that have supported the VIZBI meeting series, especially NIH, EMBO, the University of Sydney, and the Garvan Institute of Medical Research. We would also like to acknowledge useful discussion and comments with Drew Berry, Christian Stolte, David James, and Sean Humphrey.

## LITERATURE CITED

1. O'Donoghue SI, Gavin A-C, Gehlenborg N, Goodsell DS, Hériché J-K, et al. 2010. Visualizing biological data—now and in the future. *Nat. Methods* 7:S2–4
2. Graber ML, Franklin N, Gordon R. 2005. Diagnostic error in internal medicine. *Arch. Intern. Med.* 165:1493–99
3. Pinto A, Brunese L. 2010. Spectrum of diagnostic errors in radiology. *World J. Radiol.* 2:377–83
4. Makary MA, Daniel M. 2016. Medical error—the third leading cause of death in the US. *BMJ* 353:i2139
5. Tufte ER. 2009. *The Visual Display of Quantitative Information*. Cheshire, CT: Graphics
6. Evanko D. 2013. Data visualization: a view of every Points of View column. *Methagora: A Blog from Nature Methods*, July 13. <http://blogs.nature.com/methagora/2013/07/data-visualization-points-of-view.html>
7. Rougier NP, Droettboom M, Bourne PE. 2014. Ten simple rules for better figures. *PLOS Comput. Biol.* 10:e1003833
8. Munzner T. 2014. *Visualization Analysis and Design*. Boca Raton, FL: CRC
9. Card SK, Mackinlay JD, Shneiderman B. 1999. *Readings in Information Visualization: Using Vision to Think*. San Francisco: Morgan Kaufmann
10. Borland D, Taylor MR II. 2007. Rainbow color map (still) considered harmful. *IEEE Comput. Graph. Appl.* 27:14–17
11. Craft M, Dobrenz B, Dornbush E, Hunter M, Morris J, et al. 2015. An assessment of visualization tools for patient monitoring and medical decision making. *Proc. Syst. Inf. Eng. Des. Symp.*, 24 Apr., Charlottesville, Va., pp. 212–17. New York: IEEE
12. Lewandowsky S, Spence I. 1989. The perception of statistical graphs. *Sociol. Methods Res.* 18:200–42

---

1. *Nature Methods* special issue on visualizing biological data, covering molecular biology, biomedical science, and evolution.

---

5. Inspirational, groundbreaking collection of historical and modern approaches to displaying quantitative data.

---

6. *Nature Methods* regularly publishes 1-page articles focused on specific visualization issues faced by life scientists.

---

7. Concise, practical guide to principles and tools for creating scientific figures.

---

8. Comprehensive overview of data visualization principles.

---

9. Definitive, annotated guide to classic papers on information visualization.

---



13. Koch K, McLean J, Segev R, Freed MA, Berry MJ 2nd, et al. 2006. How much the eye tells the brain. *Curr. Biol.* 16:1428–34
14. Healey CG, Enns JT. 2012. Attention and visual memory in visualization and computer graphics. *IEEE Trans. Vis. Comput. Graph.* 18:1170–88
15. Ball R, North C. 2007. Realizing embodied interaction for visual analytics through large displays. *Comput. Graph.* 31:380–400
16. Cleveland WS, Diaconis P, McGill R. 1982. Variables on scatterplots look more highly correlated when the scales are increased. *Science* 216:1138–41
17. Heer J, Kong N, Agrawala M. 2009. Sizing the horizon: the effects of chart size and layering on the graphical perception of time series visualizations. *Proc. Int. Conf. Human Factors Comput. Syst., Boston, Mass., 4–9 Apr.*, pp. 1303–12. New York: Assoc. Comput. Mach.
18. Inselberg A. 1997. Multidimensional detective. *Proc. IEEE Symp. Inf. Vis., Phoenix, Ariz., 21 Oct.*, pp. 100–7. New York: IEEE
19. Hegarty M. 2011. The cognitive science of visual-spatial displays: implications for design. *Top. Cogn. Sci.* 3:446–74
20. Nielsen CB, Cantor M, Dubchak I, Gordon D, Wang T. 2010. Visualizing genomes: techniques and challenges. *Nat. Methods* 7:S5–15
21. Procter JB, Barton GJ, Thompson J, Westhof E, Creevey C, Letunic I. 2010. Visualization of multiple alignments, phylogenies and gene family evolution. *Nat. Methods* 7:S16–25
22. O'Donoghue SI, Goodsell DS, Frangakis AS, Jossinet F, Laskowski R, et al. 2010. Visualization of macromolecular structures. *Nat. Methods* 7:S42–55
23. Gehlenborg N, O'Donoghue SI, Baliga NS, Goessmann A, Hibbs MA, et al. 2010. Visualization of omics data for systems biology. *Nat. Methods* 7:S56–68
24. Walter T, Shattuck D, Baldock R, Bastin M, Carpenter AE, et al. 2010. Visualization of image data from cells to organisms. *Nat. Methods* 7:S26–41
25. Ware C. 2004. ***Information Visualization: Perception for Design***. San Francisco: Morgan Kaufmann
26. Soegaard M, Rikke Friis D, eds. 2013. *The Encyclopedia of Human-Computer Interaction*. Aarhus, Den.: Interact. Des. Found. 2nd ed. <https://www.interaction-design.org/literature/book/the-encyclopedia-of-human-computer-interaction-2nd-ed>
27. Johnson GT, Hertig S. 2014. A guide to the visual analysis and communication of biomolecular structural data. *Nat. Rev. Mol. Cell Biol.* 15:690–98
28. Inselberg A. 2009. ***Parallel Coordinates: Visual Multidimensional Geometry and Its Applications***. New York: Springer
29. Levy SE, Myers RM. 2016. Advancements in next-generation sequencing. *Annu. Rev. Genom. Hum. Genet.* 17:95–115
30. Pavlopoulos GA, Malliarakis D, Papanikolaou N, Theodosiou T, Enright AJ, Iliopoulos I. 2015. Visualizing genome and systems biology: technologies, tools, implementation techniques and trends, past, present and future. *GigaScience* 4:1–27
31. Shneiderman B. 1996. The eyes have it: a task by data type taxonomy for information visualizations. *Proc. IEEE Symp. Vis. Lang., Boulder, Colo., 3–6 Sept.*, pp. 336–43. Los Alamitos, NM: IEEE Comput. Soc.
32. Ernst J, Kellis M. 2012. ChromHMM: automating chromatin-state discovery and characterization. *Nat. Methods* 9:215–16
33. Pabinger S, Dander A, Fischer M, Snajder R, Sperk M, et al. 2014. A survey of tools for variant analysis of next-generation genome sequencing data. *Briefings Bioinform.* 15:256–78
34. Schroeder MP, Gonzalez-Perez A, Lopez-Bigas N. 2013. Visualizing multidimensional cancer genomics data. *Genome Med.* 5:9
35. Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragozy T, et al. 2009. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science* 326:289–93
36. Serra F, Di Stefano M, Spill YG, Cuartero Y, Goodstadt M, et al. 2015. Restraint-based three-dimensional modeling of genomes and genomic domains. *FEBS Lett.* 589:2987–95
37. Ay F, Noble WS. 2015. Analysis methods for studying the 3D architecture of the genome. *Genome Biol.* 16:183

---

25. Outline of key principles and methods for interactive display of visual information.

---

27. Visual analysis and communication guide for biomolecular data; also relevant to other biomedical data.

---

28. Definitive guide to the theory and practice of using parallel coordinates to explore high-dimensional data.

---

38. Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, et al. 2016. Juicebox provides a visualization system for Hi-C contact maps with unlimited zoom. *Cell Syst.* 3:99–101
39. Zhou X, Li D, Lowdon RF, Costello JF, Wang T. 2014. methylC track: visual integration of single-base resolution DNA methylation data on the WashU EpiGenome Browser. *Bioinformatics* 30:2206–7
40. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, et al. 2009. Circos: an information aesthetic for comparative genomics. *Genome Res.* 19:1639–45
41. Taberlay PC, Achinger-Kawecka J, Lun ATL, Buske FA, Sabir KS, et al. 2016. Three-dimensional disorganisation of the cancer genome occurs coincident with long range genetic and epigenetic alterations. *Genome Res.* 26:719–31
42. Shalon D, Smith SJ, Brown PO. 1996. A DNA microarray system for analyzing complex DNA samples using two-color fluorescent probe hybridization. *Genome Res.* 6:639–45
43. Wills QF, Livak KJ, Tipping AJ, Enver T, Goldson AJ, et al. 2013. Single-cell gene expression analysis reveals genetic associations masked in whole-tissue experiments. *Nat. Biotechnol.* 31:748–52
44. Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, et al. 2008. The transcriptional landscape of the yeast genome defined by RNA sequencing. *Science* 320:1344–49
45. Gierlinski M, Cole C, Schofield P, Schurch NJ, Sherstnev A, et al. 2015. Statistical models for RNA-seq data derived from a two-condition 48-replicate experiment. *Bioinformatics* 31:3625–30
46. Wilkinson L, Friendly M. 2009. The history of the cluster heat map. *Am. Stat.* 63:179–84
47. Wong B. 2010. Points of view: color coding. *Nat. Methods* 7:573
48. Pereverzeva M, Murray SO. 2014. Luminance gradient configuration determines perceived lightness in a simple geometric illusion. *Front. Hum. Neurosci.* 8:977
49. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, et al. 2014. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.* 32:381–86
50. Cao J, Packer JS, Ramani V, Cusanovich DA, Huynh C, et al. 2017. Comprehensive single-cell transcriptional profiling of a multicellular organism by combinatorial indexing. *Science* 357:661–67
51. Berman H, Henrick K, Nakamura H. 2003. Announcing the worldwide Protein Data Bank. *Nat. Struct. Biol.* 10:980
52. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, et al. 2004. UCSF Chimera—a visualization system for exploratory research and analysis. *J. Comput. Chem.* 25:1605–12
53. Kozlikova B, Krone M, Lindow N, Falk M, Baaden M, et al. 2015. Visualization of biomolecular structures: state of the art. *Proc. Eurograph. Conf. Vis., Cagliari, Ital., 25–29 May*, ed. R Borgo, F Ganovelli, I Viola, pp. 61–82. Geneva: Eurograph. Assoc.
54. Kendrew J, Dickerson R, Strandberg B, Hart R, Davies D, et al. 1960. Structure of myoglobin: a three-dimensional Fourier synthesis at 2 Å resolution. *Nature* 185:422–27
55. Farrugia L. 2012. WinGX and ORTEP for Windows: an update. *J. Appl. Crystallogr.* 45:849–54
56. Chung JC, Harris MR, Brooks FP, Fuchs H, Kelley MT, et al. 1989. Exploring virtual worlds with head-mounted displays. *Proc. Non-Hologr. Three-Dimens. Vis. Disp. Technol., Los Angeles, Calif., 15–20 Jan.*, ed. SS Fisher, WE Robbins, pp. 15–20. Bellingham, Wash.: SPIE
57. Humphrey W, Dalke A, Schulten K. 1996. VMD: visual molecular dynamics. *J. Mol. Graph.* 14:33–38
58. Gillet A, Sanner M, Stoffer D, Olson A. 2005. Tangible interfaces for structural molecular biology. *Structure* 13:483–91
59. Sabir KS, Stolte C, Tabor B, O'Donoghue SI. 2013. The Molecular Control Toolkit: controlling 3D molecular graphics via gesture and voice. *Proc. IEEE Symp. Biol. Data Vis., Atlanta, Ga., 13–14 Oct.*, ed. J Roerdink, J Kennedy, pp. 49–56. New York: IEEE
60. Gillet A, Sanner M, Stoffer D, Goodsell D, Olson A. 2004. Augmented reality with tangible auto-fabricated models for molecular biology applications. *Proc. IEEE Visualization, Austin, Tex., 10–15 Oct.*, ed. H Rushmeier, G Turk, JJ van Wijk, pp. 245–41. New York: IEEE
61. Heinrich J, Vuong J, Hammang CJ, Wu A, Rittenbruch M, et al. 2016. Evaluating viewpoint entropy for ribbon representation of protein structure. *Comput. Graph. Forum* 35:181–90
62. Lv Z, Tek A, Da Silva F, Empereur-Mot C, Chavent M, Baaden M. 2013. Game on, science—how video game technology may help biologists tackle visualization challenges. *PLOS ONE* 8:e57990
63. O'Donoghue SI, Sabir KS, Kalemantov M, Stolte C, Wellmann B, et al. 2015. Aquaria: simplifying discovery and insight from protein structures. *Nat. Methods* 12:98–99

64. Buja A, McDonald JA, Michalak J, Stuetzle W. 1991. Interactive data visualization using focusing and linking. *Proc. IEEE Conf. Vis., San Diego, Calif., 22–25 Oct.*, ed. GM Nielson, L Rosenblum, pp. 156–63. New York: IEEE
65. Levitt M. 2009. Nature of the protein universe. *PNAS* 106:11079–84
66. Perdigão N, Heinrich J, Stolte C, Sabir KS, Buckley MJ, et al. 2015. Unexpected features of the dark proteome. *PNAS* 112:15898–903
67. Rysavy SJ, Beck DA, Daggett V. 2014. Dymeomics: data-driven methods and models for utilizing large-scale protein structure repositories for improving fragment-based loop prediction. *Protein Sci.* 23:1584–95
68. Bai X-C, McMullan G, Scheres SH. 2015. How cryo-EM is revolutionizing structural biology. *Trends Biochem. Sci.* 40:49–57
69. Johnson GT, Autin L, Al-Alusi M, Goodsell DS, Sanner MF, Olson AJ. 2015. cellPACK: a virtual mesoscope to model and visualize structural systems biology. *Nat. Methods* 12:85–91
70. Ghosh S, Matsuoka Y, Asai Y, Hsin K-Y, Kitano H. 2011. Software for systems biology: from tools to integrated platforms. *Nat. Rev. Genet.* 12:821–32
71. Kitano H. 2002. Systems biology: a brief overview. *Science* 295:1662–64
72. Schomburg D, Michal G. 2012. *Biochemical Pathways: An Atlas of Biochemistry and Molecular Biology*. Hoboken, NJ: Wiley
73. Bader GD, Cary MP, Sander C. 2006. Pathguide: a pathway resource list. *Nucleic Acids Res.* 34:D504–6
74. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, et al. 2003. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* 13:2498–504
75. Bastian M, Heymann S, Jacomy M. 2009. Gephi: an open source software for exploring and manipulating networks. *Proc. Int. Conf. Weblogs Soc. Media, San Jose, Calif., 17–20 May*, pp. 361–62. Menlo Park, CA: Assoc. Adv. Artif. Intell.
76. Kobourov SG. 2013. Force-directed drawing algorithms. In *Handbook of Graph Drawing and Visualization*, ed. R Tamassia, pp. 383–408. Boca Raton, FL: CRC
77. Barsky A, Gardy JL, Hancock RE, Munzner T. 2007. Cerebral: a Cytoscape plugin for layout of and interaction with biological networks using subcellular localization annotation. *Bioinformatics* 23:1040–42
78. Holten D, Van Wijk JJ. 2009. Force-directed edge bundling for graph visualization. *Comput. Graph. Forum* 28:983–90
79. Barsky A, Munzner T, Gardy J, Kincaid R. 2008. Cerebral: visualizing multiple experimental conditions on a graph with biological context. *IEEE Trans. Vis. Comput. Graph.* 14:1253–60
80. Zhou H-X, Rivas G, Minton AP. 2008. Macromolecular crowding and confinement: biochemical, biophysical, and potential physiological consequences. *Annu. Rev. Biophys.* 37:375–97
81. Von Landesberger T, Kuijper A, Schreck T, Kohlhammer J, van Wijk JJ, et al. 2011. Visual analysis of large graphs: state-of-the-art and future research challenges. Presented at *Comput. Graph. Forum* 30:1719–49
82. Kwon O-H, Crnovrsanin T, Ma K-L. 2017. What would a graph look like in this layout? A machine learning approach to large graph visualization. *IEEE Trans. Vis. Comput. Graph.* 24:478–88
83. Aebersold R, Mann M. 2016. Mass-spectrometric exploration of proteome structure and function. *Nature* 537:347–55
84. Humphrey SJ, Azimifar SB, Mann M. 2015. High-throughput phosphoproteomics reveals in vivo insulin signaling dynamics. *Nat. Biotechnol.* 33:990–95
85. Ma DK, Stolte C, Krycer JR, James DE, O'Donoghue SI. 2015. SnapShot: insulin/IGF1 signaling. *Cell* 161:948.e1
86. Burgess A, Vuong J, Rogers S, Malumbres M, O'Donoghue SI. 2017. SnapShot: phosphoregulation of mitosis. *Cell* 169:1358.e1
87. Sydor AM, Czymmek KJ, Puchner EM, Mennella V. 2015. Super-resolution microscopy: from single molecules to supramolecular assemblies. *Trends Cell Biol.* 25:730–48
88. Reynaud EG, Peychl J, Huisken J, Tomancak P. 2015. Guide to light-sheet microscopy for adventurous biologists. *Nat. Methods* 12:30–34
89. Walter T, Shattuck DW, Baldock R, Bastin ME, Carpenter AE, et al. 2010. Visualization of image data from cells to organisms. *Nat. Methods* 7:S26–41

90. Rossner M, Yamada KM. 2004. What's in a picture? The temptation of image manipulation. *J. Cell Biol.* 166:11–15
91. Burel J-M, Besson S, Blackburn C, Carroll M, Ferguson RK, et al. 2015. Publishing and sharing multi-dimensional image data with OMERO. *Mamm. Genome* 26:441–47
92. Bernhardt S, Nicolau SA, Soler L, Doignon C. 2017. The status of augmented reality in laparoscopic surgery as of 2016. *Med. Image Anal.* 37:66–90
93. Maier-Hein L, Vedula SS, Speidel S, Navab N, Kikinis R, et al. 2017. Surgical data science for next-generation interventions. *Nat. Biomed. Eng.* 1:691–96
94. Ding W, Neher R. 2014. *panX*. Pangenome Vis. Tool. <http://pangenome.tuebingen.mpg.de/>
95. Parks DH, Porter M, Churcher S, Wang S, Blouin C, et al. 2009. GenGIS: a geospatial information system for genomic data. *Genome Res.* 19:1896–904
96. Argimon S, Abudahab K, Goater RJ, Fedosejev A, Bhai J, et al. 2016. Microreact: visualizing and sharing data for genomic epidemiology and phylogeography. *Microb. Genom* 2:e000093
97. Paten B, Novak AM, Eizenga JM, Garrison E. 2017. Genome graphs and the evolution of genome inference. *Genome Res.* 27:665–76
98. Wong B. 2011. Points of view: color blindness. *Nat. Methods* 8:441
99. Cromey DW. 2010. Avoiding twisted pixels: ethical guidelines for the appropriate use and manipulation of scientific digital images. *Sci. Eng. Ethics* 16:639–67
100. McGill G. 2008. Molecular movies. . . Coming to a lecture near you. *Cell* 133:1127–32
101. Iwasa JH. 2015. Bringing macromolecular machinery to life using 3D animation. *Curr. Opin. Struct. Biol.* 31:84–88
102. Van Noorden R, Maher B, Nuzzo R. 2014. The top 100 papers. *Nature* 514:550–53
103. Shneiderman B, Plaisant C, Cohen M, Jacobs S, Elmqvist N, et al. 2018. *Designing the User Interface: Strategies for Effective Human-Computer Interaction*. Boston: Pearson. 6th ed.
104. Sanyal J, Zhang S, Bhattacharya G, Amburn P, Moorhead R. 2009. A user study to compare four uncertainty visualization methods for 1D and 2D datasets. *IEEE Trans. Vis. Comput. Graph.* 15:1209–18
105. Callaway E. 2016. The visualizations transforming biology. *Nature* 535:187–88
106. Van Wijk JJ. 2005. The value of visualization. *Proc. IEEE Visualization, Minneap., Minn., 23–28 Oct.*, pp. 79–86. New York: IEEE
107. Heer J, Bostock M. 2010. Crowdsourcing graphical perception: using mechanical turk to assess visualization design. *Proc. CHI Conf. Human Factors Comput. Syst., Atlanta, Ga., 10–15 Apr.*, pp. 203–12. New York: Assoc. Comput. Mach.
108. Blaschek T, Kurzhals K, Raschke M, Burch M, Weiskopf D, Ertl T. 2014. State-of-the-art of visualization for eye tracking data. *Proc. EuroVis Eurograph. Conf. Vis., Swansea, Wales, 9–13 June*. Geneva: Eurograph. Assoc.
109. Anderson EW, Potter KC, Matzen LE, Shepherd JF, Preston GA, Silva CT. 2011. A user study of visualization effectiveness using EEG and cognitive load. *Comput. Graph. Forum* 30:791–800
110. Gehlenborg N, Wong B. 2012. Mapping quantitative data to color: data structure informs choice of color maps. *Nat. Methods* 9:769–70
111. Wong B. 2011. Points of view: avoiding color. *Nat. Methods* 8:525
112. Tufte ER. 1990. *Envisioning Information*. Cheshire, CT: Graphics
113. Gehlenborg N, Wong B. 2012. Points of view: into the third dimension. *Nat. Methods* 9:851
114. Kabsch W, Mannherz HG, Suck D, Pai EF, Holmes KC. 1990. Atomic structure of the actin:DNase I complex. *Nature* 347:37–44
115. Hunter JD. 2007. Matplotlib: a 2D graphics environment. *Comput. Sci. Eng.* 9:90–95
116. Wickham H. 2016. *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag
117. Anscombe FJ. 1973. Graphs in statistical analysis. *Am. Statistician* 27:17–21
118. Matejka J, Fitzmaurice G. 2017. Same stats, different graphs: generating datasets with varied appearance and identical statistics through simulated annealing. *Proc. CHI Conf. Human Factors Comput. Syst., Denver, Colo., 6–11 May*, pp. 1290–94. New York: Assoc. Comput. Mach.
119. Fry B. 2008. *Visualizing Data: Exploring and Explaining Data with the Processing Environment*. Sebastopol, CA: O'Reilly Media

120. Moreland K. 2016. Why we use bad color maps and what you can do about it. *IS&T Int. Symp. Electron. Imaging, San Francisco, Calif., 14–18 Feb.*, ed. BE Rogowitz, TN Pappas, D de Ridder, pp. 1–6(6). Springfield, VA: Soc. Imaging Sci. Technol.
121. Mackinlay JD. 1986. Automating the design of graphical presentations of relational information. *ACM Trans. Graph.* 5:110–15
122. Gehlenborg N, Wong B. 2012. Points of view: heat maps. *Nat. Methods* 9:213
123. Yates A, Akanni W, Amode MR, Barrell D, Billis K, et al. 2015. Ensembl 2016. *Nucleic Acids Res.* 44:D710–16
124. Down TA, Piipari M, Hubbard TJ. 2011. Dalliace: interactive genome viewing on the web. *Bioinformatics* 27:889–90
125. Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, et al. 2014. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* 159:1665–80
126. van der Maaten L, Hinton G. 2008. Visualizing data using t-SNE. *J. Mach. Learn. Res.* 9:2579–605
127. Wattenberg M, Viégas F, Johnson I. 2016. How to use t-SNE effectively. *Distill*. Updated on 13 Oct. 2016. <http://doi.org/10.23915/distill.00002>
128. Haghverdi L, Buettner F, Theis FJ. 2015. Diffusion maps for high-dimensional single-cell analysis of differentiation data. *Bioinformatics* 31:2989–98
129. Coifman RR, Lafon S, Lee AB, Maggioni M, Nadler B, et al. 2005. Geometric diffusions as a tool for harmonic analysis and structure definition of data: diffusion maps. *PNAS* 102:7426–31
130. McCarthy DJ, Campbell KR, Lun AT, Wills QF. 2017. Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* 33:1179–86
131. Camp JG, Sekine K, Gerber T, Loeffler-Wirth H, Binder H, et al. 2017. Multilineage communication regulates human liver bud development from pluripotency. *Nature* 546:533–38
132. Richardson JS, Richardson D, Tweedy N, Gernert K, Quinn T, et al. 1992. Looking at proteins: representations, folding, packing, and design. Biophysical Society National Lecture, 1992. *Biophys. J.* 63:1185–1209
133. Heinrich J, Kaur S, O'Donoghue SI. 2015. Evaluating the effectiveness of color to convey uncertainty in macromolecular structures. *Proc. IEEE Symp. Big Data Visual Anal., Hobart, Aust., 22–25 Sept.*, ed. U Engelke, J Heinrich, T Bednarsz, K Klein, QV Nguyen, pp. 1–18. New York: IEEE
134. Dosztanyi Z, Csizmek V, Tompa P, Simon I. 2005. IUPred: web server for the prediction of intrinsically unstructured regions of proteins based on estimated energy content. *Bioinformatics* 21:3433–34
135. Rose AS, Bradley AR, Valasatava Y, Duarte JM, Prlić A, Rose PW. 2016. Web-based molecular graphics for large complexes. *Proc. Int. Conf. Web3D Technol., 21st, Anaheim, Calif., 22–24 July*, pp. 185–86. New York: Assoc. Comput. Mech.
136. Zhao G, Perilla JR, Yufenyuy EL, Meng X, Chen B, et al. 2013. Mature HIV-1 capsid structure by cryo-electron microscopy and all-atom molecular dynamics. *Nature* 497:643–46
137. El Omari K, De Mesmaeker J, Karia D, Ginn H, Bhattacharya S, Mancini EJ. 2012. Structure of the DNA-bound T-box domain of human TBX1, a transcription factor associated with the DiGeorge syndrome. *Proteins* 80:655–60
138. Isberg V, Mordalski S, Munk C, Rataj K, Harpsøe K, et al. 2015. GPCRdb: an information system for G protein-coupled receptors. *Nucleic Acids Res.* 44:D356–64
139. Dunker AK, Lawson JD, Brown CJ, Williams RM, Romero P, et al. 2001. Intrinsically disordered protein. *J. Mol. Graph. Model.* 19:26–59
140. Humphrey SJ, Yang G, Yang P, Fazakerley DJ, Stöckli J, et al. 2013. Dynamic adipocyte phosphoproteome reveals that Akt directly regulates mTORC2. *Cell Metab.* 17:1009–20
141. Szklarczyk D, Franceschini A, Wyder S, Forslund K, Heller D, et al. 2015. STRING v10: protein–protein interaction networks, integrated over the tree of life. *Nucleic Acids Res.* 43:D447–52
142. Lu CP, Hager GD, Mjolsness E. 2000. Fast and globally convergent pose estimation from video images. *IEEE Trans. Pattern Anal. Mach. Intell.* 22:610–22
143. Porter IM, McClelland SE, Khoudoli GA, Hunter CJ, Andersen JS, et al. 2007. Bod1, a novel kinetochore protein required for chromosome biorientation. *J. Cell Biol.* 179:187–97

144. Nolden M, Zelzer S, Seitel A, Wald D, Muller M, et al. 2013. The Medical Imaging Interaction Toolkit: challenges and advances: 10 years of open-source development. *Int. J. Comput. Assist. Radiol. Surg.* 8:607–20
145. Simpfendorfer T, Baumhauer M, Muller M, Gutt CN, Meinzer HP, et al. 2011. Augmented reality visualization during laparoscopic radical prostatectomy. *J. Endourol.* 25:1841–45
146. Han MV, Zmasek CM. 2009. phyloXML: XML for evolutionary biology and comparative genomics. *BMC Bioinform.* 10:356
147. Darling ACE, Mau B, Blattner FR, Perna NT. 2004. Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Res.* 14:1394–403
148. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, et al. 2010. QIIME allows analysis of high-throughput community sequencing data. *Nat. Methods* 7:335–36
149. Rodgers P, Stapleton G, Chapman P. 2015. Visualizing sets with linear diagrams. *ACM Trans. Comput.-Hum. Interact.* 22:27
150. Caporaso JG, Lauber CL, Costello EK, Berg-Lyons D, Gonzalez A, et al. 2011. Moving pictures of the human microbiome. *Genome Biol.* 12:R50

# Contents

Big Data Approaches for Modeling Response and Resistance to Cancer Drugs <i>Peng Jiang, William R. Sellers, and X. Shirley Liu</i> .....	1
From Tissues to Cell Types and Back: Single-Cell Gene Expression Analysis of Tissue Architecture <i>Xi Chen, Sarah A. Teichmann, and Kerstin B. Meyer</i> .....	29
Advances in Electronic Phenotyping: From Rule-Based Definitions to Machine Learning Models <i>Juan M. Banda, Martin Seneviratne, Tina Hernandez-Boussard, and Nigam H. Shah</i> .....	53
Defining Phenotypes from Clinical Data to Drive Genomic Research <i>Jamie R. Robinson, Wei-Qi Wei, Dan M. Roden, and Joshua C. Denny</i> .....	69
Alignment-Free Sequence Analysis and Applications <i>Jie Ren, Xin Bai, Yang Young Lu, Kujin Tang, Ying Wang, Gesine Reinert, and Fengzhu Sun</i> .....	93
Privacy Policy and Technology in Biomedical Data Science <i>April Moreno Arellano, Wenrui Dai, Shuang Wang, Xiaoqian Jiang, and Lucila Ohno-Machado</i> .....	115
Opportunities and Challenges of Whole-Cell and -Tissue Simulations of the Outer Retina in Health and Disease <i>Philip J. Luthert, Luis Serrano, and Christina Kiel</i> .....	131
Network Analysis as a Grand Unifier in Biomedical Data Science <i>Patrick McGillivray, Declan Clarke, William Meyerson, Jing Zhang, Donghoon Lee, Mengting Gu, Sushant Kumar, Holly Zhou, and Mark Gerstein</i> .....	153
Deep Learning in Biomedical Data Science <i>Pierre Baldi</i> .....	181
Computational Methods for Understanding Mass Spectrometry-Based Shotgun Proteomics Data <i>Pavel Sinitcyn, Jan Daniel Rudolph, and Jürgen Cox</i> .....	207
Data Science Issues in Studying Protein-RNA Interactions with CLIP Technologies <i>Anob M. Chakrabarti, Nejc Haberman, Arne Praznik, Nicholas M. Luscombe, and Jernej Ule</i> .....	235

Large-Scale Analysis of Genetic and Clinical Patient Data <i>Marylyn D. Ritchie</i> .....	263
Visualization of Biomedical Data <i>Seán I. O'Donoghue, Benedetta Frida Baldi, Susan J. Clark, Aaron E. Darling, James M. Hogan, Sandeep Kaur, Lena Maier-Hein, Davis J. McCarthy, William J. Moore, Esther Stenau, Jason R. Swedlow, Jenny Vuong, and James B. Procter</i> .....	275
A Census of Disease Ontologies <i>Melissa Haendel, Julie McMurry, Rose Relevo, Chris Mungall, Peter Robinson, and Christopher G. Chute</i> .....	305

### Errata

An online log of corrections to *Annual Review of Biomedical Data Science* articles may be found at <http://www.annualreviews.org/errata/biodatasci>