

1

# TOPOLOGICAL ARTIST MODEL

2

HANNAH AIZENMAN

3

A DISSERTATION PROPOSAL SUBMITTED TO

4

THE GRADUATE FACULTY IN COMPUTER SCIENCE IN PARTIAL FULFILLMENT OF THE

5

REQUIREMENTS FOR THE DEGREE OF DOCTOR OF PHILOSOPHY,

6

THE CITY UNIVERSITY OF NEW YORK

7

COMMITTEE MEMBERS:

8

DR. MICHAEL GROSSBERG (ADVISOR), DR. ROBERT HARALICK, DR. LEV MANOVICH,

9

DR. HUY VO, DR. MARCUS HANWELL

10

JUNE 2021

# Abstract

A critical aspect of data visualization is that the graphical representation of data is expected to match the properties of the data; this fails when order is not preserved in representations of ordinal data or scale for numerical data. In this work, we propose that the mathematical notions of equivariance and topology formalizes the expectation of matching properties. We developed a model we call the topological artist model (TAM) in which data and graphics can be viewed as sections of fiber bundles. This model allows for (1) decomposing the translation of data fields (variables) into visual channels via an equivariant map on the fibers and (2) a topology-preserving map of the base spaces that translates the dataset connectivity into graphical elements. Furthermore, our model supports an algebraic sum operation such that more complex visualizations can be built from simple ones. We illustrate the application of the model through case studies of a scatter plot, line plot, and heatmap. We show that this model can be implemented with a small prototype.

To demonstrate the practical value of our model, we propose a model driven re-architecture of the artist layer of the Python visualization library Matplotlib. We can show that we can concretely represent the base spaces, make use of programming types for the fiber, and build on Matplotlib's existing infrastructure for the rendering. In addition to providing a way to ensure the library preserves structure, the functional decomposition of the artist in the model could improve modularity, maintainability, and point to ways in which the library provides better support concurrency and interactivity. The thesis will follow through on this proposal to explore how to further develop our model, showing how it can support Matplotlib's current diverse range of data visualizations while providing a better platform for domain-specific visualization library developers.

# Contents

35	<b>Abstract</b>	<b>ii</b>
36	<b>1 Introduction</b>	<b>1</b>
37	<b>2 Background</b>	<b>2</b>
38	2.1 Tools . . . . .	3
39	2.2 Data . . . . .	4
40	2.3 Visualization . . . . .	7
41	2.4 Contribution . . . . .	10
42	<b>3 Topological Artist Model</b>	<b>10</b>
43	3.1 Data Space $E$ . . . . .	11
44	3.1.1 Variables in Fiber Space $F$ . . . . .	12
45	3.1.2 Measurement Scales: Monoid Actions . . . . .	14
46	3.1.3 Continuity of the Data $K$ . . . . .	15
47	3.1.4 Data $\tau$ . . . . .	18
48	3.1.5 Applications to Data Containers . . . . .	19
49	3.2 Graphic Space $H$ . . . . .	20
50	3.2.1 Idealized Display $D$ . . . . .	20
51	3.2.2 Continuity of the Graphic $S$ . . . . .	21
52	3.2.3 Graphic $\rho$ . . . . .	23
53	3.3 Artist . . . . .	24
54	3.3.1 Visual Fiber Bundle $V$ . . . . .	25
55	3.3.2 Visual Encoders $\nu$ . . . . .	27
56	3.3.3 Graphic Assembler $Q$ . . . . .	29
57	3.3.4 Assembly $Q$ . . . . .	31
58	3.3.5 Assembly factory $\hat{Q}$ . . . . .	34
59	3.3.6 Sheafs . . . . .	36
60	3.3.7 Composition of Artists: $+$ . . . . .	36

61	3.3.8 Equivalence class of artists $A'$ . . . . .	37
62	<b>4 Prototype Implementation: Matplottoy</b>	<b>38</b>
63	4.1 Artist Class $A'$ . . . . .	39
64	4.2 Encoders $\nu$ . . . . .	44
65	4.3 Data $E$ . . . . .	45
66	4.4 Case Study: Penguins . . . . .	50
67	<b>5 Discussion</b>	<b>54</b>
68	5.1 Limitations . . . . .	55
69	5.2 Future Work . . . . .	55
70	<b>6 Conclusion</b>	<b>57</b>

# 1 Introduction

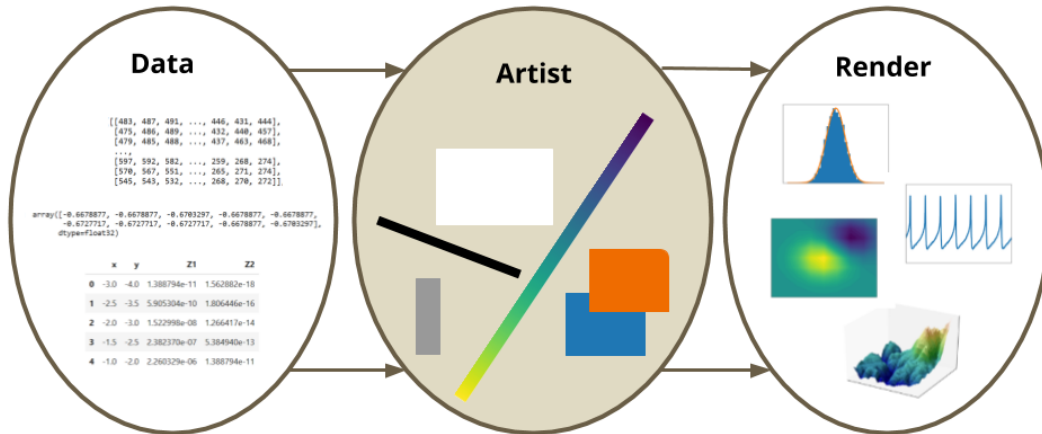


Figure 1: Visualization is a mapping from data into visual encodings that are then rendered into images. In our model, this visual encoding stage is called the artist.

The work presented in this paper is motivated by a need for a library of visualization components that developers could use to build complex, domain specific tools tuned to the semantics and structure carried in domain specific data. While many researchers have identified and described important aspects of visualization, they have specialized in such different ways as to not provide a model general enough to natively support the full range of data and visualization types many general purpose modern visualization tools may need to support. The core architecture also needs to be robust to the big data needs of many visualization practitioners, and therefore support distributed and streaming data needs. To support both exploratory and confirmatory visualization[1], this tool needs to support 2D and 3D, static, dynamic and interactive visualizations.

Specifically, this work was driven by a rearchitecture of the Python visualization library Matplotlib[2] to meet modern data visualization needs. We aim to take advantage of developments in software design, data structures, and visualization to improve the consistency, composibility, and discoverability of the API. To do so, this work first presents a mathematical description of how data is transformed into graphic representations, as shown in figure 1. As with other mathematical formalisms of visualization [3–6], a mathematical framework

88 provides a way to formalize the properties and structure of the visualization. In contrast to  
89 the other formalisms, the model presented here is focused on the components that build a  
90 visualization rather than the visualization itself.

91 In other words this model is not intended to be evaluative, it is intended to be a reference  
92 specification for visualization library API. To make this model as implementation indepen-  
93 dent as possible, we propose fairly general mathematical abstractions of the data container  
94 such that we do not need to assume the data has any specific structure, such as a relational  
95 database. We reuse this structure for the graphic as that allows us to specifically discuss  
96 how structure is preserved. We take a functional approach because functional paradigms  
97 encourage writing APIs that are flexible, concise and predictable due to the lack of side  
98 effects [7]. Furthermore, by structuring the API in terms of composition of the smallest  
99 units of transformation for which we can define correctness, a functional paradigm naturally  
100 leads to a library of highly modular components that are composable in such a way that  
101 by definition the composition is also correct. This allows us to ensure that domain specific  
102 visualizations built on top of these components are also correct without needing knowl-  
103 edge of the domain. As with the other mathematical formalisms of visualization, we factor  
104 out the rendering into a separate stage; but, our framework describes how these rendering  
105 instructions are generated.

106 In this work, we present a framework for understanding visualization as equivariant maps  
107 between topological spaces. Using this mathematical formalism, we can interpret and extend  
108 prior work and also develop new tools. We validate our model by using it to re-design artist  
109 and data access layer of Matplotlib, a general purpose visualization tool.

## 110 2 Background

111 One of the reasons we developed a new formalism rather than adopting the architecture  
112 of an existing library is that most information visualization software design patterns, as  
113 categorized by Heer and Agrawala[8], are tuned to very specific data structures. This in  
114 turn restricts the design space of visual algorithms that display information (the visualization

types the library supports) since the algorithms are designed such that the structure of data is assumed, as described in Tory and Möller’s taxonomy [9]. In proposing a new architecture, we contrast the trade offs libraries make, describe different types of data continuity, and discuss metrics by which a visualization library is traditionally evaluated.

## 2.1 Tools

One extensive family of relational table based libraries are those based on Wilkinson’s Grammar of Graphics (GoG) [10], including ggplot[11], protovis[12] and D3 [13], vega[14] and altair[15]. The restriction to tables in turn restricts the native design space to visualizations suited to tables. Since the data space and graphic space is very well defined in this grammar, it lends itself to a declarative interface [16]. This grammar oriented approach allows users to describe how to compose visual elements into a graphical design [17], while we are proposing a framework for building those elements. An example of this distinction is that the GoG grammar includes computation and aggregation of the table as part of the grammar, while we propose that most computations are specific to domains and only try to describe them when they are specifically part of the visual encoding - for example mapping data to a color. Disentangling the computation from the visual transforms allows us to determine whether the visualization library needs to handle them or if they can be more efficiently computed by the data container.

A different class of user facing tools are those that support images, such as ImageJ[18] or Napari[19]. These tools often have some support for visualizing non image components of a complex data set, but mostly in service to the image being visualized. These tools are ill suited for general purpose libraries that need to support data other than images because the architecture is oriented towards building plugins into the existing system [20] where the image is the core data structure. Even the digital humanities oriented ImageJ macro ImagePlot[21], which supports some non-image aggregate reporting charts, is still built around image data as the primary input.

There are also visualization tools where there is no single core structure, and instead internally carry around many different representations of data. Matplotlib, has this struc-

ture, as does VTK [22, 23] and its derivatives such as MayaVi[24] and extensions such as ParaView[25] and the infoviz themed Titan[26]. Where GoG and ImageJ type libraries have very consistent APIs for their visualization tools because the data structure is the same, the APIs for visualizations in VTK and Matplotlib are significantly dependent on the structure of the data it expects. This in turn means that every new type of visualization must carry implicit assumptions about data structure in how it interfaces with the input data. This has lead to poor API consistency and brittle code as every visualization type has a very different point of view on how the data is structured. This API choice particularly breaks down when the same dataset is fed into visualizations with different assumptions about structure or into a dashboard consisting of different types of visualization[27, 28] because there is no consistent way to update the data and therefore no consistent way of guaranteeing that the views stay in sync. Our model is a structure dependent formalism, but then also provides a core representation of that structure that is abstract enough to provide a common interface for many different types of visualization.

## 2.2 Data

Discrete and continuous data and their attributes form a discipline independent design space [29], so one of the drivers of this work was to facilitate building libraries that could natively support domain specific data containers that do not make assumptions about data continuity. As shown in figure 2, there are many types of connectivity. A database typically consists of unconnected records, while an image is an implicit 2D grid and a network is some sort of explicitly connected graph. These data structures typically contain not only the measurements or values of the data, but also domain specific semantic information such as that the data is a map or an image that a modern visualization library could exploit if this information was exposed to the API.



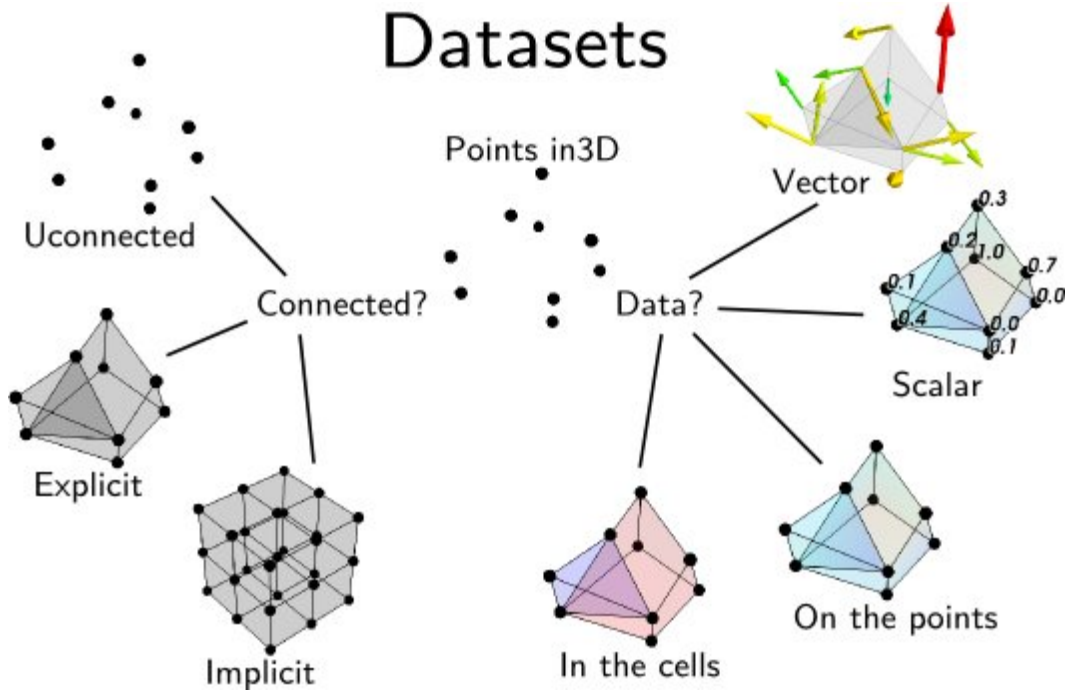


Figure 2: One way to describe data is by the connectivity of the points in the dataset. A database for example is often discrete unconnected points, while an image is an implicitly connected 2D grid. This image is from the Data Representation chapter of the MayaVi 4.7.2 documentation.[\[30\]](#)

As shown in figure 2, there are many distinct ways of encoding each specific type of structure, while as mentioned in section 2.1 APIs are clearer when structured around a common data representation. Fiber bundles were proposed by Butler as one such representation because they encode the continuity of the data separately from the types of variables and are flexible enough to support discrete and ND continuous datasets [\[31, 32\]](#). Since Butler’s model lacks a robust way of describing variables, we fold in Spivak’s Simplicial formulation of databases [\[33, 34\]](#) so that we can encode a schema like description of the data in the fiber bundle. In this work we will refer to the points of the dataset as *records* to indicate that a point can be a vector of heterogenous elements. Each *component* of the record is a single object, such as a temperature measurement, a color value, or an image. We also generalize *component* to mean all objects in the dataset of a given type, such as

178 all temperatures or colors or images. The way in which these records are connected is the  
179 *connectivity*, *continuity*, or more generally *topology*.

**definitions**

**records** points, observations, entries

**components** variables, attributes, fields

**connectivity** how the records are connected to each other

180 Often this topology has metadata associated with it, describing for example dependent  
181 variables and the independent variables they are dependent on, or information about the  
182 structure of the data such as when and where the measurement was taken. Building on  
183 the idea of metadata as *keys* and their associated *values* proposed by Munzner [35], we  
184 propose that information rich metadata are part of the components and instead the values  
185 are keyed on coordinate free structural ids. In contrast to Munzner’s model where the  
186 semantic meaning of the key is tightly coupled to the position of the value in the dataset,  
187 our model considers keys to be a pure reference to topology. This allows the metadata  
188 to be altered, for example by changing the coordinate systems or time resolution, without  
189 imposing new semantics on the underlying structure.

190 **2.3 Visualization**

	<i>Points</i>	<i>Lines</i>	<i>Areas</i>	<i>Best to show</i>
<i>Shape</i>		<i>possible, but too weird to show</i>	<i>cartogram</i>	<i>qualitative differences</i>
<i>Size</i>			<i>cartogram</i>	<i>quantitative differences</i>
<i>Color Hue</i>				<i>qualitative differences</i>
<i>Color Value</i>				<i>quantitative differences</i>
<i>Color Intensity</i>				<i>qualitative differences</i>
<i>Texture</i>				<i>qualitative &amp; quantitative differences</i>

Figure 3: Retinal variables are a codification of how position, size, shape, color and texture are used to illustrate variations in the components of a visualization. The best to show column describes which types of information can be expressed in the corresponding visual encoding. This tabular form of Bertin’s retinal variables is from Understanding Graphics [36] who reproduced it from Krygier and Wood’s *Making Maps: A Visual Guide to Map Design for GIS* [37]

191 Visual representations of data, by definition, reflect something of the underlying structure  
 192 and semantics[38], whether through direct mappings from data into visual elements or via  
 193 figurative representations that have meaning due to their similarity in shape to external  
 194 concepts [39]. The components of a visual representation were first codified by Bertin[40].  
 195 As illustrated in figure 3, Bertin proposes that there are classes of visual encodings such as  
 196 shape, color, and texture that when mapped to from specific types of measurement, quan-

197 titative or qualitative, will preserve the properties of that measurement type. For example,  
 198 that nominal data mapped to hue preserves the selectivity of the nominal measurements.  
 199 Furthermore he proposes that the visual encodings be composited into graphical marks that  
 200 match the connectivity of the data - for example discrete data is a point, 1D continuous  
 201 is the line, and 2D data is the area mark. A general form of marks are glyphs, which are  
 202 graphical objects that convey one or more attributes of the data entity mapped to it[41, 42]  
 203 and minimally need to be differentiable from other visual elements [43]. The set of encoding  
 204 relations from data to visual representation is termed the graphical design by Mackinlay [3,  
 205 44] and the design rendered in an idealized abstract space is what throughout this paper we  
 206 will refer to as a graphic.

207     The measure of how much of the structure of the data the graphic encodes is a concept  
 208 Mackinlay termed expressiveness, while the graphic’s effectiveness describes how much de-  
 209 sign choices are made in deference to perceptual saliency [42, 45–47]. When the properties  
 210 of the representation match the properties of the data, then the visualization is easier to  
 211 understand according to Norman’s Naturalness Principal[48]. These ideas are combined into  
 212 Tufte’s notion of graphical integrity, which is that a visual representation of quantitative  
 213 data must be directly proportional to the numerical quantities it represents (Lie Principal),  
 214 must have the same number of visual dimensions as the data, and should be well labeled and  
 215 contextualized, and not have any extraneous visual elements [49]. This notion of matching  
 216 is explicitly formalized by Mackinlay as a structure preserving mapping of a binary operator  
 217 from one domain to another [44]. A functional dependency framework for evaluating visual-  
 218 izations was proposed by Sugibuchi et al [5], and an algebraic basis for visualization design  
 219 and evaluation was proposed by Kindlmann and Scheidegger[4]. Vickers et al. propose a  
 220 category theory framework[6] that extends structural preservation to layout, but is focused  
 221 strictly on the design layer like the other mathematical frameworks.

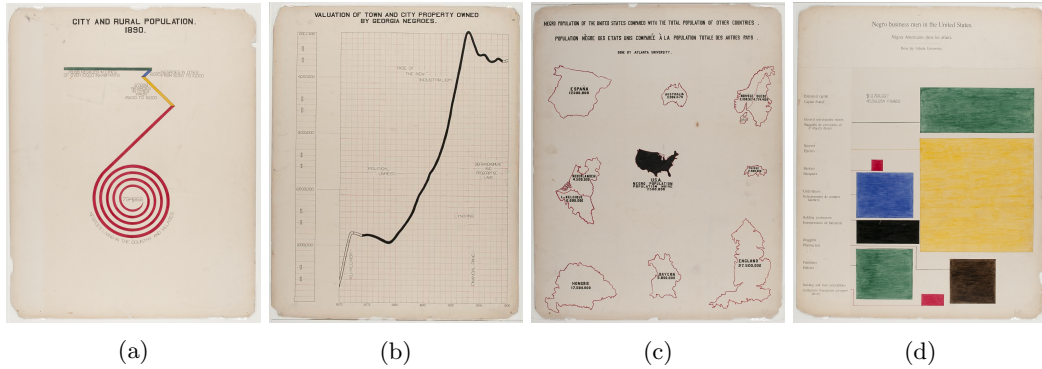


Figure 4: Du Bois’ data portraits[50] of post reconstruction Black American life exemplify that the fundamental characteristics of data visualization is that the visual elements vary in proportion to the source data. In figure 4a, the length of each segment maps to population; in figure 4b, the line changes color to indicate a shift in the political environment; in figure 4c the countries are scaled to population size; and figure 4d is a treemap where the area of the rectangle is representative of the number of businesses in each field. The images here are from the Prints and Photographs collection of the Library of Congress [51–54]

One example of highly expressive visualizations are the data portraits by Du Bois shown in figure 4. While the Du Bois charts are different from the usual scatter, line, and plot charts, they conform to the constraint that a graphic is a structure preserving map from data to visual representation. Figure 4a is semantically similar to a bar chart in that the lengths of the segments are mapped to the values, but in this chart the segments are stacked together. Figure 4b is a multicolored line chart where the color shifts are at periods of political significance. In figure 4c, Du Bois combines a graphical representation where glyph size varies by population with a figurative representation of those glyphs as the countries the data is from, which means that the semantic and numerical properties of the data are preserved in the graph. Figure ?? is simply a treemap[55] with space between the marks. Since the Du Bois data portraits meet the criteria of a faithful visual representation, we propose a mathematical framework and implementation that allows us to express the Du Bois charts and common chart types with equal fidelity.

## 2.4 Contribution

This work presents a mathematical model of the transformation from data to graphic representation and a proof of concept implementation. Specifically, the contributions of this work are

1. a formal description of the topology preserving relationship between data and graphic via continuous maps
2. a formal description of the property preservation from data component to visual representation as equivariant maps that carry a homomorphism of monoid actions
3. abstraction of data structure using fiber bundles with schema like fibers to encode components and topology
4. algebraic sum operator such that more complex visualizations can be built from simple ones
5. a functional oriented visualization tool architecture built on the mathematical model to demonstrate the utility of the model
6. a prototype of the architecture built on Matplotlib’s infrastructure to demonstrate the feasibility of the model

In contrast to mathematical models of visualization that aim to evaluate visualization design, we propose a topological framework for building tools to build visualizations. We defer judgement of expressivity and effectiveness to developers building domain specific tools, but provide them the framework to do so.

## 3 Topological Artist Model

As discussed in the introduction, visualization is generally defined as structure preserving maps from a data object to a graphic object. In order to formalize this statement, we describe the connectivity of the records using topology and define the structure on the components in

terms of the monoid actions on the component types. By formalizing structure in this way, we can evaluate the extent to which a visualization preserves the structure of the data it is representing and build structure preserving visualization tools. We introduce the notion of an artist  $\mathcal{A}$  as an equivariant map from data to graphic

$$\mathcal{A} : \mathcal{E} \rightarrow \mathcal{H} \quad (1)$$

that carries a homomorphism of monoid actions  $\varphi : M \rightarrow M'$  [56], which are discussed in detail in section 3.1.2. Given  $M$  on data  $\mathcal{E}$  and  $M'$  on graphic  $\mathcal{H}$ , we propose that artists  $\mathcal{A}$  are equivariant maps

$$\mathcal{A}(m \cdot r) = \varphi(m) \cdot \mathcal{A}(r) \quad (2)$$

such that applying a monoid action  $m \in M$  to the data  $r \in \mathcal{E}$  input to  $\mathcal{A}$  is equivalent to applying a monoid action  $\varphi(m) \in M'$  to the graphic  $\mathcal{A}(r) \in \mathcal{H}$  output of the artist.

We model the data  $\mathcal{E}$ , graphic  $\mathcal{H}$ , and intermediate visual encoding  $\mathcal{V}$  stages of visualization as topological structures that encapsulate types of variables and continuity; by doing so we can develop implementations that keep track of both in ways that let us distribute computation while still allowing assembly and dynamic update of the graphic. To explain which structure the artist is preserving, we first describe how we model data (3.1), graphics (3.2), and intermediate visual characteristics (3.3) as fiber bundles. We then discuss the equivariant maps between data and visual characteristics (3.3.2) and visual characteristics and graphics (3.3.3) that make up the artist.

### 3.1 Data Space $E$

Building on Butler’s proposal of using fiber bundles as a common data representation structure for visualization data[31, 32], a fiber bundle is a tuple  $(E, K, \pi, F)$  defined by the projection map  $\pi$

$$F \hookrightarrow E \xrightarrow{\pi} K \quad (3)$$

that binds the components of the data in  $F$  to the continuity represented in  $K$ . The fiber bundle models the properties of data component types  $F$  (3.1.1), the continuity of records  $K$  (3.1.3), the collections of records  $\tau$  (3.1.4), and the space  $E$  of all possible datasets with these components and continuity.

By definition fiber bundles are locally trivial[57, 58], meaning that over a localized neighborhood we can dispense with extra structure on  $E$  and focus on the components and continuity. We use fiber bundles as the data model because they are inclusive enough to express all the types of data described in section 2.2.

### 3.1.1 Variables in Fiber Space $F$

To formalize the structure of the data components, we use notation introduced by Spivak [34] that binds the components of the fiber to variable names. This allows us to describe the components in a schema like way. Spivak constructs a set  $\mathbb{U}$  that is the disjoint union of all possible objects of types  $\{T_0, \dots, T_m\} \in \mathbf{DT}$ , where  $\mathbf{DT}$  are the data types of the variables in the dataset. He then defines the single variable set  $\mathbb{U}_\sigma$

$$\begin{array}{ccc} \mathbb{U}_\sigma & \longrightarrow & \mathbb{U} \\ \pi_\sigma \downarrow & & \downarrow \pi \\ C & \xrightarrow{\sigma} & \mathbf{DT} \end{array} \quad (4)$$

which is  $\mathbb{U}$  restricted to objects of type  $T$  bound to variable name  $c$ . The  $\mathbb{U}_\sigma$  lookup is by name to specify that every component is distinct, since multiple components can have the same type  $T$ . Given  $\sigma$ , the fiber for a one variable dataset is

$$F = \mathbb{U}_{\sigma(c)} = \mathbb{U}_T \quad (5)$$

where  $\sigma$  is the schema binding variable name  $c$  to its datatype  $T$ . A dataset with multiple variables has a fiber that is the cartesian cross product of  $\mathbb{U}_\sigma$  applied to all the columns:

$$F = \mathbb{U}_{\sigma(c_1)} \times \dots \times \mathbb{U}_{\sigma(c_i)} \times \dots \times \mathbb{U}_{\sigma(c_n)} \quad (6)$$



which is equivalent to

$$F = F_0 \times \dots \times F_i \times \dots \times F_n \quad (7)$$

276 which allows us to decouple  $F$  into components  $F_i$ .

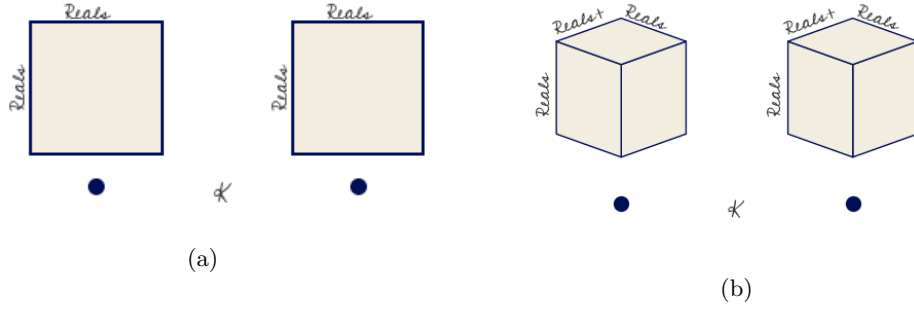


Figure 5: These two datasets have the same base space  $K$  of discrete points, but figure 5a has fiber  $F = \mathbb{R} \times \mathbb{R}$  which is (time, temperature) while figure 5b has fiber  $\mathbb{R} \times \mathbb{R}^+ \times \mathbb{R}$  which is (time, wind=(speed, direction))

For example, the data in figure 5a is a pair of times and °K temperature measurements taken at those times. Time is a positive number of type `datetime` which can be resolved to floats  $\mathbb{U}_{\text{datetime}} = \mathbb{R}$ . Temperature values are real positive numbers  $\mathbb{U}_{\text{float}} = \mathbb{R}^+$ . The fiber is

$$\mathbb{U} = \mathbb{R} \times \mathbb{R}^+ \quad (8)$$

where the first component  $F_0$  is the set of values specified by  $(c = \text{time}, T = \text{datetime}, \mathbb{U}_\sigma = \mathbb{R})$  and  $F_1$  is specified by  $(c = \text{temperature}, T = \text{float}, \mathbb{U}_\sigma = \mathbb{R}^+)$  and is the set of values  $\mathbb{U}_\sigma = \mathbb{R}^+$ . In figure 5b, temperature is replaced with wind. This wind variable is of type `wind` and has two components speed and direction  $\{(s, d) \in \mathbb{R}^2 \mid 0 \leq s, 0 \leq d \leq 360\}$ . Therefore, the fiber is

$$F = \mathbb{R}^+ \times \mathbb{R}^2 \quad (9)$$

277 such that  $F_1$  is specified by  $(c = \text{wind}, T = \text{wind}, \mathbb{U}_\sigma = \mathbb{R}^2)$ . As illustrated in figure 5,  
 278 Spivak's framework provides a consistent way to describe potentially complex components  
 279 of the input data.

### 3.1.2 Measurement Scales: Monoid Actions

Implementing expressive visual encodings requires formally describing the structure on the components of the fiber, which we define by the actions of a monoid on the component. In doing so, we specify the properties of the component that must be preserved in a graphic representation. While structure on a set of values is often described algebraically as operations or through the actions of a group, for example Steven’s scales [59], we generalize to monoids to support more component types. Monoids are also commonly found in functional programming because they specify compositions of transformations [60, 61].

A monoid [62]  $M$  is a set with an associative binary operator  $*$  :  $M \times M \rightarrow M$ . A monoid has an identity element  $e \in M$  such that  $e * a = a * e = a$  for all  $a \in M$ . As defined on a component of  $F$ , a left monoid action [63, 64] of  $M_i$  is a set  $F_i$  with an action  $\bullet : M \times F_i \rightarrow F_i$  with the properties:

**associativity** for all  $f, g \in M_i$  and  $x \in F_i$ ,  $f \bullet (g \bullet x) = (f * g) \bullet x$

**identity** for all  $x \in F_i, e \in M_i$ ,  $e \bullet x = x$

As with the fiber  $F$  the total monoid space  $M$  is the cartesian product

$$M = M_0 \times \dots \times M_i \times \dots \times M_n \quad (10)$$

of each monoid  $M_i$  on  $F_i$ . The monoid is also added to the specification of the fiber  $(c_i, T_i, \mathbb{U}_\sigma M_i)$

Steven’s described the measurement scales[59, 65] in terms of the monoid actions on the measurements: nominal data is permutable, ordinal data is monotonic, interval data is translatable, and ratio data is scalable [66]. For example, given an arbitrary interval scale fiber component  $(c = \text{temperature}, T = \text{float}, \mathbb{U}_\sigma = \mathbb{R})$  with with arbitrary monoid translation actions chosen for this example:

- monoid operator addition  $*$  = +
- monoid operations:  $f : x \mapsto x + 1^\circ\text{C}$ ,  $g : x \mapsto x + 2^\circ\text{C}$

297 • monoid action operator composition • = ◦

By structure preservation, we mean that monoid actions are composable. For the translation actions described above on the temperature fiber, this means that they satisfy the condition

$$\begin{array}{ccc}
 \mathbb{R} & & \\
 x+1^\circ \downarrow & \searrow (x+1^\circ C) \circ (x+2^\circ C) & \\
 \mathbb{R} & \xrightarrow{x+2^\circ C} & \mathbb{R}
 \end{array} \tag{11}$$

298 where  $1^\circ C$  and  $2^\circ C$  are valid distances between two temperatures  $x$ . What this diagram  
 299 means is that either the fiber could be shifted by  $1^\circ C$  (vertical line) then by  $2^\circ C$  (horizontal),  
 300 or the two shifts could be combined such that in this case the fiber is shifted by  $3^\circ C$  (diagonal)  
 301 and these two paths yield the same temperature.

302 While many component types will be one of the measurement scale types, we gen-  
 303 eralize to monoids specifically for the case of partially ordered set. Given a set  $W =$   
 304  $\{mist, drizzle, rain\}$ , then the map  $f : W \rightarrow W$  defined by

305 1.  $f(rain) = drizzle$ ,

306 2.  $f(drizzle) = mist$

307 3.  $f(mist) = mist$

308 is order preserving such that  $mist \leq drizzle \leq rain$  but has no inverse since  $drizzle$  and  
 309  $mist$  go to the same value  $mist$ . Therefore order preserving maps do not form a group, and  
 310 instead we generalize to monoids to support partial order component types. Defining the  
 311 monoid actions on the components serves as the basis for identifying the invariance[4] that  
 312 must be preserved in the visual representation of the component. We propose equivariance  
 313 of monoid actions individually on the fiber to visual component maps and on the graphic  
 314 as a whole.

### 315 3.1.3 Continuity of the Data $K$

316 The base space  $K$  is way to express how the records in  $E$  are connected to each other, for  
 317 example if they are discrete points or if they lie in a 2D continous surface. Connectivity

318 type is assumed in the choice of visualization, for example a line plot implies 1D continuous  
 319 data, but an explicit representation allows for verifying that the topology of the graphic  
 320 representation is equivalent to the topology of the data.



Figure 6: The topological base space  $K$  encodes the connectivity of the data space, for example if the data is independent points or on a plane or a sphere

321 As illustrated in figure 6,  $K$  is akin to an indexing space into  $E$  that describes the  
 322 structure of  $E$ .  $K$  can have any number of dimensions and can be continuous or discrete.

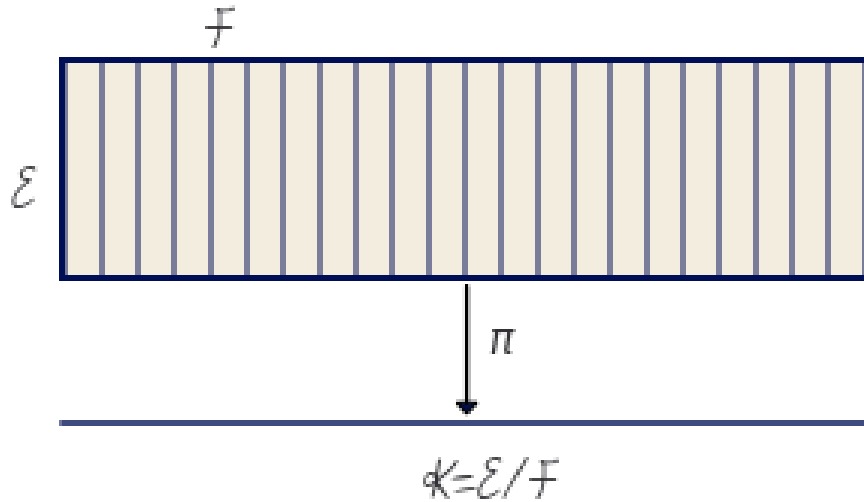


Figure 7: The base space  $E$  is divided into fiber segments  $F$ . The base space  $K$  acts as an index into the records in the fibers.

Formally  $K$  is the quotient space [67] of  $E$  meaning it is the finest space[68] such that every  $k \in K$  has a corresponding fiber  $F_k$ [67]. In figure 7,  $E$  is a rectangle divided by vertical fibers  $F$ , so the minimal  $K$  for which there is always a mapping  $\pi : E \rightarrow K$  is the closed interval  $[0, 1]$ . As with fibers and monoids, we can decompose the total space into components  $\pi : E_i \rightarrow K$  where

$$\pi : E_1 \oplus \dots \oplus E_i \oplus \dots \oplus E_n \rightarrow K \quad (12)$$

323 which is a decomposition of  $F$ . The  $K$  remains the same because the connectivity of records  
 324 does not change just because there are fewer elements in each record.

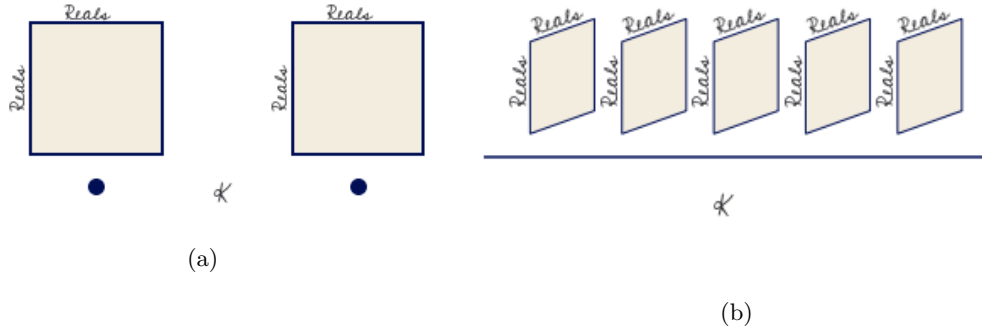


Figure 8: These two datasets have the same (time, temperature) fiber. In figure 8a the total space  $E$  is discrete over points  $k \in K$ , meaning the records in the fiber are also discrete. In figure 8b  $E$  lies over the continuous interval  $K$ , meaning the records in the fiber are sampled from a continuous space.

325 The datasets in figure 8 have the same fiber of (temperature, time). In figure 8a the  
 326 fibers lie over discrete  $K$  such that the records in the datasets in the fiber bundles are  
 327 discrete. The same fiber in figure 8b lies over a continuous interval  $K$  such that the records  
 328 are samples from a continuous function defined on  $K$ . By encoding this continuity in the  
 329 model as  $K$  the data model now explicitly carries information about its structure such  
 330 that the implicit assumptions of the visualization algorithms are now explicit. The explicit

331 topology is a concise way of distinguishing visualizations that appear identical, for example  
 332 heatmaps and images.

#### 333 3.1.4 Data $\tau$

While the projection function  $\pi : E \rightarrow K$  ties together the base space  $K$  with the fiber  $F$ , a section  $\tau : K \rightarrow E$  encodes a dataset. A section function takes as input location  $k \in K$  and returns a record  $r \in E$ . For example, in the special case of a table [34],  $K$  is a set of row ids,  $F$  is the columns, and the section  $\tau$  returns the record  $r$  at a given key in  $K$ . For any fiber bundle, there exists a map

$$\begin{array}{ccc} F & \hookrightarrow & E \\ & & \pi \downarrow \nearrow \tau \\ & & K \end{array} \quad (13)$$

such that  $\pi(\tau(k)) = k$ . The set of all global sections is denoted as  $\Gamma(E)$ . Assuming a trivial fiber bundle  $E = K \times F$ , the section is

$$\tau(k) = (k, (g_{F_0}(k), \dots, g_{F_n}(k))) \quad (14)$$

where  $g : K \rightarrow F$  is the index function into the fiber. This formulation of the section also holds on locally trivial sections of a non-trivial fiber bundle. Because we can decompose the bundle and the fiber, we can decompose  $\tau$  as

$$\tau = (\tau_0, \dots, \tau_i, \dots, \tau_n) \quad (15)$$

334 where each section  $\tau_i$  is a variable or set of variables. This allows for accessing the data  
 335 component wise in addition to accessing the data in terms of its location over  $K$ .

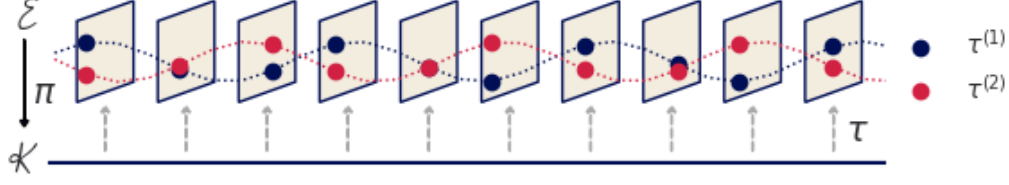


Figure 9: Fiber (time, temperature) with an interval  $K$  basespace. The sections  $\tau^{(1)}$  and  $\tau^{(2)}$  are constrained such that the time variable must be monotonic, which means each section is a timeseries of temperature values. They are included in the global set of sections  $\tau^{(1)}, \tau^{(2)} \in \Gamma(E)$

In the example in figure 9, the fiber is (time, temperature) as described in figure 5 and the base space is the interval  $K$ . The section  $\tau^{(1)}$  resolves to a series of monotonically increasing in time records of (time, temperature) values. Section  $\tau^{(2)}$  returns a different timeseries of (time, temperature) values. Both sections are included in the global set of sections  $\tau^{(1)}, \tau^{(2)} \in \Gamma(E)$ .

### 3.1.5 Applications to Data Containers

This model provides a common formalism for widely used data containers without sacrificing the semantic structure embedded in each container. For example, the section can be any instance of a univariate numpy array[69] that stores an image. This could be a section of a fiber bundle where  $K$  is a 2D continuous plane and the  $F$  is  $(\mathbb{R}^3, \mathbb{R}, \mathbb{R})$  where  $\mathbb{R}^3$  is color, and the other two components are the x and y positions of the sampled data in the image. This position information is already implicitly encoded in the array as the index and the resolution of the image being stored. Instead of an image, the numpy array could also store a 2D discrete table. The fiber would not change, but the  $K$  would now be 0D discrete points. These different choices in topology indicate, for example, what sorts of interpolation would be appropriate when visualizing the data.

There are also many types of labeled containers that can richly be described in this framework because of the schema like structure of the fiber. For example, a pandas series which stores a labeled list, or a dataframe[70] which stores a relational table. A series could

store the values of  $\tau^{(1)}$  and a second series could be  $\tau^{(2)}$ . We could also fatten the fiber to hold two temperature series, such that a section would be an instance of a dataframe with a time column and two temperature columns. While the series and dataframe explicitly have a time index column, they are components in our model and the index is assumed to be data independent references such as hashvalues, virtual memory locations, or random number keys.

Where this model particularly shines are N dimensional labeled data structures. For example, an xarray[71] data that stores temperature field could have a  $K$  that is a continuous volume and the components would be the temperature and the time, latitude, and longitude the measurements were sampled at. A section can also be an instance of a distributed data container, such as a dask array [72]. As with the other containers,  $K$  and  $F$  are defined in terms of the index and dtypes of the components of the array. Because our framework is defined in terms of the fiber, continuity, and sections, rather than the exact values of the data, our model does not need to know what the exact values are until the renderer needs to fill in the image.

## 3.2 Graphic Space $H$

We introduce a graphic bundle to hold the essential information necessary to render a graphical design constructed by the artist. As with the data, we can represent the target graphic as a section  $\rho$  of a bundle  $(H, S, \pi, D)$ . The graphic bundle  $H$  consists of a base  $S$  (3.2.1) that is a thickened form of  $K$  a fiber  $D$  (3.2.2) that is an idealized display space, and sections  $\rho$  (3.2.3) that encode a graphic where the visual characteristics are fully specified.

### 3.2.1 Idealized Display $D$

To fully specify the visual characteristics of the image, we construct a fiber  $D$  that is an infinite resolution version of the target space. Typically  $H$  is trivial and therefore sections can be thought of as mappings into  $D$ . In this work, we assume a 2D opaque image  $D = \mathbb{R}^5$  with elements

$$(x, y, r, g, b) \in D \tag{16}$$



such that a rendered graphic only consists of 2D position and color. To support overplotting and transparency, the fiber could be  $D = \mathbb{R}^7$  such that  $(x, y, z, r, g, b, a) \in D$  specifies the target display. By abstracting the target display space as  $D$ , the model can support different targets, such as a 2D screen or 3D printer.

### 3.2.2 Continuity of the Graphic $S$

Just as the  $K$  encodes the connectivity of the records in the data, we propose an equivalent  $S$  that encodes the connectivity of the rendered elements of the graphic. For example, consider a  $S$  that is mapped to the region of a 2D display space that represents  $K$ . For some visualizations,  $K$  may be lower dimension than  $S$ . For example, a point that is 0D in  $K$  cannot be represented on screen unless it is thickened to 2D to encode the connectivity of the pixels that visually represent the point. This thickening is often not necessary when the dimensionality of  $K$  matches the dimensionality of the target space, for example if  $K$  is 2D and the display is a 2D screen. We introduce  $S$  to thicken  $K$  in a way which preserves the structure of  $K$ .

Formally, we require that  $K$  be a deformation retract[73] of  $S$  so that  $K$  and  $S$  have the same homotopy. The surjective map  $\xi : S \rightarrow K$

$$\begin{array}{ccc} E & & H \\ \pi \downarrow & & \pi \downarrow \\ K & \xleftarrow{\xi} & S \end{array} \quad (17)$$

goes from region  $s \in S_k$  to its associated point  $s$ . This means that if  $\xi(s) = k$ , the record at  $k$  is copied over the region  $s$  such that  $\tau(k) = \xi^* \tau(s)$  where  $\xi^* \tau(s)$  is  $\tau$  pulled back over  $S$ .

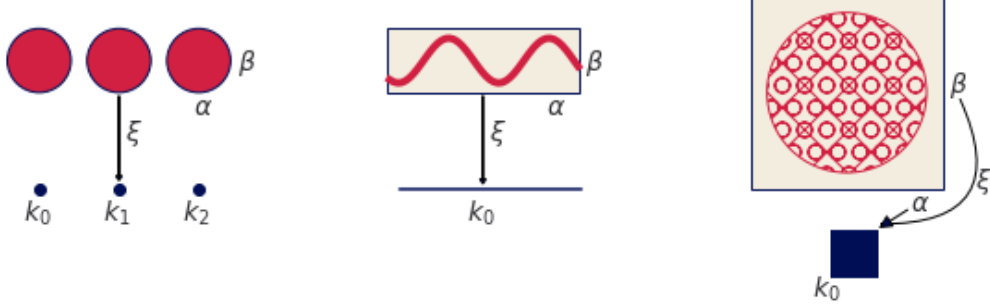


Figure 10: The scatter and line graphic base spaces have one more dimension of continuity than  $K$  so that  $S$  can encode physical aspects of the glyph, such as shape (a circle) or thickness. The image has the same dimension in  $S$  as in  $K$ .

When  $K$  is discrete points and the graphic is a scatter plot, each point  $k \in K$  corresponds to a 2D disk  $S_k$  as shown in figure 10. In the case of 1D continuous data and a line plot, the region  $\beta$  over a point  $\alpha_i$  specifies the thickness of the line in  $S$  for the corresponding  $\tau$  on  $k$ . The image has the same dimensions in data space and graphic space such that no extra dimensions are needed in  $S$ .

The mapping function  $\xi$  provides a way to identify the part of the visual transformation that is specific to the connectivity of the data rather than the values; for example it is common to flip a matrix when displaying an image. The  $\xi$  mapping is also used by interactive visualization components to look up the data associated with a region on screen. One example is to fill in details in a hover tooltip, another is to convert region selection (such as zooming) on  $S$  to a query on the data to access the corresponding record components on  $K$ .

405 **3.2.3 Graphic  $\rho$**

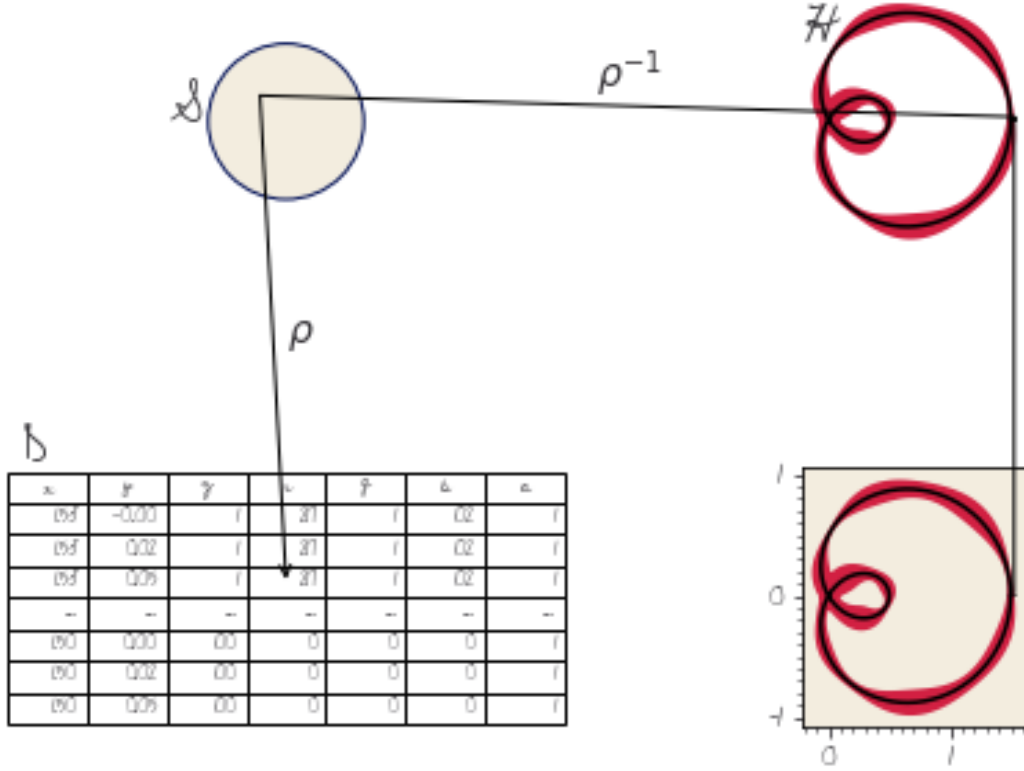


Figure 11: To render a graphic, a pixel  $p$  is selected in the display space, which is defined in the same coordinates as the  $x$  and  $y$  components in  $D$ . The inverse mapping  $\rho_{xy}(p)$  returns a region  $S_p \subset S$ .  $\rho(S_p)$  returns the list of elements  $(x, y, r, g, b) \in D$  that lie over  $S_p$ . The integral over the  $(r, g, b)$  elements is the color of the pixel.

406 This section describes how we go from a graphic in an idealized prerender space to a rendered  
 407 image, where the graphic is the section  $\rho : S \rightarrow H$ . It is sufficient to sketch out how an  
 408 arbitrary pixel would be rendered, where a pixel  $p$  in a real display corresponds to a region  
 409  $S_p$  in the idealized display. To determine the color of the pixel, we aggregate the color values  
 410 over the region via integration.

411 For a 2D screen, the pixel is defined as a region  $p = [y_{top}, y_{bottom}, x_{right}, x_{left}]$  of the  
 412 rendered graphic. Since the  $x$  and  $y$  in  $p$  are in the same coordinate system as the  $x$  and  $y$

413 components of  $D$  the inverse map of the bounding box  $S_p = \rho_{xy}^{-1}(p)$  is a region  $S_p \subset S$ .  
 414 To compute the color, we integrate on  $S_p$

$$r_p = \iint_{S_p} \rho_r(s) ds^2 \quad (18)$$

$$g_p = \iint_{S_p} \rho_g(s) ds^2 \quad (19)$$

$$b_p = \iint_{S_p} \rho_b(s) ds^2 \quad (20)$$

415 As shown in figure 11, a pixel  $p$  in the output space is selected and inverse mapped into  
 416 the corresponding region  $S_p \subset S$ . This triggers a lookup of the  $\rho$  over the region  $S_p$ , which  
 417 yields the set of elements in  $D$  that specify the  $(r, g, b)$  values corresponding to the region  
 418  $p$ . The color of the pixel is then obtained by taking the integral of  $\rho_{rgb}(S_p)$ .

419 In general,  $\rho$  is an abstraction of rendering. In very broad strokes  $\rho$  can be a specification  
 420 such as PDF[74], SVG[75], or an OpenGL scene graph[76]. Alternatively,  $\rho$  can be a rendering  
 421 engine such as cairo[77] or AGG[78]. Implementation of  $\rho$  is out of scope for this work,

### 422 3.3 Artist

We propose that the transformation from data to visual representation can be described as a structure preserving map from one topological space to another. We name this map the artist as that is the analogous part of the Matplotlib[79] architecture that builds visual elements. The topological artist  $A$  is a monoid equivariant sheaf map from the sheaf on a data bundle  $E$  which is  $\mathcal{O}(E)$  to the sheaf on the graphic bundle  $H$ ,  $\mathcal{O}(H)$ .

$$A : \mathcal{O}(E) \rightarrow \mathcal{O}(H) \quad (21)$$

423 Sheafs are a mathematical object with restriction maps that define how to glue  $\tau$  over local  
 424 neighborhoods  $U \subseteq K$ , discussed in section ??, such that the  $A$  maps are consistent over  
 425 continuous regions of  $K$ . While  $A$  can usually construct graphical elements solely with the

data in  $\tau$ , some visualizations, such as line, may also need some finite number  $n$  of derivatives, which is captured by the jet bundle  $\mathcal{J}^n$  [80, 81] with  $\mathcal{J}^0(E) = E$ . In this work, we at most need  $\mathcal{J}^2(E)$  which is the value at  $\tau$  and its first and second derivatives; therefore the artist takes as input the jet bundle  $E' = \mathcal{J}^2(E)$ .

Specifically,  $A$  is the equivariant map from  $E'$  to a specific graphic  $\rho \in \Gamma(H)$

$$\begin{array}{ccccc}
 E' & \xrightarrow{\nu} & V & \xleftarrow{\xi^*} & \xi^*V & \xrightarrow{Q} & H \\
 & \searrow \pi & \downarrow \pi & & \xi^* \pi \downarrow & \swarrow \pi & \\
 & & K & \xleftarrow{\xi} & S & & 
 \end{array} \tag{22}$$

where the input can be point wise  $\tau(k) \mid k \in K$ . The encoders  $\nu : E' \rightarrow V$  convert the data components to visual components(3.2.2). The continuity map  $\xi : S \rightarrow K$  then pulls back the visual bundle  $V$  over  $S$ (3.3.2). Then the assembly function  $Q : \xi^*V \rightarrow H$  composites the fiber components of  $\xi^*V$  into a graphic in  $H$ (3.3.3). This functional decomposition of the visualization artist facilitates building reusable components at each stage of the transformation because the equivariance constraints are defined on  $\nu$ ,  $Q$ , and  $\xi$ .

### 3.3.1 Visual Fiber Bundle $V$

We introduce a visual bundle  $V$  to store the visual representations the artist needs to assemble into a graphic. The visual bundle  $(V, K, \pi, P)$  has section  $\mu : V \rightarrow K$  that resolves to a visual variable in the fiber  $P$ . The visual bundle  $V$  is the latent space of possible parameters of a visualization type, such as a scatter or line plot. We define  $P$  in terms of the parameters of a visualization libraries compositing functions; for example table 1 is a sample of the fiber space for Matplotlib [2].

$\nu_i$	$\mu_i$	$\text{codomain}(\nu_i) \subset P_i$
position	x, y, z, theta, r	$\mathbb{R}$
size	linewidth, markersize	$\mathbb{R}^+$
shape	markerstyle	$\{f_0, \dots, f_n\}$
color	color, facecolor, markerfacecolor, edgecolor	$\mathbb{R}^4$
texture	hatch	$\mathbb{N}^{10}$
	linestyle	$(\mathbb{R}, \mathbb{R}^{+^{n, n\%2=0}})$

Table 1: Some possible components of the fiber  $P$  for a visualization function implemented in Matplotlib

444 A section  $\mu$  is a tuple of visual values that specifies the visual characteristics of a part of  
 445 the graphic. For example, given a fiber of  $\{xpos, ypos, color\}$  one possible section could be  
 446  $\{.5, .5, (255, 20, 147)\}$ . The  $\text{codomain}(\nu_i)$  determines the monoid actions on  $P_i$ . These fiber  
 447 components are implicit in the library, by making them explicit as components of the fiber  
 448 we can build consistent definitions and expectations of how these parameters behave.

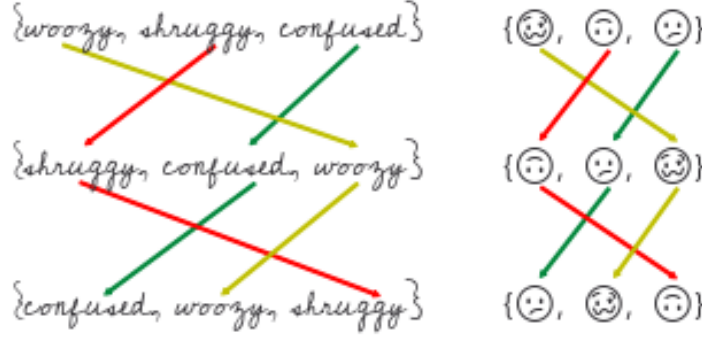


Figure 12: In this artist,  $\nu$  maps the strings to the emojis. This  $\nu$  is equivariant because the monoid actions (which are represented by the colored arrows) are the same on both the  $\tau$  input and  $\mu$  output sets.

As introduced in section 2.3, there are many ways to visually represent data components. We define the visual transformers  $\nu$

$$\{\nu_0, \dots, \nu_n\} : \{\tau_0, \dots, \tau_n\} \mapsto \{\mu_0, \dots, \mu_n\} \quad (23)$$

450 as the set of equivariant maps  $\nu_i : \tau_i \mapsto \mu_i$ . Given  $M_i$  is the monoid action on  $E_i$  and that  
 451 there is a monoid  $M_i'$  on  $V_i$ , then there is a monoid homomorphism from  $\varphi : M_i \rightarrow M_i'$   
 452 that  $\nu$  must preserve. As mentioned in section 3.1.2, we choose monoid actions as the basis  
 453 for equivariance because they define the structure on the fiber components.

A validly constructed  $\nu$  is one where the diagram of the monoid transform  $m$  commutes such that

$$\begin{array}{ccc} E_i & \xrightarrow{\nu_i} & V_i \\ m_r \downarrow & & \downarrow m_v \\ E_i & \xrightarrow{\nu_i} & V_i \end{array} \quad (24)$$

In general, the data fiber  $F_i$  cannot be assumed to be of the same type as the visual fiber  $P_i$  and the actions of  $M$  on  $F_i$  cannot be assumed to be the same as the actions of  $M'$  on  $P$ ; therefore an equivariant  $\nu_i$  must satisfy the constraint

$$\nu_i(m_r(E_i)) = \varphi(m_r)(\nu_i(E_i)) \quad (25)$$

such that  $\varphi$  maps a monoid action on data to a monoid action on visual elements. However, we can construct a monoid action of  $M$  on  $P_i$  that is compatible with a monoid action of  $M$  on  $F_i$ . We can compose the monoid actions on the visual fiber  $M' \times P_i \rightarrow P_i$  with the homomorphism  $\varphi$  that takes  $M$  to  $M'$ . This allows us to define a monoid action on  $P$  of  $M$  that is  $(m, v) \rightarrow \varphi(m) \bullet v$ . Therefore, without a loss of generality, we can assume that an action of  $M$  acts on  $F_i$  and on  $P_i$  compatibly such that  $\varphi$  is the identity function.

On example of an equivariant  $\nu$  is illustrated in figure 12, which is a mapping from **Strings** to symbols. The data is an example of a Steven's nominal measurement set, which is defined as having on it permutation group actions

$$\text{if } r_1 \neq r_2 \text{ then } \nu(r_1) \neq \nu(r_2) \quad (26)$$

such that shuffling the words must have an equivalent shuffle of the symbols they are mapped to. This is illustrated in the identical actions, represented by the colored arrows, on the words and emojis. To preserve ordinal and partial order monoid actions,  $\nu$  must be a monotonic function such that given  $r_1, r_2 \in E_i$ ,

$$\text{if } r_1 \leq r_2 \text{ then } \nu(r_1) \leq \nu(r_2) \quad (27)$$

the visual encodings must also have some sort of ordering. For interval scale data,  $\nu$  is equivariant under translation monoid actions if

$$\nu(x + c) = \nu(x) + c \quad (28)$$



while for ratio data, there must be equivalent scaling

$$\nu(xc) = \nu(x) * c \quad (29)$$

We therefore can test if a  $\nu$  is equivariant by testing the actions under which it must commute. For example, we define a transform  $\nu_i(x) = .5$  on interval data. This means it must commute under translation, for example  $t(x) = x + 2$ . Testing this constraint

$$\nu(t(r + 2)) \stackrel{?}{=} \nu(r) + 2 \quad (30)$$

$$.5 \neq .5 + 2 \quad (31)$$

we find that the  $\nu$  defined here does not commute and is therefore invalid. The constraints on  $\nu$  can be embedded into our artist such that the  $\nu$  functions can test for equivariance and also provide guidance on constructing new  $\nu$  functions.

### 3.3.3 Graphic Assembler $Q$

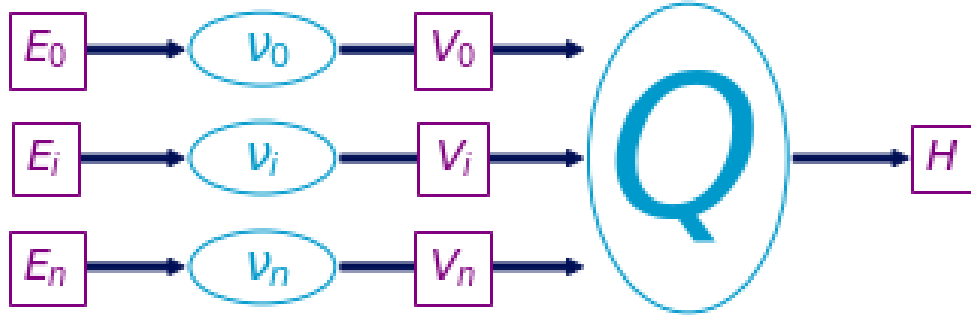


Figure 13:  $\nu_i$  functions convert data  $\tau_i$  to visual characteristics  $\mu_i$ , then  $Q$  assembles  $\mu_i$  into a graphic  $\rho$  such that there is a map  $\xi$  preserving the continuity of the data.  $\rho$  applied to a region of connected components  $S_j$  generates a part of a graphic, for example the point graphical mark.

464 As shown in figure 13, the assembly function  $Q$  combines the fiber  $F_i$  wise  $\nu$  transforms into  
 465 a graphic in  $H$ . Together,  $\nu$  and  $Q$  are a map-reduce operation: map the data into their  
 466 visual encodings, reduce the encodings into a graphic. As with  $\nu$  the constraint on  $Q$  is  
 467 that for every monoid action on the input  $\mu$  there is corresponding monoid action on the  
 468 output  $\rho$ .

While  $\rho$  generates the entire graphic, we will restrict the discussion of  $Q$  to generation of sections of a glyph. We formally describe a glyph as  $Q$  applied to the regions  $k$  that map back to a set of path connected components  $J \subset K$  as input:

$$J = \{j \in K \text{ s. t. } \exists \gamma \text{ s.t. } \gamma(0) = k \text{ and } \gamma(1) = j\} \quad (32)$$

where the path[82]  $\gamma$  from  $k$  to  $j$  is a continuous function from the interval  $[0,1]$ . We define the glyph as the graphic generated by  $Q(S_j)$

$$H \xrightleftharpoons[\rho(S_j)]{} S_j \xrightleftharpoons[\xi^{-1}(J)]{\xi(s)} J_k \quad (33)$$

469 such that for every glyph there is at least one corresponding region on  $K$ . This is in keeping  
 470 with the definition of glyph as any differentiable element put forth by Ziemkiewicz and  
 471 Kosara[43]. The primitive point, line, and area marks[40, 83] are specially cased glyphs.

472 It is on sections of these glyphs that we define the equivariant map as  $Q : \mu \mapsto \rho$  and an  
 473 action on the subset of graphics  $Q(\Gamma(V)) \in \Gamma(H)$  that  $Q$  can generate. We then define the  
 474 constraint on  $Q$  such that if  $Q$  is applied to  $\mu, \mu'$  that generate the same  $\rho$  then the output  
 475 of both sections acted on by the same monoid  $m$  must be the same. While it may seem  
 476 intuitive that visualizations that generate the same glyph should consistently generate the  
 477 same glyph given the same input, we formalize this constraint such that it can be specified  
 478 as part of the implementation of  $Q$ .

Lets call the visual representations of the components  $\Gamma(V) = X$  and the graphic  $Q(\Gamma(V)) = Y$ . If for elements of the monoid  $m \in M$  and for all  $\mu, \mu' \in X$ , we define

the monoid action on  $X$  so that it is by definition equivariant

$$Q(\mu) = Q(\mu') \implies Q(m \circ \mu) = Q(m \circ \mu') \quad (34)$$

479 then a monoid action on  $Y$  can be defined as  $m \circ \rho = \rho'$ . The transformed graphic  $\rho'$  is  
 480 equivariant to a transform on the visual bundle  $\rho' = Q(m \circ \mu)$  on a section that  $\mu \in Q^{-1}(\rho)$   
 481 that must be part of generating  $\rho$ .

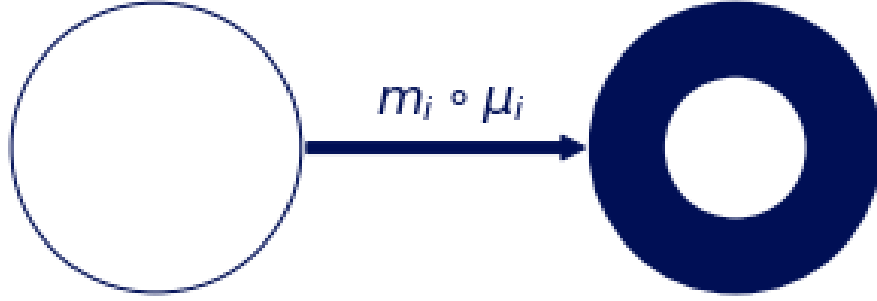


Figure 14: These two glyphs are generated by the same  $Q$  function. The monoid action  $m_i$  on edge thickness  $\mu_i$  of the first glyph yields the thicker edge  $\mu_i'$  in the second glyph.

482 The glyph in figure 14 has the following characteristics  $P$  specified by  $(xpos, ypos, color, thickness)$   
 483 such that one section is  $\mu = (0, 0, 0, 1)$  and  $Q(\mu) = \rho$  generates a piece of the thin hollow  
 484 circle. The equivariance constraint on  $Q$  is that the action  $m = (e, e, e, x + 2)$ , where  $e$  is  
 485 identity, translates  $\mu$  to  $\mu' = (e, e, e, 3)$ . The corresponding action on  $\rho$  causes  $Q(\mu')$  to be  
 486 the thicker circle in figure 14.

#### 487 3.3.4 Assembly $Q$

488 In this section we formulate the minimal  $Q$  that will generate distinguishable graphical  
 489 marks: non-overlapping scatter points, a non-infinitely thin line, and an image.

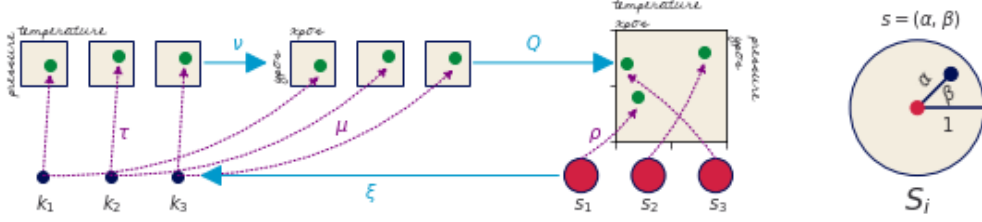


Figure 15: The data is discrete points (temperature, time). Via  $\nu$  these are converted to (xpos, ypos) and pulled over discrete  $S$ . These values are then used to parameterize  $\rho$  which returns a color based on the parameters (xpos,ypos) and position  $\alpha, \beta$  on  $S_k$  that  $\rho$  is evaluated on.

The scatter plot in figure 15 can be defined as  $Q(xpos, ypos)(\alpha, \beta)$  where color  $\rho_{RGB} = (0, 0, 0)$  is defined as part of  $Q$  and  $s = (\alpha, \beta)$  defines the region on  $S$ . The position of this swatch of color can be computed relative to the location on the disc  $S_k$  as shown in figure 15:

$$x = size * \alpha \cos(\beta) + xpos \quad (35)$$

$$y = size * \alpha \sin(\beta) + ypos \quad (36)$$

490 such that  $\rho(s) = (x, y, 0, 0, 0)$  colors the point (x,y) black. Here *size* can either be defined  
 491 inside  $Q$  or it could also be a parameter in  $V$  that is passed along with (xpos, ypos). As  
 492 seen in figure 15, a scatter has a direct mapping from a region on  $S_k$  to its corresponding  $k$ .

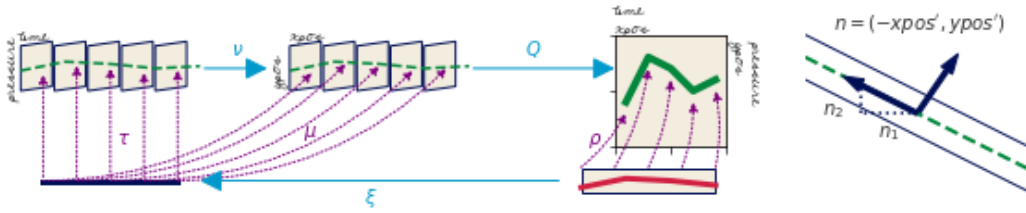


Figure 16: The line fiber (*time*, *temp*) is thickened with the derivative (*time'*, *temperature'*) because that information will be necessary to figure out the tangent to the point to draw a line. This is because the line needs to be pushed perpendicular to the tangent of (xpos, ypos). The data is converted to visual characteristics (xpos, ypos). The  $\alpha$  coordinates on  $S$  specifies the position of the line, the  $\beta$  coordinate specifies thickness.

In contrast to the scatter, the line plot  $Q(xpos, \hat{n}_1, ypos, \hat{n}_2)(\alpha, \beta)$  shown in fig 16 has a  $\xi$  function that is not only parameterized on  $k$  but also on the  $\alpha$  distance along  $k$  and corresponding region in  $S$ . The line also exemplifies the need for the jet since the line needs to know the tangent of the data to draw an envelope above and below each  $(xpos, ypos)$  such that the line appears to have a thickness. The magnitude of the slope is

$$|n| = \sqrt{n_1^2 + n_2^2} \quad (37)$$

such that the normal is

$$\hat{n}_1 = \frac{n_1}{|n|}, \quad \hat{n}_2 = \frac{n_2}{|n|} \quad (38)$$

which yields components of  $\rho$

$$x = xpos(\xi(\alpha)) + width * \beta \hat{n}_1(\xi(\alpha)) \quad (39)$$

$$y = ypos(\xi(\alpha)) + width * \beta \hat{n}_2(\xi(\alpha)) \quad (40)$$

493 where  $(x,y)$  look up the position  $\xi(\alpha)$  on the data. At that point, we also look up the the  
 494 derivatives  $\hat{n}_1, \hat{n}_2$  which are then multiplied by a *width* parameter to specify the thickness.  
 495 As with the *size* parameter in scatter, *width* can be defined in  $Q$  or as a component of  $V$ .

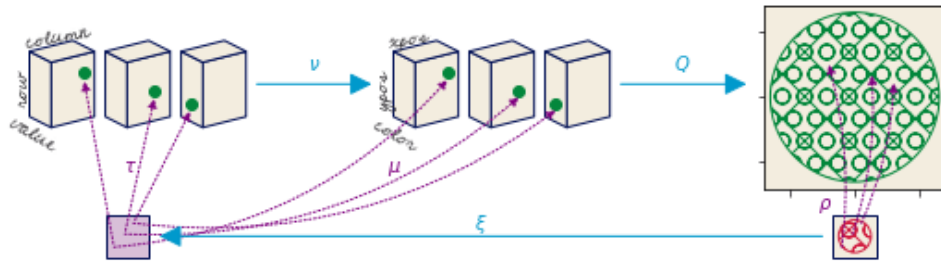


Figure 17: The only visual parameter an image requires is color since  $\xi$  encodes the mapping between position in data and position in graphic.

The image  $Q(xpos, ypos, color)$  in figure 17 is a direct lookup into  $\xi : S \rightarrow K$ . Since  $K$  is 2D continuous space, the indexing variables  $(\alpha, \beta)$  define the distance along the space. This

is then used by  $\xi$  to map into  $K$  to lookup the color values

$$R = R(\xi(\alpha, \beta)) \quad (41)$$

$$G = G(\xi(\alpha, \beta)) \quad (42)$$

$$B = B(\xi(\alpha, \beta)) \quad (43)$$

496 that the data values have been mapped into. In the case of an image, the indexing mapper  
 497  $\xi$  may do some translating to a convention expected by  $Q$ , for example reorientng the array  
 498 such that the first row in the data is at the bottom of the graphic.

### 499 3.3.5 Assembly factory $\hat{Q}$

500 The graphic base space  $S$  is not accessible in many architectures, including Matplotlib;  
 501 instead we can construct a factory function  $\hat{Q}$  over  $K$  that can build a  $Q$ . As shown in  
 502 eq 22,  $Q$  is a bundle map  $Q : \xi^*V \rightarrow H$  where  $\xi^*V$  and  $H$  are both bundles over  $S$ .

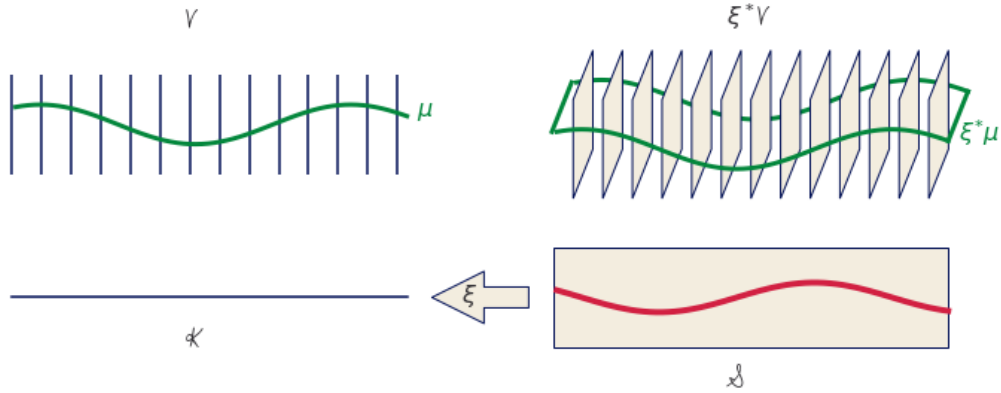


Figure 18: The pullback of the visual bundle  $\xi^*V$  is the replication of a  $\mu$  over all points  $s$  that map back to a single  $k$ . Because the  $\mu$  is the same, we can construct a  $\hat{Q}$  on  $\mu$  over  $k$  that will fabricate the  $Q$  for the equivalent region of  $s$  associated to that  $k$

The preimage of the continuity map  $\xi^{-1}(k) \subset S$  is such that many graphic continuity points  $s \in S_K$  go to one data continuity point  $k$ ; therefore, by definition the pull back of  $\mu$

$$\xi^*V|_{\xi^{-1}(k)} = \xi^{-1}(k) \times P \quad (44)$$

503 copies the visual fiber  $P$  over the the points  $s$  in graphic space  $S$  that correspond to one  $k$   
 504 in data space  $K$ . This set of points  $s$  are the preimage  $\xi^{-1}(k)$  of  $k$ .

505 This copying is illustrated in figure 18, where the 1D fiber  $P \hookrightarrow V$  over  $K$  is copied  
 506 repeatedly to become the 2D fiber  $P^*\mu \hookrightarrow \xi^*V$  with identical components over  $S$ . Given  
 507 the section  $\xi^*\mu$  pulled back from  $\mu$  and the point  $s \in \xi^{-1}(k)$ , there is a direct map from  $\mu$  on  
 508 a point  $k$ , there is a direct map from the visual section over data base space  $(k, \mu(k)) \mapsto$   
 509  $(s, \xi^*\mu(s))$  to the visual section  $\xi^*\mu$  over graphic base space. This map means that the  
 510 pulled back section  $\xi^*\mu(s) = \xi^*(\mu(k))$  is the section  $\mu$  copied over all  $s$ . This means that  
 511  $\xi^*\mu$  is identical for all  $s$  where  $\xi(s) = k$ , which is illustrated in figure 18 as each dot on  $P$  is  
 512 equivalent to the line intersection  $P^*\mu$ .

Given the equivalence between  $\mu$  and  $\xi^*\mu$  defined above, the reliance on  $S$  can be factored out. When  $Q$  maps visual sections into graphics  $Q : \Gamma(\xi^*V) \rightarrow \Gamma(H)$ , if we restrict  $Q$  input to the pulled back visual section  $\xi^*\mu$  then

$$\rho(s) := Q(\xi^*\mu)(s) \quad (45)$$

the graphic section  $\rho$  evaluated on a visual region  $s$  is defined as the assembly function  $Q$  with input pulled back visual section  $\xi^*\mu$  also evaluated on  $s$ . Since the pulled back visual section  $\xi^*\mu$  is the visual section  $\mu$  copied over every graphic region  $s \in \xi^{-1}(k)$ , we can define a  $Q$  factory function

$$\hat{Q}(\mu(k))(s) := Q((\xi^*\mu)(s)) \quad (46)$$

513 where the assembly function  $\hat{Q}$  that takes as input the visual section on data  $\mu$  is defined  
 514 to be the assembly function  $Q$  that takes as input the copied section  $\xi^*\mu$  such that both  
 515 functions are evaluated over the same location  $\xi^{-1}(k) = s$  in the base space  $S$ .

516 Factoring out  $s$  from equation 46 yields  $\hat{Q}(\mu(k)) = Q(\xi^*\mu)$  where  $Q$  is no longer bound  
 517 to input but  $\hat{Q}$  is still defined in terms of  $K$ . In fact,  $\hat{Q}$  is a map from visual space to  
 518 graphic space  $\hat{Q} : \Gamma(V) \rightarrow \Gamma(H)$  locally over  $k$  such that it can be evaluated on a single  
 519 visual record  $\hat{Q} : \Gamma(V_k) \rightarrow \Gamma(H|_{\xi^{-1}(k)})$ . This allows us to construct a  $\hat{Q}$  that only depends  
 520 on  $K$ , such that for each  $\mu(k)$  there is part of  $\rho|_{\xi^{-1}(k)}$ . The construction of  $\hat{Q}$  allows us  
 521 to retain the functional map reduce benefits of  $Q$  without having to majorly restructure  
 522 the existing pipeline for libraries that delegate the construction of  $\rho$  to a back end such as  
 523 Matplotlib.

### 524 3.3.6 Sheafs

The restriction maps of a sheaf describe how local  $\tau$  can be glued into larger sections [84, 85]. As part of the definition of local triviality, there is an open neighborhood  $U \subset K$  for every  $k \in K$ . We can define the inclusion map  $\iota : U \rightarrow K$  which pulls  $E$  over  $U$

$$\begin{array}{ccc} \iota^*E & \xrightarrow{\iota^*} & E \\ \pi \downarrow \Big) \iota^*\tau & & \pi \downarrow \Big) \tau \\ U & \xrightarrow{\iota} & K \end{array} \quad (47)$$

525 such that the pulled back  $\iota^*\tau$  only contains records over  $U \subset K$ . By gluing  $\iota^*\tau$  together, the  
 526 sheaf is putting a continuous structure on local sections which allows for defining a section  
 527 over a subset in  $K$ . That section over subset  $K$  maps to the graphic generated by  $A$  for  
 528 visualizations such as sliding windows[86, 87] streaming data, or navigation techniques such  
 529 as pan and zoom[88].

### 530 3.3.7 Composition of Artists: +

531 To build graphics that are the composites of multiple artists, we define a simple addition  
 532 operator that is the disjoint union of fiber bundles  $E$ . For example, a scatter plot  $E_1$  and a  
 533 line plot  $E_2$  have different  $K$  that are mapped to separate  $S$ . To fully display both graphics,  
 534 the composite graphic  $A_1 + A_2$  needs to include all records on both  $K_1$  and  $K_2$ , which are the  
 535 sections on the disjoint union  $K_1 \sqcup K_2$ . This in turn yields disjoint graphics  $S_1 \sqcup S_2$  rendered



to the same image. Constraints can be placed on the disjoint union such as that the fiber components need to have the same  $\nu$ position encodings or that the position  $\mu$  need to be in a specified range. There is a second type of composition where  $E_1$  and  $E_2$  share a base space  $K_2 \hookrightarrow K_1$  such that the the artists can be considered to be acting on different components of the same section. This type of composition is important for creating visualizations where elements need to update together in a consistent way, such as multiple views [89, 90] and brush-linked views[91, 92].

### 3.3.8 Equivalence class of artists $A'$

It is impractical to implement an artist for every single graphic; instead we implement an approximation of an the equivalence class of artists  $\{A \in A' : A_1 \equiv A_2\}$ . Roughly, equivalent artists have the same fiber bundle  $V$  and same assembly function  $Q$  but act on different sections  $\mu$ , but we will formalize the definition of the equivalence class in future work. As a first pass for implementation purposes, we identify a minimal  $P$  associated with each  $A'$  that defines what visual characteristics of the graphic must originate in the data such that the graphic is identifiable as a given chart type.

For example, a scatter plot of red circles is the output of one artist, a scatter plot of green squares the output of another. These two artists are equivalent since their only difference is in the literal visual encodings (color, shape). Shape and color could also be defined in  $Q$  but the position must come from the fiber  $P = (xpos, ypos)$  since fundamentally a scatter plot is the plotting of one position against another[38]. We also use this criteria to identify derivative types, for example the bubble chart[49] is a type of scatter where by definition the glyph size is mapped from the data. The criteria for equivalence class membership serves as the basis for evaluating invariance[4].

## 559 4 Prototype Implementation: Matplottoy

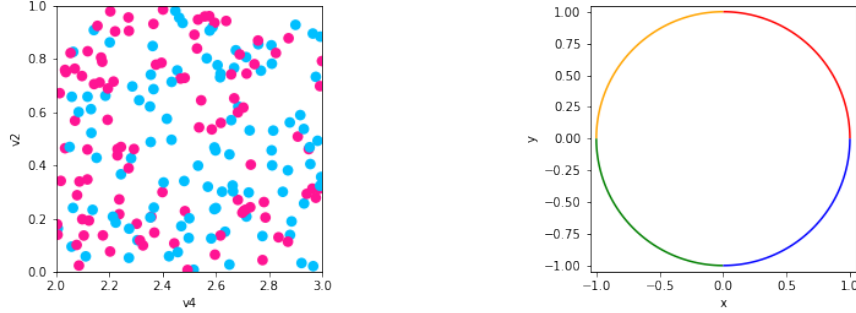


Figure 19: Scatter plot and line plot implemented using prototype artists and data models, building on Matplotlib rendering.

560 To prototype our model, we implemented the artist classes for the scatter and line plots  
 561 shown in figure 19 because they differ in every attribute: different visual channels  $\nu$  that  
 562 composite to different marks  $Q$  with different continuities  $\xi$ . We make use of the Matplotlib  
 563 figure and axes artists [2, 79] so that we can initially focus on the data to graphic transfor-  
 564 mations. We also exploit the Matplotlib transform stack to transform data coordinates into  
 565 screen coordinates. To generate the images in figure 19, we instantiate **fig**, **ax** artists that  
 566 will contain the new **Point**, **Line** primitive objects we implemented based on our topology  
 567 model.

<pre> 1  fig, ax = plt.subplots() 2  artist = Point(data, transforms) 3  ax.add_artist(artist) </pre>	<pre> 1  fig, ax = plt.subplots() 2  artist = Line(data, transforms) 3  ax.add_artist(artist) </pre>
---	--

568 We then add the `Point` and `Line` artist that construct the scatter and line graphics.  
569 These artists are implemented as the equivalence class  $A'$  with the aesthetic configurations  
570 factored out into a `transforms` dictionary that specifies the visual bundle  $V$ . The equivalence  
571 classes  $A'$  map well to Python classes since the functional aspects- $\nu$ ,  $\hat{Q}$ , and  $\xi$ - are completely  
572 reusable in a consistent composition, while the visual values in  $V$  are what change between  
573 different artists belonging to the same class  $A'$ . The `data` object is an abstraction of a  
574 data bundle  $E$  with a specified section  $\tau$ . Implementing  $H$  and  $\rho$  are out of scope for this  
575 prototype because they are part of the rendering process. We also did not implement any  
576 form of  $\xi$  because the scatter, line, and bar plots prototyped here directly broadcast from  $k$   
577 to  $s$ , unlike for example an image which may need to be rotated.

## 578 4.1 Artist Class $A'$

579 The artist is the piece of the Matplotlib architecture that constructs an internal representa-  
580 tion of the graphic that the render then uses to draw the graphic. In the prototype artist,  
581 `transform` is a dictionary of the form `{parameter:(variable, encoder)}` where parame-  
582 ter is a component in  $P$ , variable is a component in  $F$ , and the  $\nu$  encoders are passed in as  
583 functions or callable objects. The data bundle  $E$  is passed in as a `data` object. By binding  
584 data and transforms to  $A'$  inside `__init__`, the `draw` method is a fully specified artist  $A$ .

---

```

1  class ArtistClass(matplotlib.artist.Artist):
2      def __init__(self, data, transforms, *args, **kwargs):
3          # properties that are specific to the graphic but not the channels
4          self.data = data
5          self.transforms = transforms
6          super().__init__(*args, **kwargs)
7
8      def assemble(self, **args):
9          # set the properties of the graphic
10
11     def draw(self, renderer):
12         # returns K, indexed on fiber then key
13         # is passed the
14         view = self.data.view(self.axes)

```

```

15         # visual channel encoding applied fiberwise
16         visual = {p: t['encoder'](view[t['name']])
17                   for p, t in self.transforms.items()}
18         self.assemble(**visual)
19         # pass configurations off to the renderer
20         super().draw(renderer)

```

---

The data is fetched in section  $\tau$  via a `view` method on the data because the input to the artist is a section on  $E$ . The `view` method takes the `axes` attribute because it provides the region in graphic coordinates  $S$  that we can use to query back into data to select a subset as discussed in section 3.3.6. The  $\nu$  functions are then applied to the data to generate the visual section  $\mu$  that here is the object `visual`. The conversion from data to visual space is simplified here to directly show that it is the encoding  $\nu$  applied to the component. In the full implementation, we allow for fixed visual parameter, such as setting a constant color for all sections, by verifying that the named component is in  $F$  before accessing the data. If the data component name is not in  $F$  this is interpreted to mean this component is a thickening of  $V$  that could be pulled back to  $E$  via an inverse identity  $\nu$ .

The components of the visual object, denoted by the Python unpacking convention `**visual` are then passed into the `assemble` function that is  $\hat{Q}$ . This assembly function is responsible for generating a representation such that it could be serialized to recreate a static version of the graphic. Although `assemble` could be implemented outside the class such that it returns an object the artist could then parse to set attributes, the attributes are directly set here to reduce indirection. This artist is not optimized because we prioritized demonstrating the separability of  $\nu$  and  $\hat{Q}$ . The last step in the artist function is handing itself off to the renderer. The extra `*arg`, `**kwargs` arguments in `__init__`, `draw` are artifacts of how these objects are currently implemented in Matplotlib.

The `Point` artist builds on `collection` artists because collections are optimized to efficiently draw a sequence of primitive point and area marks. In this prototype, the scatter marker shape is fixed as a circle, and the only visual fiber components are x and y position, size, and the facecolor of the marker. We only show the `assemble` function here because the `__init__`, `draw` are identical the prototype artist.

---

```

1 class Point(mcollections.Collection):
2     def assemble(self, x, y, s, facecolors='C0' ):
3         # construct geometries of the circle glyphs in visual coordinates
4         self._paths = [mpath.Path.circle(center=(xi,yi), radius=si)
5                         for (xi, yi, si) in zip(x, y, s)]
6         # set attributes of glyphs, these are vectorized
7         # circles and facecolors are lists of the same size
8         self.set_facecolors(facecolors)

```

---

609 The **view** method repackages the data as a fiber component indexed table of vertices. Even  
610 though the **view** is fiber indexed, each vertex at an index  $k$  has corresponding values in  
611 section  $\tau(k_i)$ . This means that all the data on one vertex maps to one glyph. To ensure the  
612 integrity of the section, **view** must be atomic. This means that the values cannot change  
613 after the method is called in draw until a new call in draw. We put this constraint on the  
614 return of the **view** method so that we do not risk race conditions.

615 This table is converted to a table of visual variables and is then passed into **assemble**.  
616 In **assemble**, the  $\mu$  components are used to construct the vector path of each circular  
617 marker with center (x,y) and size x and set the colors of each circle. This is done via the  
618 **Path.circle** object. As mentioned in sections ?? and 3.3.3, this assembly function could  
619 as easily be implemented such that it was fed one  $\tau(k)$  at a time.

620 The main difference between the **Point** and **Line** objects is in the **assemble** function  
621 because line has different continuity from scatter and is represented by a different type of  
622 graphical mark.

---

```

1 class Line(mcollections.LineCollection):
2     def assemble(self, x, y, color='C0'):
3         #assemble line marks as set of segments
4         segments = [np.vstack((vx, vy)).T for vx, vy
5                     in zip(x, y)]
6         self.set_segments(segments)
7         self.set_color(color)

```

---

623 In the **Line** artist, **view** returns a table of edges. Each edge consists of (x,y) points sampled  
624 along the line defined by the edge and information such as the color of the edge. As with  
625 **Point**, the data is then converted into visual variables. In **assemble**, this visual represen-  
626 tation is composed into a set of line segments, where each segment is the array generated  
627 by `np.vstack((vx, vy))`. Then the colors of each line segment are set. The colors are  
628 guaranteed to correspond to the correct segment because of the atomicity constraint on  
629 **view**.

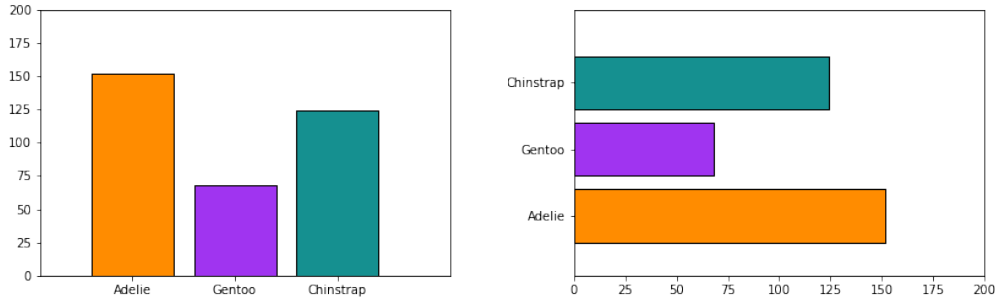


Figure 20: Frequency of Penguin types visualized as discrete bars.

630 The bar charts in figure 20 are generated with a **Bar** artist. The artist has required  
631 visual parameters  $P$  of (position, length), and an additional parameter **orientation** which  
632 controls whether the bars are arranged vertically or horizontally. This parameter only applies  
633 holistically to the graphic and never to individual data parameters, and highlights how the  
634 model encourages explicit differentiation between parameters in  $V$  and graphic parameters  
635 applied directly to  $\hat{Q}$ .

---

```

1 class Bar(mcollections.Collection):
2     def __init__(self, data, transforms, orientation='v', *args, **kwargs):
3         """
4         orientation: str, optional
5             v: bars aligned along x axis, heights on y
6             h: bars aligned along y axis, heights on x
7         """
8         self.orientation = orientation
9         super().__init__(*args, **kwargs)
```

```

10     self.data = data
11     self.transforms = copy.deepcopy(transforms)
12
13     def assemble(self, position, length, floor=0, width=0.8,
14                 facecolors='CO', edgecolors='k', offset=0):
15         #set some defaults
16         width = itertools.repeat(width) if np.isscalar(width) else width
17         floor = itertools.repeat(floor) if np.isscalar(floor) else (floor)
18
19         # offset is passed through via assemblers such as multigroup,
20         # not supposed to be directly tagged to position
21         position = position + offset
22
23     def makeBars(xval, xoff, yval, yoff):
24         return [(x, y), (x, y+yo), (x+xo, y+yo), (x+xo, y), (x, y)]
25         for (x, xo, y, yo) in zip(xval, xoff, yval, yoff)]
26     #build bar glyphs based on graphic parameter
27     if self.orientation in {'vertical', 'v'}:
28         verts = makeBars(position, width, floor, length)
29     elif self.orientation in {'horizontal', 'h'}:
30         verts = makeBars(floor, length, position, width)
31
32     self._paths = [mpath.Path(xy, closed=True) for xy in verts]
33     self.set_edgecolors(edgecolors)
34     self.set_facecolors(facecolors)
35
36     def draw(self, renderer, *args, **kwargs):
37         view = self.data.view(self.axes)
38         visual = {}
39         for (p, t) in self.transforms.items():
40             if isinstance(t, dict):
41                 if t['name'] in self.data.FB.F and 'encoder' in t:
42                     visual[p] = t['encoder'](view[t['name']])
43                 elif 'encoder' in t: # constant value
44                     visual[p] = t['encoder'](t['name'])
45                 elif t['name'] in self.data.FB.F: # identity
46                     visual[p] = view[t['name']]
47             else: # no transform
48                 visual[p] = t

```

```

49         self.assemble(**visual)
50         super().draw(renderer, *args, **kwargs)

```

---

The `draw` method here has a more complex unpacking of visual encodings to support passing in visual component data directly. This is vastly simplifies building composite objects as the alternative would be higher order functions that take as input the transforms passed in by the user. This construction supports a constant visual parameter, an identity transform where the value is the same in  $E$  and  $V$ , and setting the visual component directly. The `assemble` function constructs bars and sets their face and edge colors. The `make_bars` function converts the input position and length to the coordinates of a rectangle of the given width. Defaults are provided for 'width' and 'floor' to make this function more reusable. Typically the defaults are used for the type of chart shown in figure 20, but these visual variables are often set when building composite versions of this chart type as discussed in section 4.4.

## 4.2 Encoders $\nu$

As mentioned above, the encoding dictionary is specified by the visual fiber component, the corresponding data fiber component, and the mapping function. The visual parameter serves as the dictionary key because the visual representation is constructed from the encoding applied to the data  $\mu = \nu \circ \tau$ . For the scatter plot, the mappings for the visual fiber components  $P = (x, y, facecolors, s)$  are defined as

---

```

1 cmap = color.Categorical({'true':'deeppink', 'false':'deepskyblue'})
2 transforms = {'x': {'name': 'v4', 'encoder': lambda x: x},
3               'y': {'name': 'v2', 'encoder': lambda x: x},
4               'facecolors': {'name': 'v3', 'encoder': cmap},
5               's': {'name': None, 'encoder': lambda _: itertools.repeat(.02)}}

```

---

where the position  $(x, y)$   $\nu$  transformers are identity functions. The size  $s$  transformer is not acting on a component of  $F$ , instead it is a  $\nu$  that returns a constant value. While size could be embedded inside the `assemble` function, it is added to the transformers to illustrate user



656 configured visual parameters that could either be constant or mapped to a component in  $F$ .  
 657 The identity and constant  $\nu$  are explicitly implemented here to demonstrate their implicit  
 658 role in the visual pipeline, but they are somewhat optimized away in `Bar`. More complex  
 659 encoders can be implemented as callable classes, such as

---

```

1 class Categorical:
2     def __init__(self, mapping):
3         # check that the conversion is to valid colors
4         assert(mcolors.is_color_like(color) for color in mapping.values())
5         self._mapping = mapping
6
7     def __call__(self, value):
8         # convert value to a color
9         return [mcolors.to_rgba(self._mapping[v]) for v in values]

```

---

660 where `__init__` can validate that the output of the  $\nu$  is a valid element of the  $P$  com-  
 661 ponent the  $\nu$  function is targeting. Creating a callable class also provides a simple way to  
 662 swap out the specific (data, value) mapping without having to reimplement the validation  
 663 or conversion logic. A test for equivariance can be implemented trivially

---

```

1 def test_nominal(values, encoder):
2     m1 = list(zip(values, encoder(values)))
3     random.shuffle(values)
4     m2 = list(zip(values, encoder(values)))
5     assert sorted(m1) == sorted(m2)

```

---

664 but is currently factored out of the artist for clarity. In this example, `is_nominal` checks  
 665 for equivariance of permutation group actions by applying the encoder to a set of values,  
 666 shuffling values, and checking that (value, encoding) pairs remain the same.

### 667 4.3 Data $E$

668 The data input into the `Artist` will often be a wrapper class around an existing data  
 669 structure. This wrapper object must specify the fiber components  $F$  and connectivity  $K$

670 and have a `view` method that returns an atomic object that encapsulates  $\tau$ . The object  
671 returned by the view must be key valued pairs of {component name : component section}  
672 where each section is a component as defined in equation 15. To support specifying the fiber  
673 bundle, we define a `FiberBundle` data class[93]

---

```

1  @dataclass
2  class FiberBundle:
3      """
4      Attributes
5      -----
6      K: {'tables': []}
7      F: {variable name: type}
8      """
9      K: dict
10     F: dict

```

---

674 that asks the user to specify how  $K$  is triangulated and the attributes of  $F$ . Python  
675 dataclasses are a good abstraction for the fiber bundle class because the `FiberBundle` class  
676 only stores data. The  $K$  is specified as tables because the `assemble` functions expect  
677 tables that match the continuity of the graphic; scatter expects a vertex table because it  
678 is discontinuous, line expects an edge table because it is 1D continuous. The fiber informs  
679 appropriate choice of  $\nu$  therefore it is a dictionary of attributes of the fiber components.

680 To generate the scatter plot in figure 19, we fully specify a dataset with random keys  
681 and values in a section chosen at random from the corresponding fiber component. The  
682 fiberbundle `FB` is a class level attribute since all instances of `VertexSimplex` come from the  
683 same fiberbundle.

---

```

1  class VertexSimplex: #maybe change name to something else
2      """Fiberbundle is consistent across all sections
3      """
4      FB = FiberBundle({'tables': ['vertex']},
5                        {'v1': float, 'v2': str, 'v3': float})
6
7      def __init__(self, sid = 45, size=1000, max_key=10**10):

```

---

```

8         # create random list of keys
9     def tau(self, k):
10         # e1 is sampled from F1, e2 from F2, etc...
11         return (k, (e1, e2, e3, e4))
12
13     def view(self, axes):
14         table = defaultdict(list)
15         for k in self.keys:
16             table['index'] = k
17             # on each iteration, add one (name, value) pair per component
18             for (name, value) in zip(self.FB.fiber.keys(), self.tau(k)[1]):
19                 table[name].append(value)
20         return table

```

---

684 The view method returns a dictionary where the key is a fiber component name and the  
 685 value is a list of values in the fiber component. The table is built one call to the section  
 686 method `tau` at a time, guaranteeing that all the fiber component values are over the same  
 687  $k$ . Table has a `get` method as it is a method on Python dictionaries. In contrast, the line  
 688 in `EdgeSimplex` is defined as the functions `_color`, `_xy` on each edge.

---

```

1 class EdgeSimplex:
2
3     FB = FiberBundle({'tables': ['vertex', 'edge']},
4                       {'x' : float, 'y': float,
5                        'color': mtypes.Color()})
6     def __init__(self, num_edges=4, num_samples=1000):
7         self.keys = range(num_edge) #edge id
8         # distance along edge
9         self.distances = np.linspace(0,1, num_samples)
10        # half generalized representation of arcs on a circle
11        self.angle_samples = np.linspace(0, 2*np.pi, len(self.keys)+1)
12
13    @staticmethod
14    def _color(edge):
15        colors = ['red', 'orange', 'green', 'blue']
16        return colors[edge%len(colors)]
17
18    @staticmethod

```

```

19     def _xy(edge, distances, start=0, end=2*np.pi):
20         # start and end are parameterizations b/c really there is
21         angles = (distances *(end-start)) + start
22         return np.cos(angles), np.sin(angles)
23
24     def tau(self, k): #will fix location on page on revision
25         x, y = self._xy(k, self.distances,
26                         self.angle_samples[k], self.angle_samples[k+1])
27         color = self._color(k)
28         return (k, (x, y, color))
29
30     def view(self, axes):
31         table = defaultdict(list)
32         for k in self.keys:
33             table['index'].append(k)
34             # (name, value) pair, value is [x0, ..., xn] for x, y
35             for (name, value) in zip(self.FB.fiber.keys(), self.tau(k, simplex)[1]):
36                 table[name].append(value)

```

---

689 Unlike scatter, the line section method `tau` returns the functions on the edge evaluated on  
 690 the interval  $[0,1]$ . By default these means each `tau` returns a list of 1000 x and y points and  
 691 the associated color. As with scatter, `view` builds a table by calling `tau` for each  $k$  Unlike  
 692 scatter, the line table is a list where each item contains a list of points. This bookkeeping  
 693 of which data is on an edge is used by the `assembly` functions to bind segments to their  
 694 visual properties.

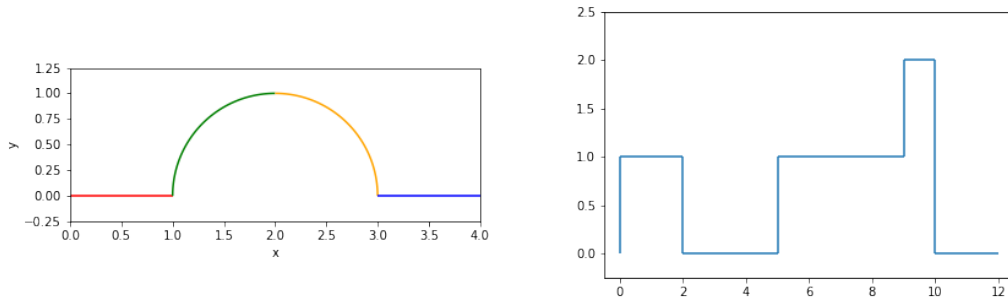


Figure 21: Continuous and discontinuous lines as defined by different data models, but generated with the same `A'Line`

695 The graphics in figure 21 are made using the Line artist and the Graphline data source

---

```

1 class GraphLine:
2     def __init__(self, FB, edge_table, vertex_table, num_samples=1000, connect=False):
3         #s set args as attributes and generate distance
4         if connect: # test connectivity if edges are continuous
5             assert edge_table.keys() == self.FB.F.keys()
6             assert is_continuous(vertex_table)
7
8     def tau(self, k):
9         # evaluates functions defined in edge table
10        return(k, (self.edges[c][k](self.distances) for c in self.FB.F.keys()))
11
12    def view(self, axes):
13        """walk the edge_vertex table to return the edge function
14        """
15        table = defaultdict(list)
16        #sort since intervals lie along number line and are ordered pair neighbors
17        for (i, (start, end)) in sorted(zip(self.ids, self.vertices), key=lambda v:v[1][0]):
18            table['index'].append(i)
19            # same as view for line, returns nested list
20            for (name, value) in zip(self.FB.F.keys(), self.tau(i, simplex)[1]):
21                table[name].append(value)
22    return table

```

---

696 where if told that the data is connected, the data source will check for that connectivity by  
697 constructing an adjacency matrix. The multicolored line is a connected graph of edges with  
698 each edge function evaluated on 1000 samples

---

```

1 simplex.GraphLine(FB, edge_table, vertex_table, connect=True)

```

---

699 while the stair chart is discontinuous and only needs to be evaluated at the edges of the  
700 interval

---

```

1 simplex.GraphLine(FB, edge_table, vertex_table, num_samples=2, connect=False)

```

---

701 such that one advantage of this model is it helps differentiate graphics that have different  
702 artists from graphics that have the same artist but make different assumptions about the  
703 source data.

## 704 4.4 Case Study: Penguins

705 For this case study, we use the Palmer Penguins dataset[94, 95] since it is multivariate and  
706 has a varying number of penguins. We use a version of the data packaged as a pandas  
707 dataframe[96] since that is a very commonly used Python labeled data structure. The  
708 wrapper is very thin because there is explicitly only one section.

---

```
1 class DataFrame:
2     def __init__(self, dataframe):
3         self.FB = FiberBundle(K = {'tables': ['vertex']},
4                                 F = dict(dataframe.dtypes))
5         self._tau = dataframe.iloc
6         self._view = dataframe
7
8     def view(self, axes=None):
9         return self._view
```

---

709 Since the aim for this wrapper is to be very generic, here the fiber is set by querying the  
710 dataframe for its metadata. The `dtypes` are a list of column names and the datatype of  
711 the values in each column; this is the minimal amount of information the model requires to  
712 verify constraints. The pandas indexer is a key valued set of discrete vertices, so there is no  
713 need to repackage for the data interface.

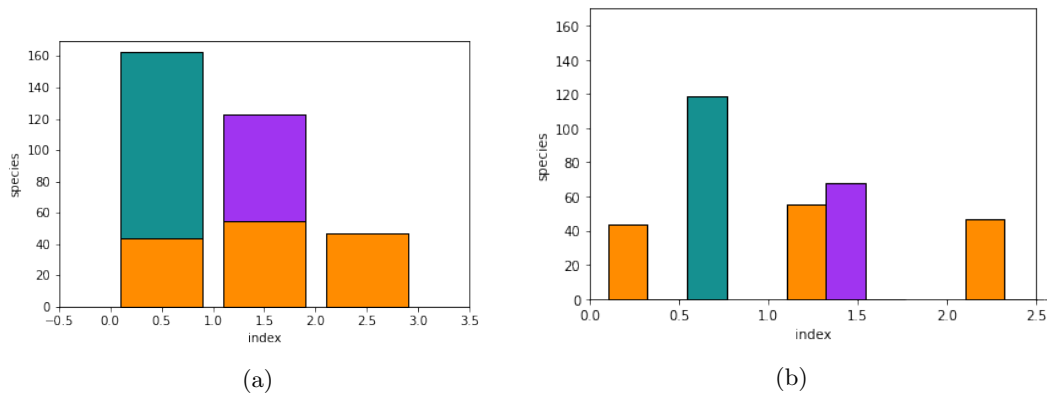


Figure 22: Penguin count disaggregated by island and species

714 The stacked and grouped bar charts in figure 22 are both out of `Bar` artists such that  
 715 the difference between `StackedBar` and `GroupedBar` is specific to the ways in which the  
 716 `Bar` are stitched together. These two artists have identical constructors and `draw` methods.  
 717 As with `Bar`, the orientation is set in the constructor. In both these artists, we separate  
 718 the transforms applied to only one component and the case `mtransforms` where the same  
 719 transform is applied to multiple components such that  $V$  has multiple components that map  
 720 to the same retinal variable.

---

```

1 class StackedBar(martist.Artist):
2     def __init__(self, data, transforms, mtransforms, orientation='v', *args, **kwargs):
3         """
4         Parameters
5         -----
6
7         orientation: str, optional
8             vertical: bars aligned along x axis, heights on y
9             horizontal: bars aligned along y axis, heights on x
10        """
11        super().__init__(*args, **kwargs)
12        self.data = data
13        self.orientation = orientation
14        self.transforms = copy.deepcopy(transforms)
15        self.mtransforms = copy.deepcopy(mtransforms)
16

```

```

17     def assemble(self):
18         view = self.data.view(self.axes)
19         self.children = [] # list of bars to be rendered
20         floor = 0
21         for group in self.mtransforms:
22             # pull out the specific group transforms
23             group['floor'] = floor
24             group.update(self.transforms)
25             bar = Bar(self.data, group, self.orientation, transform=self.axes.transData)
26             self.children.append(bar)
27             floor += view[group['length']]['name']
28
29
30     def draw(self, renderer, *args, **kwargs):
31         # all the visual conversion gets pushed to child artists
32         self.assemble()
33         #self._transform = self.children[0].get_transform()
34         for artist in self.children:
35             artist.draw(renderer, *args, **kwargs)

```

---

721 Since all the visual transformation is passed through to Bar, the draw method does not  
722 do any visual transformations. In StackedBar the view is used to adjust the floor for  
723 every subsequent bar chart since a stacked bar chart is bar chart area marks concatenated  
724 together in the length parameter. In contrast, GroupedBar does not even need the view, but  
725 instead keeps track of the relative position of each group of bars in the visual only variable  
726 offset.

---

```

1 class GroupedBar(martist.Artist):
2     def assemble(self):
3         self.children = [] # list of bars to be rendered
4         ngroups = len(self.mtransforms)
5
6         for gid, group in enumerate(self.mtransforms):
7             group.update(self.transforms)
8             width = group.get('width', .8)
9             group['width'] = width/ngroups
10            group['offset'] = gid/ngroups*width

```



```

11         bar = Bar(self.data, group, self.orientation, transform=self.axes.transData)
12         self.children.append(bar)

```

---

727 Since the only difference between these two glyphs is in the composition of `Bar`, they take  
 728 in the exact same transform specification dictionaries. The `transform` dictionary dictates  
 729 the position of the group, in this case by island the penguins are found on.

---

```

1     transforms = {'position': {'name': 'island',
2                               'encoder': position.Nominal({'Biscoe': 0.1, 'Dream': 1.1, 'Torgersen': 2.1})}}
3     group_transforms = [{'length': {'name': 'Adelie'},
4                             'facecolors': {'name': "Adelie_s", 'encoder': cmap}},
5                          {'length': {'name': 'Chinstrap'},
6                              'facecolors': {'name': "Chinstrap_s", 'encoder': cmap}},
7                          {'length': {'name': 'Gentoo'},
8                              'facecolors': {'name': "Gentoo_s", 'encoder': cmap}}]

```

---

730 `group_transforms` describes the group, and takes a list of dictionaries where each dictionary  
 731 is the aesthetics of each group. That `position` and `length` are required parameters is  
 732 enforced in the creation of the `Bar` artist. These means that these two artists have identical  
 733 function signatures

---

```

1     artistSB = bar.StackedBar(bt, ts, group_transforms)
2     artistGB = bar.GroupedBar(bt, ts, group_transforms)

```

---

734 but differ in assembly  $\hat{Q}$ . By decomposing the architecture into data, visual encoding,  
 735 and assembly steps, we are able to build components that are more flexible and also more self  
 736 contained than the existing code base. While very rough, this API demonstrates that the  
 737 ideas presented in the math framework are implementable. For example, the `draw` function  
 738 that maps most closely to  $A$  is functional, with state only being necessary for bookkeeping  
 739 the many inputs that the function requires. In choosing a functional approach, if not  
 740 implementation, we provide a framework for library developers to build reusable encoder  
 741  $\nu$  assembly  $\hat{Q}$  and artists  $A$ . We argue that if these functions are built such that they

are equivariant with respect to monoid actions and the graphic topology is a deformation retraction of the data topology, then the artist by definition will be a structure and property preserving map from data to graphic.

## 5 Discussion

This work contributes a mathematical description of the mapping  $A$  from data to visual representation. Combining Butler’s proposal of a fiber bundle model of visualization data with Spivak’s formalism of schema lets this mode; support a variety of datasets, including discrete relational tables,, multivariate high resolution spatio temporal datasets, and complex networks. Decomposing the artist into encoding  $\nu$ , assembly  $Q$ , and reindexing  $\xi$  provides the specifications for producing visualization where the structure and properties match those of the input data. These specifications are that the graphic must have continuity equivalent to the data, and that the visual characteristics of the graphics are equivariant to their corresponding components under monoid actions. This model defines these constraints on the transformation function such that they are not specific to any one type of encoding or visual characteristic. Encoding the graphic space as a fiber bundle provides a structure rich abstraction of the target graphical design in the target display space.

The toy prototype built using this model validates that is usable for a general purpose visualization tool since it can be iteratively integrated into the existing architecture rather than starting from scratch. Factoring out glyph formation into assembly functions allows for much more clarity in how the glyphs differ. This prototype demonstrates that this framework can generate the fundamental marks: point (scatter plot), line (line chart), and area (bar chart). Furthermore, the grouped and stacked bar examples demonstrate that this model supports composition of glyphs into more complex graphics. These composite examples also rely on the fiber bundles section base book keeping to keep track of which components contribute to the attributes of the glyph. Implementing this example using a Pandas dataframe demonstrates the ease of incorporating existing widely used data containers rather than requiring users to conform to one stands.

## 769 5.1 Limitations

770 So far this model has only been worked out for a single data set tied to a primitive mark,  
771 but it should be extensible to compositing datasets and complex glyphs. The examples and  
772 prototype have so far only been implemented for the static 2D case, but nothing in the math  
773 limits to 2D and expansion to the animated case should be possible because the model is  
774 formalized in terms of the sheaf. While this model supports equivariance of figurative glyphs  
775 generated from parameters of the data[97, 98], it currently does not have a way to evaluate  
776 the semantic accuracy of the figurative representation. Effectiveness is out of scope for this  
777 model because it is not part of the structure being preserved, but potentially a developer  
778 building a domain specific library with this model could implement effectiveness criteria in  
779 the artists. Also, even though the model is designed to be backend independent, it has only  
780 really been tested against the AGG backend. It is especially unknown how this framework  
781 interfaces with high performance rendering libraries such as OpenGL[76]. Because this model  
782 has been limited to the graphic design space, it does not address the critical task of laying  
783 out the graphics in the image

784 This model and the associated prototype is deeply tied to Matplotlib’s existing archi-  
785 tecture. While the model is expected to generalize to other libraries, such as those built on  
786 Mackinlay’s APT framework, this has not been worked through. In particular, Mackinlay’s  
787 formulation of graphics as a language with semantic and syntax lends itself a declarative in-  
788 terface[7], which Heer and Bostock use to develop a domain specific visualization language  
789 that they argue makes it simpler for designers to construct graphics without sacrificing  
790 expressivity [16]. Similarly, the model presented in this work formulates visualization as  
791 equivariant maps from data space to visual space, and is designed such that developers can  
792 build software libraries with data and graphic topologies tuned to specific domains.

## 793 5.2 Future Work

794 While the model and prototype demonstrate that generation of simple marks from the data,  
795 there is a lot of work left to develop a model that underpins a minimally viable library.

Foremost is implementing a data object that encodes data with a 2D continuous topology and an artist that can consume data with a 2D topology to visualize the image[9, 99, 100] and also encoding a separate heatmap[101, 102] artist that consumes 1D discrete data. A second important proof of concept artist is a boxplot[103] because it is a graphic that assumes computation on the data side and the glyph is built from semantically defined components and a list of outliers. The model supports simple composition of glyphs by overlaying glyphs at the same position, but more work is needed to define an operator where the fiber bundles have shared  $S_2 \hookrightarrow S_1$  such that fibers could be pulled back over the subset. While the model's simple addition supports axes as standalone artists with overlapping visual position encoding, the complex operator would allow for binding together data that needs to be updated together. Additionally, implementing the complex addition operator and explicit graphic to data maps would allow for developing a mathematical formalism and prototype of how interactivity would work in this model. In summary, the proposed scope of work for the dissertation is

- expansion of the mathematical framework to include complex addition
- formalization of definition of equivalence class  $A'$
- implementation of artist with explicit  $\xi$
- specification of interactive visualization
- mathematical formulation of a graphic with axes labeling
- implementation of new prototype artists that do not inherit from Matplotlib artists
- provisional mathematics and implementation of user level composite artists
- proof of concept domain specific user facing library

Other potential tasks for future work is implementing a data object for a non-trivial fiber bundle and exploiting the models section level formalism to build distributed data source models and concurrent artists. This could be pushed further to integrate with topological[104] and functional [105] data analysis methods. Since this model formalizes notions

of structure preservation, it can serve as a good base for tools that assess quality metrics[106] or invariance [4] of visualizations with respect to graphical encoding choices. While this paper formulates visualization in terms of monoidal action homomorphisms between fiberbundles, the model lends itself to a categorical formulation[107, 108] that could be further explored.

## 6 Conclusion

An unofficial philosophy of Matplotlib is to support making whatever kinds of plots a user may want, even if they seem nonsensical to the development team. The topological framework described in this work provides a way to facilitate this graph creation in a rigorous manner; any artist that meets the equivariance criteria described in this work by definition generates a graphic representation that matches the structure of the data being represented. We leave it to domain specialists to define the structure they need to preserve and the maps they want to make, and hopefully make the process easier by untangling these components into separate constrained maps and providing a fairly general data and display model.

## References

- [1] John W. Tukey. “We Need Both Exploratory and Confirmatory”. In: *The American Statistician* 34.1 (Feb. 1980), pp. 23–25. ISSN: 0003-1305. DOI: [10.1080/00031305.1980.10482706](https://doi.org/10.1080/00031305.1980.10482706).
- [2] J. D. Hunter. “Matplotlib: A 2D Graphics Environment”. In: *Computing in Science Engineering* 9.3 (May 2007), pp. 90–95. ISSN: 1558-366X. DOI: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55).
- [3] Jock Mackinlay. “Automating the Design of Graphical Presentations of Relational Information”. In: *ACM Transactions on Graphics* 5.2 (Apr. 1986), pp. 110–141. ISSN: 0730-0301. DOI: [10.1145/22949.22950](https://doi.org/10.1145/22949.22950).

- [4] Gordon Kindlmann and Carlos Scheidegger. “An Algebraic Process for Visualization Design”. In: *IEEE transactions on visualization and computer graphics* 20.12 (2014), pp. 2181–2190.
- [5] T. Sugibuchi, N. Spyrtatos, and E. Siminenko. “A Framework to Analyze Information Visualization Based on the Functional Data Model”. In: *2009 13th International Conference Information Visualisation*. 2009, pp. 18–24. DOI: [10.1109/IV.2009.56](https://doi.org/10.1109/IV.2009.56).
- [6] P. Vickers, J. Faith, and N. Rossiter. “Understanding Visualization: A Formal Approach Using Category Theory and Semiotics”. In: *IEEE Transactions on Visualization and Computer Graphics* 19.6 (June 2013), pp. 1048–1061. ISSN: 1941-0506. DOI: [10.1109/TVCG.2012.294](https://doi.org/10.1109/TVCG.2012.294).
- [7] Kenneth C. Louden. *Programming Languages : Principles and Practice*. English. Pacific Grove, Calif: Brooks/Cole, 2010. ISBN: 978-0-534-95341-6 0-534-95341-7.
- [8] J. Heer and M. Agrawala. “Software Design Patterns for Information Visualization”. In: *IEEE Transactions on Visualization and Computer Graphics* 12.5 (2006), pp. 853–860. DOI: [10.1109/TVCG.2006.178](https://doi.org/10.1109/TVCG.2006.178).
- [9] M. Tory and T. Moller. “Rethinking Visualization: A High-Level Taxonomy”. In: *IEEE Symposium on Information Visualization*. Oct. 2004, pp. 151–158. DOI: [10.1109/INFVIS.2004.59](https://doi.org/10.1109/INFVIS.2004.59).
- [10] Leland Wilkinson. *The Grammar of Graphics*. en. 2nd ed. Statistics and Computing. New York: Springer-Verlag New York, Inc., 2005. ISBN: 978-0-387-24544-7.
- [11] Hadley Wickham. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN: 978-3-319-24277-4.
- [12] M. Bostock and J. Heer. “Protovis: A Graphical Toolkit for Visualization”. In: *IEEE Transactions on Visualization and Computer Graphics* 15.6 (Nov. 2009), pp. 1121–1128. ISSN: 1941-0506. DOI: [10.1109/TVCG.2009.174](https://doi.org/10.1109/TVCG.2009.174).
- [13] M. Bostock, V. Ogievetsky, and J. Heer. “D<sup>3</sup> Data-Driven Documents”. In: *IEEE Transactions on Visualization and Computer Graphics* 17.12 (Dec. 2011), pp. 2301–2309. ISSN: 1941-0506. DOI: [10.1109/TVCG.2011.185](https://doi.org/10.1109/TVCG.2011.185).

- [14] Arvind Satyanarayan, Kanit Wongsuphasawat, and Jeffrey Heer. “Declarative Interaction Design for Data Visualization”. en. In: *Proceedings of the 27th Annual ACM Symposium on User Interface Software and Technology*. Honolulu Hawaii USA: ACM, Oct. 2014, pp. 669–678. ISBN: 978-1-4503-3069-5. DOI: [10.1145/2642918.2647360](https://doi.org/10.1145/2642918.2647360).
- [15] Jacob VanderPlas et al. “Altair: Interactive Statistical Visualizations for Python”. en. In: *Journal of Open Source Software* 3.32 (Dec. 2018), p. 1057. ISSN: 2475-9066. DOI: [10.21105/joss.01057](https://doi.org/10.21105/joss.01057).
- [16] Jeffrey Heer and Michael Bostock. “Declarative Language Design for Interactive Visualization”. In: *IEEE Transactions on Visualization and Computer Graphics* 16.6 (Nov. 2010), pp. 1149–1156. ISSN: 1077-2626. DOI: [10.1109/TVCG.2010.144](https://doi.org/10.1109/TVCG.2010.144).
- [17] Krist Wongsuphasawat. *Navigating the Wide World of Data Visualization Libraries (on the Web)*. 2021.
- [18] Caroline A Schneider, Wayne S Rasband, and Kevin W Eliceiri. “NIH Image to ImageJ: 25 Years of Image Analysis”. In: *Nature Methods* 9.7 (July 2012), pp. 671–675. ISSN: 1548-7105. DOI: [10.1038/nmeth.2089](https://doi.org/10.1038/nmeth.2089).
- [19] Nicholas Sofroniew et al. *Napari/Napari: 0.4.5rc1*. Zenodo. Feb. 2021. DOI: [10.5281/zenodo.4533308](https://doi.org/10.5281/zenodo.4533308).
- [20] *Writing Plugins*. en. [https://imagej.net/Writing\\_plugins](https://imagej.net/Writing_plugins).
- [21] Software Studies. *Culturevis/Imageplot*. Jan. 2021.
- [22] Marcus D. Hanwell et al. “The Visualization Toolkit (VTK): Rewriting the Rendering Code for Modern Graphics Cards”. en. In: *SoftwareX* 1-2 (Sept. 2015), pp. 9–12. ISSN: 23527110. DOI: [10.1016/j.softx.2015.04.001](https://doi.org/10.1016/j.softx.2015.04.001).
- [23] Berk Geveci et al. “VTK”. In: *The Architecture of Open Source Applications* 1 (2012), pp. 387–402.
- [24] P. Ramachandran and G. Varoquaux. “Mayavi: 3D Visualization of Scientific Data”. In: *Computing in Science Engineering* 13.2 (Mar. 2011), pp. 40–51. ISSN: 1558-366X. DOI: [10.1109/MCSE.2011.35](https://doi.org/10.1109/MCSE.2011.35).

- [25] James Ahrens, Berk Geveci, and Charles Law. “Paraview: An End-User Tool for Large Data Visualization”. In: *The visualization handbook* 717.8 (2005).
- [26] Brian Wylie and Jeffrey Baumes. “A Unified Toolkit for Information and Scientific Visualization”. In: *Proc.SPIE*. Vol. 7243. Jan. 2009. DOI: [10.1117/12.805589](https://doi.org/10.1117/12.805589).
- [27] A. Sarikaya et al. “What Do We Talk About When We Talk About Dashboards?”. In: *IEEE Transactions on Visualization and Computer Graphics* 25.1 (Jan. 2019), pp. 682–692. ISSN: 1941-0506. DOI: [10.1109/TVCG.2018.2864903](https://doi.org/10.1109/TVCG.2018.2864903).
- [28] Stephen Few and Perceptual Edge. “Dashboard Confusion Revisited”. In: *Perceptual Edge* (2007), pp. 1–6.
- [29] Z. Pousman, J. Stasko, and M. Mateas. “Casual Information Visualization: Depictions of Data in Everyday Life”. In: *IEEE Transactions on Visualization and Computer Graphics* 13.6 (Nov. 2007), pp. 1145–1152. ISSN: 1941-0506. DOI: [10.1109/TVCG.2007.70541](https://doi.org/10.1109/TVCG.2007.70541).
- [30] *Data Representation in Mayavi — Mayavi 4.7.2 Documentation*. <https://docs.enthought.com/mayavi/mayavi/d>
- [31] D. M. Butler and M. H. Pendley. “A Visualization Model Based on the Mathematics of Fiber Bundles”. en. In: *Computers in Physics* 3.5 (1989), p. 45. ISSN: 08941866. DOI: [10.1063/1.168345](https://doi.org/10.1063/1.168345).
- [32] David M. Butler and Steve Bryson. “Vector-Bundle Classes Form Powerful Tool for Scientific Visualization”. en. In: *Computers in Physics* 6.6 (1992), p. 576. ISSN: 08941866. DOI: [10.1063/1.4823118](https://doi.org/10.1063/1.4823118).
- [33] David I Spivak. *Databases Are Categories*. en. Slides. June 2010.
- [34] David I Spivak. “SIMPLICIAL DATABASES”. en. In: (), p. 35.
- [35] Tamara Munzner. “Ch 2: Data Abstraction”. In: *CPSC547: Information Visualization, Fall 2015-2016* ().
- [36] Connie Malamed. *Information Display Tips*. <https://understandinggraphics.com/visualizations/information-display-tips/>. Blog. Jan. 2010.



- [37] John Krygier and Denis Wood. *Making Maps: A Visual Guide to Map Design for GIS*. English. 1 edition. New York: The Guilford Press, Aug. 2005. ISBN: 978-1-59385-200-9.
- [38] Michael Friendly. “A Brief History of Data Visualization”. en. In: *Handbook of Data Visualization*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 15–56. ISBN: 978-3-540-33036-3 978-3-540-33037-0. DOI: [10.1007/978-3-540-33037-0\\_2](https://doi.org/10.1007/978-3-540-33037-0_2).
- [39] L. Byrne, D. Angus, and J. Wiles. “Acquired Codes of Meaning in Data Visualization and Infographics: Beyond Perceptual Primitives”. In: *IEEE Transactions on Visualization and Computer Graphics* 22.1 (Jan. 2016), pp. 509–518. ISSN: 1077-2626. DOI: [10.1109/TVCG.2015.2467321](https://doi.org/10.1109/TVCG.2015.2467321).
- [40] Jacques Bertin. *Semiology of Graphics : Diagrams, Networks, Maps*. English. Redlands, Calif.: ESRI Press, 2011. ISBN: 978-1-58948-261-6 1-58948-261-1.
- [41] C. Ware. *Information Visualization: Perception for Design*. Interactive Technologies. Elsevier Science, 2019. ISBN: 978-0-12-812876-3.
- [42] Tamara Munzner. *Visualization Analysis and Design*. AK Peters Visualization Series. CRC press, Oct. 2014. ISBN: 978-1-4665-0891-0.
- [43] Caroline Ziemkiewicz and Robert Kosara. “Embedding Information Visualization within Visual Representation”. In: *Advances in Information and Intelligent Systems*. Ed. by Zbigniew W. Ras and William Ribarsky. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 307–326. ISBN: 978-3-642-04141-9. DOI: [10.1007/978-3-642-04141-9\\_15](https://doi.org/10.1007/978-3-642-04141-9_15).
- [44] JOCK D. MACKINLAY. “AUTOMATIC DESIGN OF GRAPHICAL PRESENTATIONS (DATABASE, USER INTERFACE, ARTIFICIAL INTELLIGENCE, INFORMATION TECHNOLOGY)”. English. PhD Thesis. 1987.
- [45] William S. Cleveland. “Research in Statistical Graphics”. In: *Journal of the American Statistical Association* 82.398 (June 1987), p. 419. ISSN: 01621459. DOI: [10.2307/2289443](https://doi.org/10.2307/2289443).

- [46] William S. Cleveland and Robert McGill. “Graphical Perception: Theory, Experimentation, and Application to the Development of Graphical Methods”. In: *Journal of the American Statistical Association* 79.387 (Sept. 1984), pp. 531–554. ISSN: 0162-1459. DOI: [10.1080/01621459.1984.10478080](https://doi.org/10.1080/01621459.1984.10478080).
- [47] John M Chambers et al. *Graphical Methods for Data Analysis*. Vol. 5. Wadsworth Belmont, CA, 1983.
- [48] *Naturalness Principle - InfoVis:Wiki*. [https://infovis-wiki.net/wiki/Naturalness\\_Principle](https://infovis-wiki.net/wiki/Naturalness_Principle).
- [49] Edward R. Tufte. *The Visual Display of Quantitative Information*. English. Cheshire, Conn.: Graphics Press, 2001. ISBN: 0-9613921-4-2 978-0-9613921-4-7 978-1-930824-13-3 1-930824-13-0.
- [50] T. W. E. B. Du Bois Center at the University of Massachusetts, W. Battle-Baptiste, and B. Rusert. *W. E. B. Du Bois’s Data Portraits: Visualizing Black America*. Princeton Architectural Press, 2018. ISBN: 978-1-61689-706-2.
- [51] *[The Georgia Negro] City and Rural Population. 1890*. eng. <https://www.loc.gov/item/2013650430/>. Image. 1900.
- [52] W. E. B. (William Edward Burghardt) Du Bois. *[The Georgia Negro] Valuation of Town and City Property Owned by Georgia Negroes*. en. <https://www.loc.gov/pictures/item/2013650441/>. Still Image. 1900.
- [53] *[A Series of Statistical Charts Illustrating the Condition of the Descendants of Former African Slaves Now in Residence in the United States of America] Negro Population of the United States Compared with the Total Population of Other Countries /*. eng. <https://www.loc.gov/item/2013650368/>. Image.
- [54] *[A Series of Statistical Charts Illustrating the Condition of the Descendants of Former African Slaves Now in Residence in the United States of America] Negro Business Men in the United States*. eng. <https://www.loc.gov/item/2014645363/>. Image.
- [55] Jeffrey Heer, Michael Bostock, and Vadim Ogievetsky. “A Tour Through the Visualization Zoo”. In: *Communications of the ACM* 53.6 (June 2010), pp. 59–67. ISSN: 0001-0782. DOI: [10.1145/1743546.1743567](https://doi.org/10.1145/1743546.1743567).

- [56] A. M. Cegarra. “Cohomology of Monoids with Operators”. In: *Semigroup Forum*. Vol. 99. Springer, 2019, pp. 67–105. ISBN: 1432-2137.
- [57] E.H. Spanier. *Algebraic Topology*. McGraw-Hill Series in Higher Mathematics. Springer, 1989. ISBN: 978-0-387-94426-5.
- [58] *Locally Trivial Fibre Bundle - Encyclopedia of Mathematics*. [https://encyclopediaofmath.org/wiki/Locally\\_trivial\\_fibre\\_bundle](https://encyclopediaofmath.org/wiki/Locally_trivial_fibre_bundle)
- [59] S. S. Stevens. “On the Theory of Scales of Measurement”. In: *Science* 103.2684 (1946), pp. 677–680. ISSN: 00368075, 10959203.
- [60] Brent A Yorgey. “Monoids: Theme and Variations (Functional Pearl)”. en. In: (), p. 12.
- [61] Michele Stieven. *A Monad Is Just a Monoid...* en. <https://medium.com/@michelestieven/a-monad-is-just-a-monoid-a02bd2524f66>. Apr. 2020.
- [62] “Monoid”. en. In: *Wikipedia* (Jan. 2021).
- [63] “Semigroup Action”. en. In: *Wikipedia* (Jan. 2021).
- [64] *Action in nLab*. [https://ncatlab.org/nlab/show/action#actions\\_of\\_a\\_monoid](https://ncatlab.org/nlab/show/action#actions_of_a_monoid).
- [65] W A Lea. “A Formalization of Measurement Scale Forms”. en. In: (), p. 44.
- [66] Eric W. Weisstein. *Similarity Transformation*. en. <https://mathworld.wolfram.com/SimilarityTransformation.html>. Text.
- [67] “Quotient Space (Topology)”. en. In: *Wikipedia* (Nov. 2020).
- [68] Professor Denis Auroux. “Math 131: Introduction to Topology”. en. In: (), p. 113.
- [69] Charles R Harris et al. “Array Programming with NumPy”. In: *Nature* 585.7825 (2020), pp. 357–362.
- [70] Jeff Reback et al. *Pandas-Dev/Pandas: Pandas 1.0.3*. Zenodo. Mar. 2020. DOI: [10.5281/zenodo.3715232](https://doi.org/10.5281/zenodo.3715232).
- [71] Stephan Hoyer and Joe Hamman. “Xarray: ND Labeled Arrays and Datasets in Python”. In: *Journal of Open Research Software* 5.1 (2017).

- [72] Matthew Rocklin. “Dask: Parallel Computation with Blocked Algorithms and Task Scheduling”. In: *Proceedings of the 14th Python in Science Conference*. Vol. 126. Citeseer, 2015.
- [73] “Retraction (Topology)”. en. In: *Wikipedia* (July 2020).
- [74] Tim Bienz, Richard Cohn, and Calif.) Adobe Systems (Mountain View. *Portable Document Format Reference Manual*. Citeseer, 1993.
- [75] A. Quint. “Scalable Vector Graphics”. In: *IEEE MultiMedia* 10.3 (July 2003), pp. 99–102. ISSN: 1941-0166. DOI: [10.1109/MMUL.2003.1218261](https://doi.org/10.1109/MMUL.2003.1218261).
- [76] George S. Carson. “Standards Pipeline: The OpenGL Specification”. In: *SIGGRAPH Comput. Graph.* 31.2 (May 1997), pp. 17–18. ISSN: 0097-8930. DOI: [10.1145/271283.271292](https://doi.org/10.1145/271283.271292).
- [77] *Cairographics.Org*. <https://www.cairographics.org/>.
- [78] *Anti-Grain Geometry* -. <http://agg.sourceforge.net/antigrain.com/index.html>.
- [79] John Hunter and Michael Droettboom. *The Architecture of Open Source Applications (Volume 2): Matplotlib*. <https://www.aosabook.org/en/matplotlib.html>.
- [80] “Jet Bundle”. en. In: *Wikipedia* (Dec. 2020).
- [81] Jana Musilová and Stanislav Hronek. “The Calculus of Variations on Jet Bundles as a Universal Approach for a Variational Formulation of Fundamental Physical Theories”. In: *Communications in Mathematics* 24.2 (Dec. 2016), pp. 173–193. ISSN: 2336-1298. DOI: [10.1515/cm-2016-0012](https://doi.org/10.1515/cm-2016-0012).
- [82] “Connected Space”. en. In: *Wikipedia* (Dec. 2020).
- [83] Sheelagh Carpendale. *Visual Representation from Semiology of Graphics by J. Bertin*. en.
- [84] Robert W. Ghrist. *Elementary Applied Topology*. Vol. 1. Createspace Seattle, 2014.
- [85] Robert Ghrist. “Homological Algebra and Data”. In: *Math. Data* 25 (2018), p. 273.

- [86] Michael S. Crouch, Andrew McGregor, and Daniel Stubbs. “Dynamic Graphs in the Sliding-Window Model”. In: *European Symposium on Algorithms*. Springer, 2013, pp. 337–348.
- [87] Chia-Shang James Chu. “Time Series Segmentation: A Sliding Window Approach”. In: *Information Sciences* 85.1 (July 1995), pp. 147–173. ISSN: 0020-0255. DOI: [10.1016/0020-0255\(95\)00021-G](https://doi.org/10.1016/0020-0255(95)00021-G).
- [88] Dmitry Nekrasovski et al. “An Evaluation of Pan & Zoom and Rubber Sheet Navigation with and without an Overview”. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. CHI '06. New York, NY, USA: Association for Computing Machinery, 2006, pp. 11–20. ISBN: 1-59593-372-7. DOI: [10.1145/1124772.1124775](https://doi.org/10.1145/1124772.1124775).
- [89] Yael Albo et al. “Off the Radar: Comparative Evaluation of Radial Visualization Solutions for Composite Indicators”. In: *IEEE Transactions on Visualization and Computer Graphics* 22.1 (Jan. 2016), pp. 569–578. ISSN: 1077-2626. DOI: [10.1109/TVCG.2015.2467322](https://doi.org/10.1109/TVCG.2015.2467322).
- [90] Z. Qu and J. Hullman. “Keeping Multiple Views Consistent: Constraints, Validations, and Exceptions in Visualization Authoring”. In: *IEEE Transactions on Visualization and Computer Graphics* 24.1 (Jan. 2018), pp. 468–477. ISSN: 1941-0506. DOI: [10.1109/TVCG.2017.2744198](https://doi.org/10.1109/TVCG.2017.2744198).
- [91] Richard A. Becker and William S. Cleveland. “Brushing Scatterplots”. In: *Technometrics* 29.2 (May 1987), pp. 127–142. ISSN: 0040-1706. DOI: [10.1080/00401706.1987.10488204](https://doi.org/10.1080/00401706.1987.10488204).
- [92] Andreas Buja et al. “Interactive Data Visualization Using Focusing and Linking”. In: *Proceedings of the 2nd Conference on Visualization '91*. VIS '91. Washington, DC, USA: IEEE Computer Society Press, 1991, pp. 156–163. ISBN: 0-8186-2245-8.
- [93] *Dataclasses — Data Classes — Python 3.9.2rc1 Documentation*. <https://docs.python.org/3/library/dataclasses>.

- [94] Kristen B. Gorman, Tony D. Williams, and William R. Fraser. “Ecological Sexual Dimorphism and Environmental Variability within a Community of Antarctic Penguins (Genus *Pygoscelis*)”. In: *PLOS ONE* 9.3 (Mar. 2014), e90081. DOI: [10.1371/journal.pone.0090081](https://doi.org/10.1371/journal.pone.0090081).
- [95] Allison Marie Horst, Alison Presmanes Hill, and Kristen B Gorman. *Palmerpenguins: Palmer Archipelago (Antarctica) Penguin Data*. Manual. 2020. DOI: [10.5281/zenodo.3960218](https://doi.org/10.5281/zenodo.3960218).
- [96] Muhammad Chenariyan Nakhaee. *Mcnakhaee/Palmerpenguins*. Jan. 2021.
- [97] F. Beck. “Software Feathers Figurative Visualization of Software Metrics”. In: *2014 International Conference on Information Visualization Theory and Applications (IVAPP)*. Jan. 2014, pp. 5–16.
- [98] Lydia Byrne, Daniel Angus, and Janet Wiles. “Figurative Frames: A Critical Vocabulary for Images in Information Visualization”. In: *Information Visualization* 18.1 (Aug. 2017), pp. 45–67. ISSN: 1473-8716. DOI: [10.1177/1473871617724212](https://doi.org/10.1177/1473871617724212).
- [99] Robert B Haber and David A McNabb. “Visualization Idioms: A Conceptual Model for Scientific Visualization Systems”. In: *Visualization in scientific computing* 74 (1990), p. 93.
- [100] Charles D Hansen and Chris R Johnson. *Visualization Handbook*. Elsevier, 2011.
- [101] Leland Wilkinson and Michael Friendly. “The History of the Cluster Heat Map”. In: *The American Statistician* 63.2 (May 2009), pp. 179–184. ISSN: 0003-1305. DOI: [10.1198/tas.2009.0033](https://doi.org/10.1198/tas.2009.0033).
- [102] Toussaint Loua. *Atlas Statistique de La Population de Paris*. J. Dejeu & cie, 1873.
- [103] Hadley Wickham and Lisa Stryjewski. “40 Years of Boxplots”. In: *The American Statistician* (2011).
- [104] C. Heine et al. “A Survey of Topology-Based Methods in Visualization”. In: *Computer Graphics Forum* 35.3 (June 2016), pp. 643–667. ISSN: 0167-7055. DOI: [10.1111/cgf.12933](https://doi.org/10.1111/cgf.12933).

- 1085 [105] James O Ramsay. *Functional Data Analysis*. Wiley Online Library, 2006.
- 1086 [106] Enrico Bertini, Andrada Tatu, and Daniel Keim. “Quality Metrics in High-  
 1087 Dimensional Data Visualization: An Overview and Systematization”. In: *IEEE*  
 1088 *Transactions on Visualization and Computer Graphics* 17.12 (2011), pp. 2203–2212.
- 1089 [107] Brendan Fong and David I. Spivak. *An Invitation to Applied Category Theory:  
 1090 Seven Sketches in Compositionality*. en. First. Cambridge University Press, July  
 1091 2019. ISBN: 978-1-108-66880-4 978-1-108-48229-5 978-1-108-71182-1. DOI: [10.1017/  
 1092 9781108668804](https://doi.org/10.1017/9781108668804).
- 1093 [108] Bartosz Milewski. “Category Theory for Programmers”. en. In: (), p. 498.