

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

This work presents a functional model of the structure-preserving maps from data to visual representation to guide the development of visualization libraries. Our model, which we call the topological equivariant artist model (TEAM), provides a means to express the constraints of preserving the data continuity in the graphic and faithfully translating the properties of the data variables into visual variables. We formalize these transformations as actions on sections of topological fiber bundles, which are mathematical structures that allow us to encode continuity as a base space, variable properties as a fiber space, and data as binding maps, called sections, between the base and fiber spaces. This abstraction allows us to generalize to any type of data structure, rather than assuming, for example, that the data is a relational table, image, data cube, or network-graph. Moreover, we extend the fiber bundle abstraction to the graphic objects that the data is mapped to. By doing so, we can track the preservation of data continuity in terms of continuous maps from the base space of the data bundle to the base space of the graphic bundle. Equivariant maps on the fiber spaces preserve the structure of the variables; this structure can be represented in terms of monoid actions, which are a generalization of the mathematical structure of Stevens' theory of measurement scales. We briefly sketch that these transformations have an algebraic structure which lets us build complex components for visualization from simple ones. We demonstrate the utility of this model through case studies of a scatter plot, line plot, and image. To demonstrate the feasibility of the model, we implement a prototype of a scatter and line plot in the context of the Matplotlib Python visualization library. We propose that the functional architecture derived from a TEAM based design specification can provide a basis for a more consistent API and better modularity, extendability, scaling and support for concurrency.

# Contents

36	<b>Abstract</b>	<b>ii</b>
37	<b>1 Introduction</b>	<b>1</b>
38	<b>2 Background</b>	<b>2</b>
39	2.1 Structure: . . . . .	3
40	2.2 Tools . . . . .	5
41	2.3 Data . . . . .	8
42	2.4 Contribution . . . . .	10
43	<b>3 Topological Artist Model</b>	<b>10</b>
44	3.1 Data Space $E$ . . . . .	11
45	3.1.1 Variables in Fiber Space $F$ . . . . .	11
46	3.1.2 Measurement Scales: Monoid Actions . . . . .	13
47	3.1.3 Continuity of the Data $K$ . . . . .	15
48	3.1.4 Data $\tau$ . . . . .	18
49	3.1.5 Sheafs . . . . .	19
50	3.1.6 Applications to Data Containers . . . . .	20
51	3.2 Graphic Space $H$ . . . . .	21
52	3.2.1 Idealized Display $D$ . . . . .	21
53	3.2.2 Continuity of the Graphic $S$ . . . . .	22
54	3.2.3 Graphic $\rho$ . . . . .	23
55	3.3 Artist . . . . .	26
56	3.3.1 Visual Fiber Bundle $V$ . . . . .	27
57	3.3.2 Visual Encoders $\nu$ . . . . .	28
58	3.3.3 Visualization Assembly . . . . .	31
59	3.3.4 Assembly $Q$ . . . . .	33
60	3.3.5 Composition of Artists: $+$ . . . . .	37
61	3.3.6 Equivalence class of artists $A'$ . . . . .	39

62	<b>4 Prototype Implementation: Matplottoy</b>	<b>40</b>
63	4.1 Artist Class $A'$ . . . . .	42
64	4.2 Encoders $\nu$ . . . . .	48
65	4.3 Data $E$ . . . . .	50
66	4.4 Case Study: Penguins . . . . .	55
67	<b>5 Discussion</b>	<b>59</b>
68	5.1 Limitations . . . . .	60
69	5.2 Future Work . . . . .	61
70	<b>6 Conclusion</b>	<b>62</b>

# 1 Introduction

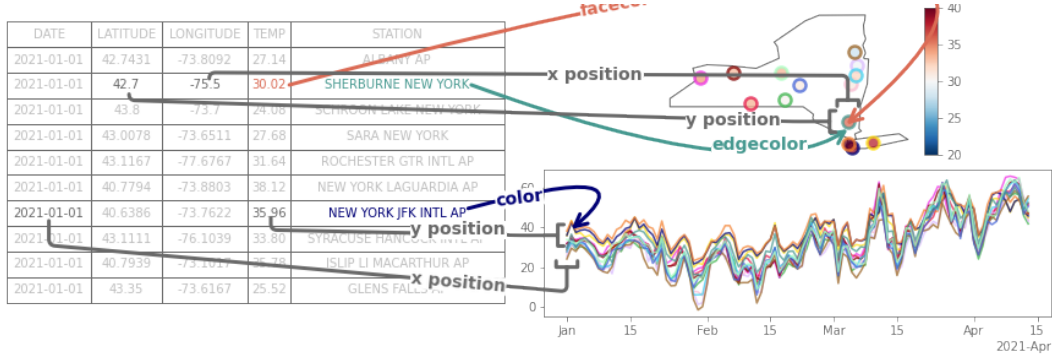


Figure 1: Visualizations are made up of transformations from data into visual representation. These functions transform individual data values to visual representation, such as date to x position or latitude to y position. These functions are composed into the assembly of all these transformations into a visual mark, such as a line or point. The same variable can be mapped in different ways, for example line is mapped to a color in the scatter plot and to y position in the line plot.

Visualization is the transformation of data into visual representation. As illustrated by Figure 1, these translations are both at the level of the individual variable and the entire record. In the case of the scatter plot, the latitude and longitude are encoded as the x and y position, respectively, while the temperature and station are represented by the face and edge colors. A row in the table is collectively encoded as a point mark. None of these encodings are fixed, as evidenced by temperature being translated into the y value in the case of the line plot. The station is now the source of the color of the entire line, and the date is the x position. As with scatter, the encodings of the individual transformations, which again are on values from the same record in the table, are composited into a line mark. It is these raw transformations from data space to visualization space that are implemented by building block level visualization libraries, named as such because the functions provided by the library can be composited in any number of ways to yield visualizations [1]. We propose that like physical building blocks, building block libraries must provide a collection of well defined pieces that can be composed in whichever ways the blocks fit together. We specify that a valid visualization block is a structure preserving transformation from data to visual space, and we define structure in terms of *continuity* and *equivariance*. We

88 then use this model to develop a design specification for the components of a building block  
89 visualization library. The notion of self contained, inherently modular, building blocks lends  
90 itself naturally to a functional paradigm of visualization [2]. We adopt a functional model  
91 for a redesign because the lack of side effects means functional architecture can be evaluated  
92 for correctness, functional programs tend to be shorter and clearer, and are well suited to  
93 distributed, concurrent, and on demand tasks[3].

94 This work is strongly motivated by the needs of the Matplotlib[4, 5] visualization library.  
95 One of the most widely used visualization libraries in Python, since 2002 new components  
96 and features have been added in a some what adhoc, sometimes hard to maintain, manner.  
97 Particularly, each new component carries its own implicit notion of how it believes the data is  
98 structured-for example if the data is a table, cube, image, or network - that is then expressed  
99 in the API for that component. In turn, this yields an inconsistent API for interfacing with  
100 the data, for example when updating streaming visualizations or constructing dashboards[6].  
101 This entangling of data model with visual transform also yields inconsistencies in how visual  
102 component transforms, e.g. shape or color, are supported. We propose that these issues can  
103 be ameliorated via a redesign of the functions that convert data to graphics, named *Artists* in  
104 Matplotlib, in a manner that reliably enforces *continuity* and *equivariance* constraints. We  
105 evaluate our functional model by implementing new artists in Matplotlib that are specified  
106 via *equivariance* and *continuity* constraints. We then use the common data model introduced  
107 by the model to demonstrate how plotting functions can be consolidated in a way that makes  
108 clear whether the difference is in expected data structure, visual component encoding, or  
109 the resulting graphic.

## 110 2 Background

111 There are many formalisms of the notion that visualization is structure preserving maps  
112 from data to visual representation, and many visualization libraries that attempt to pre-  
113 serve structure in some manner; this work bridges the formalism and implementation in a  
114 functional manner with a topological approach at a building blocks library level to propose

a new model of the constraints visual transformations must satisfy such that they can be composed to produce visualize representations that can be considered equivalent to the data being represented.

## 2.1 Structure:

Visual representations of data, by definition, reflect something of the underlying structure and semantics[7], whether through direct mappings from data into visual elements or via figurative representations that have meaning due to their similarity in shape to external concepts [8]. The components of a visual representation were first codified by Bertin[9], who introduced a notion of structure preservation that we formally describe in terms of *equivariance* and *continuity*.

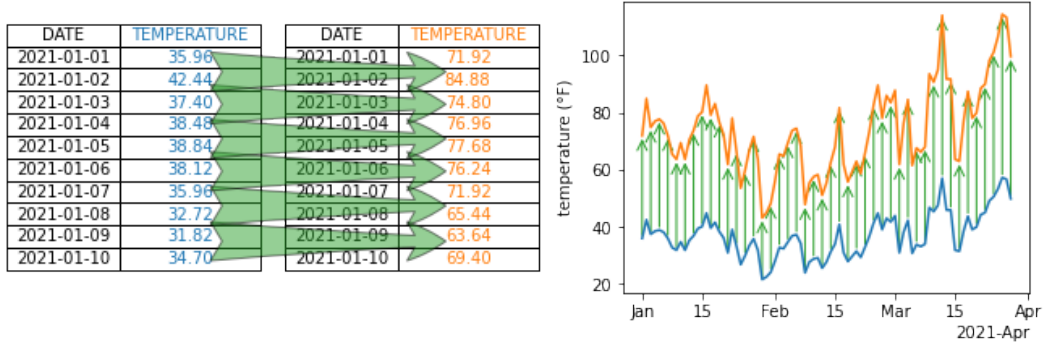


Figure 2: The data in blue is scaled by a factor of two, yielding the data in orange. To preserve *equivariance*, the blue line plot representation of the unscaled data is also scaled by a factor of two, yielding the orange line plot that is equivalent to the scaled data.

Bertin proposes that there are classes of visual encodings-such as position, shape, color, and texture-that when mapped to from specific types of measurement, quantitative or qualitative, will preserve the properties of that measurement type. For example, in Figure 2, the data and visual representation are scaled by equivalent factors of two, resulting in the change illustrated in the shift from blue to orange data and lines. The idea of equivariance is formally defined as the mapping of a binary operator from the data domain to the visual domain in Mackinlay’s *A Presentation Tool*(APT) model[10, 11]. The algebraic model of

132 visualization proposed by Kindlmann and Scheidegger uses equivariance to refer generally  
 133 to invertible binary transformations[12], which are mathematical groups [13]. Our model  
 134 defines *equivariance* in terms of monoid actions, which are a more restrictive set than all bi-  
 135 nary operations and more general than groups. As with the algebraic model, our model also  
 136 defines structure preservation as commutative mappings from data space to representation  
 137 space to graphic space, but our model uses topology to explicitly include continuity.

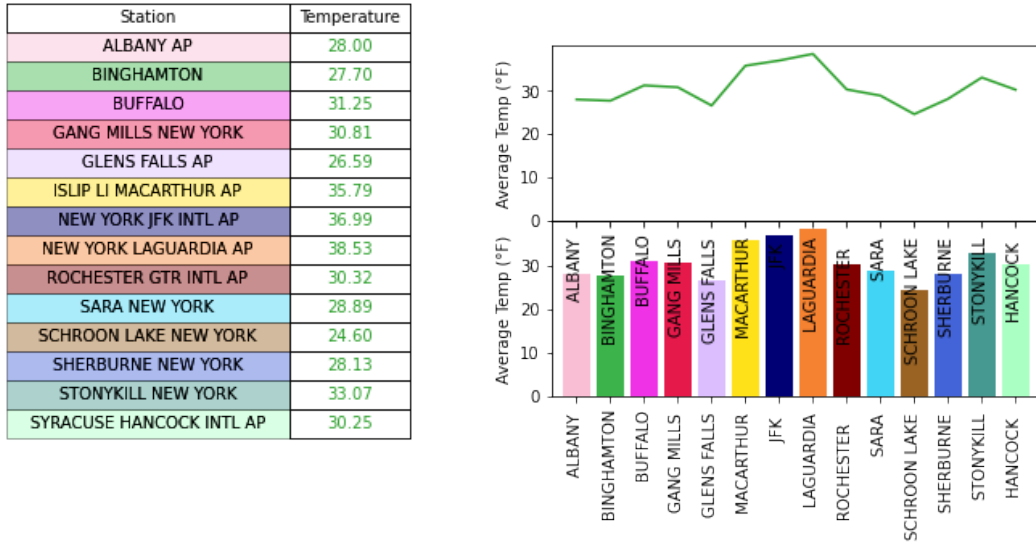


Figure 3: The line plot does not preserve *continuity* because it implies that the average temperature at each station lie along a 1D continuous line, while the bar plot preserves *continuity* by representing the average temperatures at each station as the discrete values they are.

138 Bertin proposes that the visual encodings be composited into graphical marks that match  
 139 the *continuity* of the data - for example discrete data is a point, 1D continuous is the line,  
 140 and 2D data is the area mark. In Figure 3, the line plot does not preserve continuity because  
 141 the line connecting the discrete categories implies that the frequency of weather events is  
 142 sampled from a continuous interval and the categories are points on that interval. But,  
 143 when the continuity is preserved, as in the bar chart, then the graphic has not introduced  
 144 new structure into the data.



## Structure

**continuity** How records in the dataset are connected to each other, e.g. discrete rows, networked nodes, points on a continuous surface

**equivariance** if an action is applied to the data or the graphic—e.g. a rotation, permutation, translation, or rescaling—there must be an equivalent action applied on the other side of the transformation.

145 The notion that a graphic should be equivalent to the data has been expressed in a  
 146 variety of ways. Informally, Norman’s Naturalness Principal[14] states that a visualization  
 147 is easier to understand when the properties of the visualization match the properties of  
 148 the data. This principal is made more concrete in Tufte’s concept of graphical integrity,  
 149 which is that a visual representation of quantitative data must be directly proportional to  
 150 the numerical quantities it represents (Lie Principal), must have the same number of visual  
 151 dimensions as the data, and should be well labeled and contextualized, and not have any  
 152 extraneous visual elements[15]. expressing, as defined by Mackinlay, is a measure how much  
 153 of the mathematical structure in the data that can be expressed in the visualizations; for  
 154 example that ordered variables can be mapped into ordered visual elements. We propose  
 155 that a graphic is an equivalent representation of the data when *continuity* and *equivariance*  
 156 are preserved.

## 2.2 Tools

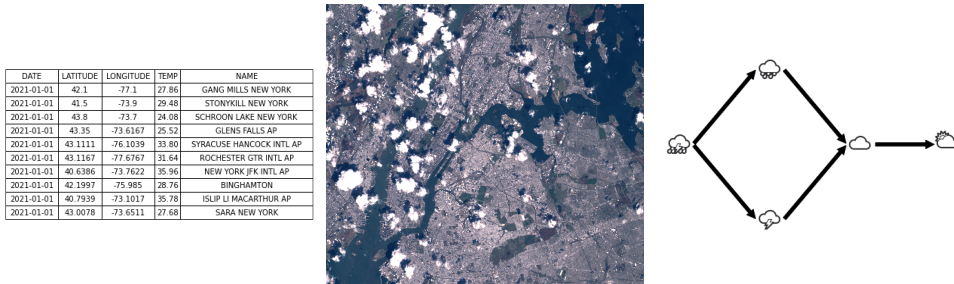


Figure 4: Visualization libraries, especially ones tied to specific domains, tend to be architected around a core data structure, such as tables, images, or networks.

158 One of the reasons we developed a new formalism rather than adopting the architecture  
159 of an existing library is that most information visualization software design patterns, as  
160 categorized by Heer and Agrawala[16], are tuned to very specific data structures. These  
161 libraries can often assume the expected data structure because they are domain specific,  
162 and that is the common data structure in that domain. For users who generally work in  
163 one domain, such as the data, networks, or graphs shown in Figure 4, this well defined data  
164 space (and corresponding visual space[17]) often yields a very coherent user experience[18].  
165 But, for developers who want to build new visualizations on top of these libraries, they must  
166 work around the existing assumptions, sometimes in ways that break the model the libraries  
167 are developed around.

168 For example, many domain specific libraries integrate computation into the visualization,  
169 for example libraries based that assume all data is a relational database. This assumption is  
170 core to tools influenced by APT, such as Tableau[19–21] and the Grammar of Graphics[22],  
171 such as ggplot[23], protovis[24], vega[25] and altair[26]. Since these libraries represent all  
172 data as a table, and computations on tables are fairly well defined[27], they can include  
173 computations on the table with a fair bit of confidence that the computation is accurate.  
174 Since most computations are specific to domains, general purpose block libraries can not  
175 make this assumption; instead a goal of this model is to identify which computations are  
176 specifically part of the visual encoding - for example mapping data to a color-and which  
177 are manipulations on the data. Disentangling the computation from the visual transforms  
178 allows us to determine whether the visualization library needs to handle them or if they can  
179 be more efficiently computed by the data container.

180 A different class of user facing tools are those that support images, such as ImageJ[28]  
181 or Napari[29]. These tools often have some support for visualizing non image components  
182 of a complex data set, but mostly in service to the image being visualized. These tools are  
183 ill suited for general purpose libraries that need to support data other than images because  
184 the architecture is oriented towards building plugins into the existing system [30] where  
185 the image is the core data structure. Even the digital humanities oriented ImageJ macro  
186 ImagePlot[31], which supports some non-image aggregate reporting charts, is still built

187 around image data as the primary input. The need to visualize and manipulate graphs has  
 188 spawned tools like Gephi[32], Graphviz[33], and Networkx[34]. As with tables and images,  
 189 extending network libraries to work with other types of data either require breaking their  
 190 internal model of how data is structured and what transformations of the data are allowable  
 191 or growing a model for other types of data structures alongside the network model.

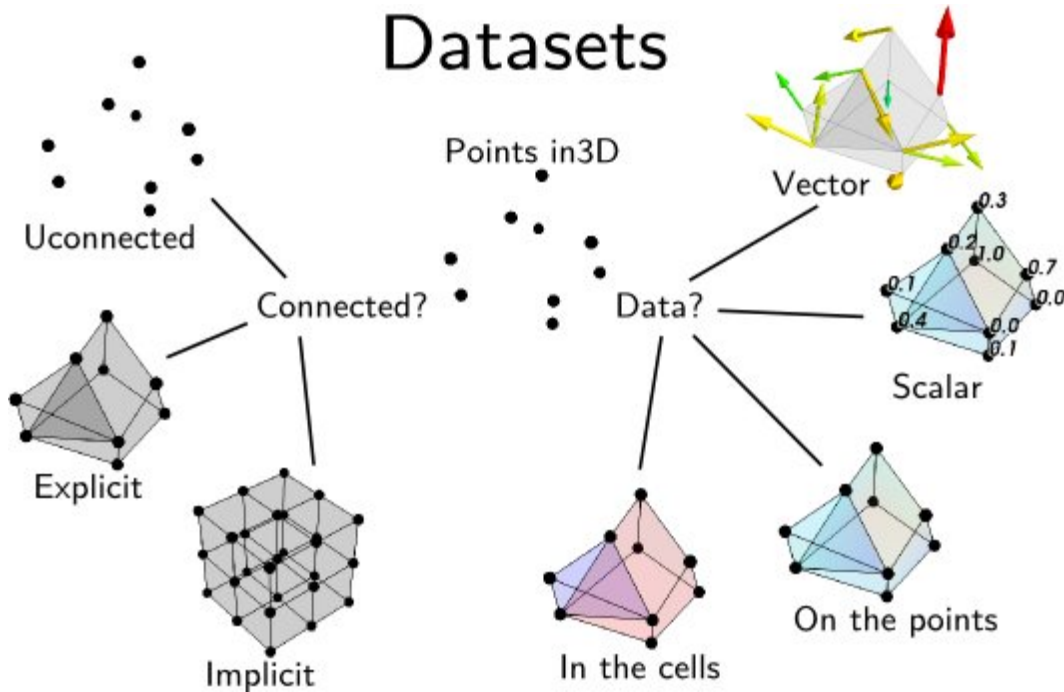


Figure 5: One way to describe data is by the connectivity of the points in the dataset. A database for example is often discrete unconnected points, while an image is an implicitly connected 2D grid. This image is from the Data Representation chapter of the MayaVi 4.7.2 documentation.[35]

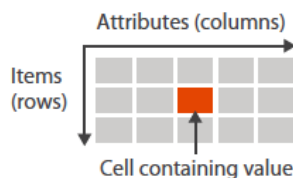
192 Many building block libraries carry multiple models of data internally because they can-  
 193 not assume a data structure. Algorithms are designed such that the structure of data is as-  
 194 sumed, as described in Tory and Möller’s taxonomy [ToryRethinkingVisualization2004],  
 195 and by definition building block libraries try to provide the components to build any sort of  
 196 visualization. Matplotlib, D3[36], and VTK [geveci2012vtk, 37] and its derivatives such as  
 197 MayaVi[38] and extensions such as ParaView[39] and the infoviz themed Titan[40]. Where  
 198 GoG and ImageJ type libraries have coherent APIs for their visualization tools because the

199 data structure is the same, the APIs for visualizations in Matplotlib, D3, and VTK are  
 200 significantly dependent on the structure of the data it expects. VTK has codified this in  
 201 terms of *continuity* based data representations, as illustrated in figure 5. This API choice  
 202 can lead to visualizations that break *continuity* when fed into visualizations with different  
 203 assumptions about structure. The lack of consistent data model can also mean no consistent  
 204 way of updating the data and therefore no way of guaranteeing that the views are in sync,  
 205 in visualizations that consist of multiple views of the same datasource, such as dash-  
 206 boards[6, 41]. To resolve this issue, our functional model takes as input a structure aware  
 207 data abstraction general enough to provide a common interface for many different types of  
 208 visualization.

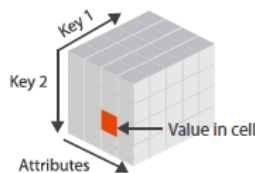
## 209 2.3 Data

210 One such general abstraction are fiber bundles, which Butler proposed as a core data struc-  
 211 ture for visualization because they encode data continuity separately from the variable  
 212 properties and are flexible enough to support discrete and ND continuous datasets [42, 43].  
 213 Since Butler’s model lacks a robust way of describing variables, we can encode a schema  
 214 like description of the data in the fiber bundle by folding in Spivak’s topological description  
 215 of data types [44, 45]. In this work we will refer to the points of the dataset as *records*  
 216 to indicate that a point can be a vector of heterogeneous elements. Each *component* of the  
 217 record is a single object, such as a temperature measurement, a color value, or an image.  
 218 We also generalize *component* to mean all objects in the dataset of a given type, such as  
 219 all temperatures or colors or images. The way in which these records are connected is the  
 220 *connectivity*, *continuity*, or more generally *topology*.

→ Tables



→ Multidimensional Table



→ Geometry (Spatial)



Figure 6: Values in a dataset have keys associated with them that describe where the value is in the dataset. These keys can be indexers or semantically meaningful; for example, in a table the keys are the variable name and the row ID. In the data cube, the keys are the row, column, and cell ID, and in the map the key is the position in the grid. Image is figure 2.8 in Munzner’s Visualization Analysis and Design[46]

The *continuity* can often be described by some variables in the dataset; this is formalized Munzner’s notion of metadata as *keys* into the data structure that return associated *values*[47]. As shown in Figure 6, keys can be labeled indexes, such as the attribute name and row ID, or physical entities such as locations on a map. We propose that information rich metadata are part of the components and instead the values are keyed on coordinate free structural ids. In contrast to Munzner’s model where the semantic meaning of the key is tightly coupled to the position of the value in the dataset, our model considers keys to be a pure reference to topology. This allows the metadata to be altered, without imposing new semantics on the underlying structure, for example by changing the coordinate systems or time resolution. This value agnostic model also supports encoding datasets where there may be multiple independent variables - such as a measure of plant growth given variations in water, sunlight, and time - without having to assume any one variable is inducing the change in growth. For building block library developers, this means the components are able to fully traverse the data structures without having to know anything about the values or the semantic meaning of the structure. Since these components are by design *equivariant* and *continuity* preserving, domain specific library developers in different domains that both rely on the same continuity, for example 2D continuity, can then safely reuse the components to build tools that can safely make domain specific assumptions.

## 2.4 Contribution

The contribution of this work is

1. formalization of the topology preserving relationship between data and graphic via continuous maps [subsubsection 3.2.2](#)
2. formalization of property preservation from data component to visual representation as monoid action equivariant maps [subsubsection 3.3.2](#)
3. functional oriented visualization architecture built on the mathematical model to demonstrate the utility of the model [subsubsection 3.3.3](#)
4. prototype of the architecture built on Matplotlib’s infrastructure to demonstrate the feasibility of the model. ??

## 3 Topological Artist Model

To guide the implementation of structure preserving building block components, we develop a mathematical formalism of visualization that specifies how these components preserve *continuity* and *equivariance*. Inspired by the somewhat analogous component in Matplotlib[5], we call the transformation from data space to graphic that these building block components implement the *artist*.

$$\mathcal{A} : \mathcal{E} \rightarrow \mathcal{H} \tag{1}$$

The *artist*  $\mathcal{A}$  is a map from the data  $\mathcal{E}$  to graphic  $\mathcal{H}$  fiber bundles. To explain how the *artist* is a structure preserving map from data to graphic, we first describe how we model data ([subsection 3.1](#)) and graphics ([subsection 3.2](#)) as topological structures that encapsulate component types and continuity. We then discuss the maps from graphic to data ([subsubsection 3.2.2](#), data components to visual components ([subsubsection 3.3.2](#)), and visual components into graphic ([subsubsection 3.3.3](#)) that make up the artist.

### 256 3.1 Data Space $E$

Building on Butler’s proposal of using fiber bundles as a common data representation structure for visualization data[42, 43], a fiber bundle is a tuple  $(E, K, \pi, F)$  defined by the projection map  $\pi$

$$F \hookrightarrow E \xrightarrow{\pi} K \quad (2)$$

257 that binds the components of the data in  $F$  to the continuity of the data encoded in  $K$ .  
 258 The fiber bundle models the properties of data component types  $F$  (subsection 3.1.1),  
 259 the continuity of records  $K$  (subsection 3.1.3), the collections of records  $\tau$  (??), and the  
 260 space  $E$  of all possible datasets with these components and continuity. By definition fiber  
 261 bundles are locally trivial[48, 49], meaning that over a localized neighborhood  $U$  the total  
 262 space is the cartesian product  $K \times F$ . We use fiber bundles as the data model because they  
 263 are inclusive enough to express all the types of structures of data described in subsection 2.2

#### 264 3.1.1 Variables in Fiber Space $F$

To formalize the structure of the data components, we use notation introduced by Spivak [45] 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  $\mathbb{U}_\sigma$

$$\begin{array}{ccc} \mathbb{U}_\sigma & \longrightarrow & \mathbb{U} \\ \pi_\sigma \downarrow & & \downarrow \pi \\ C & \xrightarrow{\sigma} & \mathbf{DT} \end{array} \quad (3)$$

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

$$F = \mathbb{U}_{\sigma(c)} = \mathbb{U}_T \quad (4)$$

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

$$F = \mathbb{U}_{\sigma(c_1)} \times \dots \times \mathbb{U}_{\sigma(c_i)} \times \dots \times \mathbb{U}_{\sigma(c_n)} \quad (5)$$

which is equivalent to

$$F = F_0 \times \dots \times F_i \times \dots \times F_n \quad (6)$$

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

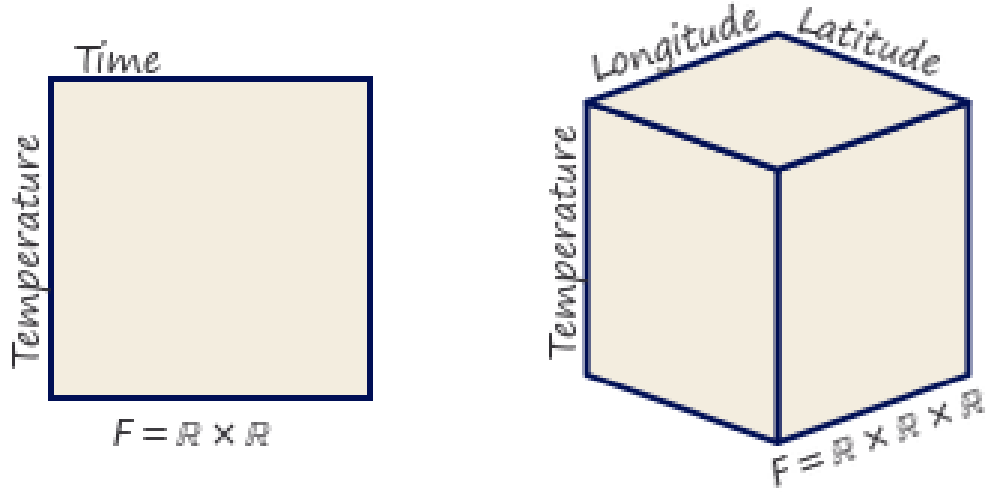


Figure 7: The fiber space is the set of all combinations of all theoretically possible values of the components. The 2D fiber  $F = \mathbb{R} \times \mathbb{R}$  encodes the properties of *time* and *temperature* components. One dimension of the fiber encodes the range of possible values for the time component of the dataset, which is a subset of the  $\mathbb{R}$ , while the other dimension encodes the range of possible values  $\mathbb{R}$  for the temperature component. This means the fiber is the set of points  $(temperature, time)$  that are all the combinations of *temperature*  $\times$  *time*. The 3D fiber encodes points at all possible combinations of *temperature*, *latitude*, and *longitude*.

For example, the records in the 2D fiber in ?? are 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

$$F = \mathbb{R} \times \mathbb{R}^+$$



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 the 3D fiber in ??, time is replaced with location. This location variable is of type **point** and has two components *latitude* and *longitude*  $\{(lat, lon) \in \mathbb{R}^2 \mid -90 \leq lat \leq 90, 0 \leq lon \leq 360\}$ . The fiber for this dataset is

$$F = \mathbb{R} \times \mathbb{R}^2 = \mathbb{R} \times \mathbb{R} \times \mathbb{R}$$

266 where the dimensionality of the fiber does not change, but the components of the fiber  
 267 can be coupled. For example, *location* can either be specified as  $(c = \text{location}, T =$   
 268 **point**,  $\mathbb{U}_\sigma = \mathbb{R}^2)$  or  $(c = \text{latitude}, T = \text{float}, \mathbb{U}_\sigma = \mathbb{R})$  and  $(c = \text{longitude}, T =$   
 269 **float**,  $\mathbb{U}_\sigma = \mathbb{R})$ .

270 As illustrated in [Figure 7](#), Spivak's framework provides a consistent way to describe  
 271 potentially complex components of the input data.

### 272 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. A monoid [50]  $M$  is a set with a binary operation  $*$  :  $M \times M \rightarrow M$  that satisfies the axioms:

**associativity** for all  $a, b, c \in M$   $(a \bullet b) \bullet c = a \bullet (b \bullet c)$

**identity** for all  $a \in M$ ,  $e \bullet a = a$

As defined on a component of  $F$ , a left monoid action [51, 52] 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  $x \in F_i, e \in M_i$ ,  $e \bullet x = x$

273 The identity and associativity properties of the action denote that the action is a monoid  
 274 homomorphism, which means that the group operation is preserved on both sides of the  
 275 action[53].

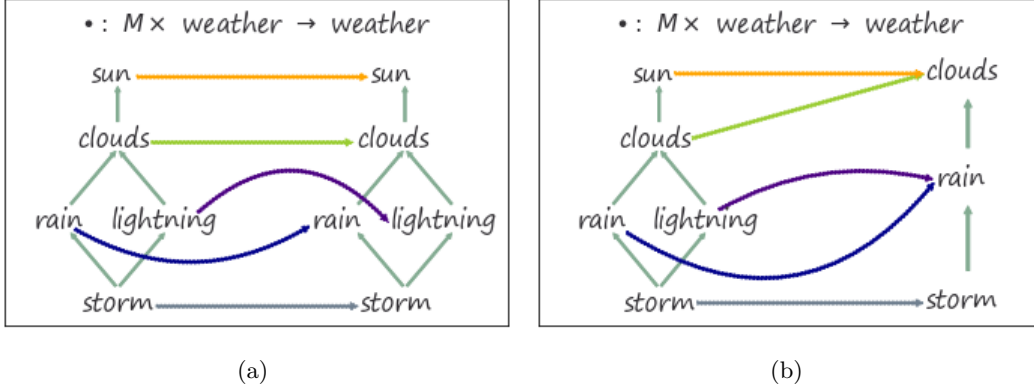


Figure 8: The action  $\bullet$  in ?? is the arrows from the partial order diagram of weather states on the left to the diagram of weather states on the right. Since the action maps the weather states to themselves, the ordering defined by the monoid  $*$  is preserved on both sides of the action. The action in ?? is monoid homomorphism because the ordering of the weather states is the same as the ordering of the elements they are mapped to. Given  $sun \geq clouds \geq rain$  on the right, the action  $sun \rightarrow clouds$ , and  $rain \rightarrow lightning$  is structure preserving because on the left  $cloud \geq rain$  so the relative ordering of elements is the same as the elements they are mapped to.

One example of monoids are partial orderings on a set, such as seen in . Each hasse diagram of the set of weather states describes an ordering on the set; the arrow goes from the lesser value to the greater one. For example,  $storm \leq rain$ . In ??, the action  $\bullet$  maps the elements of a set of weather states into itself by mapping them into other elements of the weather states. The action in Figure 8a, represented as the arrows between the hasse diagrams of the weather states, maps the weather states to themselves; therefore the

ordering of the weather states is identical on both sides of the action and it is therefore homomorphic. The action  $\bullet$  in Figure 8b is a monotone map[54]

$$if\ a \leq b\ then\ \bullet(a) \leq \bullet(b) \mid a, b \in F_i$$

where the structure the action preserves is the relative, rather than exact, ordering. Since groups are monoids with invertible operations, this definition of structure is also broad enough to include the Steven's measurment scales[55, 56]. Monoids are also commonly found in functional programming since the core property of monoids is composability [57].

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

$$M = M_0 \times \dots \times M_i \times \dots \times M_n \quad (7)$$

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

### 3.1.3 Continuity of the Data $K$



Figure 9: The topological base space  $K$  encodes the continuity of the data space, for example if the data is discrete points or lies on a plane or a sphere

The base space  $K$  acts as an indexing space, as emphasized by Butler[42, 43], to express how the records in  $E$  are connected to each other. As shown in Figure 9,  $K$  can have any number of dimensions and can be continuous or discrete. The base space also does not describe anything about the dataset besides the continuity. While the base space may have components to identify the continuity, such as *time*, *latitutde*, *longitude*, these labels

288 are indexed into from  $K$  the same as any other component. This is similar to the notion of  
 289 structural *keys* with associated *values* proposed by Munzner[46], but our model treats keys as  
 290 a pure reference to topology. Decoupling the keys from their semantics allows the metadata  
 291 to be altered; this provides for coordinate agnostic representation of the continuity and  
 292 facilitates encoding of data where the independent variable may not be clear. For example  
 293 the amount of snow on the ground is dependent on time of day and how much snow has  
 294 fallen, and changing the coordinate system or time resolution should have no effect on how  
 295 the records are connected to each other.

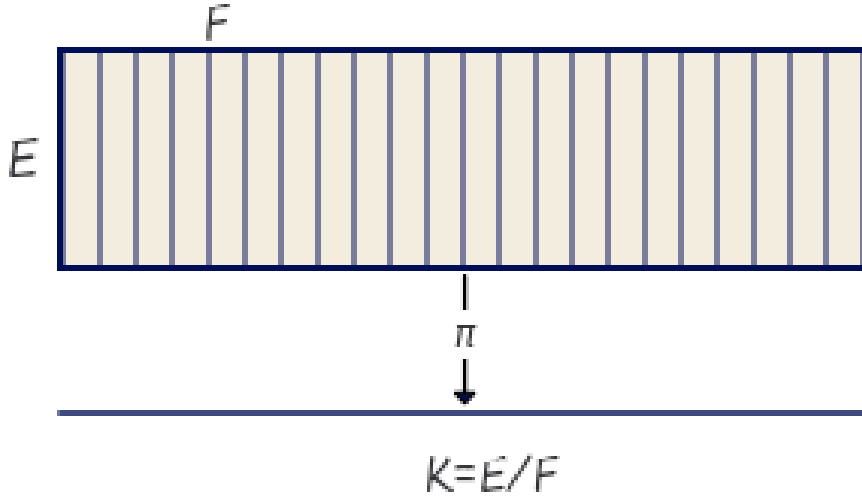


Figure 10: The base space  $E$  is divided into fiber segments  $F$ . The base space  $K$  acts as an index into the records in the fibers.

296 Formally  $K$  is the quotient space [58] of  $E$  meaning it is the finest space[59] such that  
 297 every  $k \in K$  has a corresponding fiber  $F_k$ [58]. In Figure 10,  $E$  is a rectangle divided by  
 298 vertical fibers  $F$ , so the minimal  $K$  for which there is always a mapping  $\pi : E \rightarrow K$  is the  
 299 closed interval  $[0, 1]$ .

As with Equation 6 and Equation 7, we can decompose the total space into component bundles  $\pi : E_i \rightarrow K$  where

$$\pi : E_1 \oplus \dots \oplus E_i \oplus \dots \oplus E_n \rightarrow K \quad (8)$$

such that  $M_i$  acts on component bundle  $E_i$ . The  $K$  remains the same because the connectivity of records does not change just because there are fewer components in each record. By encoding this continuity in the model as  $K$  the data model now explicitly carries information about its structure such that the implicit assumptions of the visualization algorithms are now explicit.

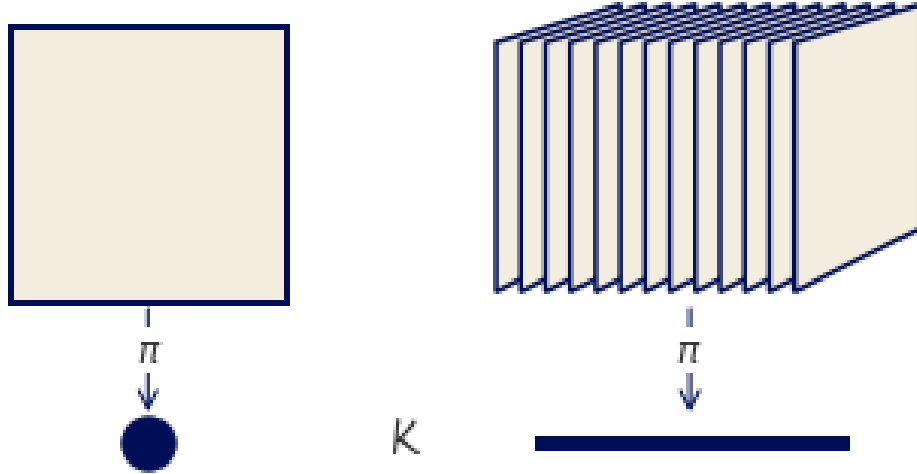


Figure 11: These two datasets have the same  $(time, temperature)$  fiber, but different continuities. The dataset on the left consists of discrete records, while the records in the dataset on the right sampled from a continuous space.

The datasets in Figure 11 have the same fiber of  $(temperature, time)$ . The dot represents a discrete base space  $K$ , meaning that every dataset encoded in the fiber bundle has discrete continuity. The line is a representation of a 1D continuity, meaning that every dataset in the fiber bundle is 1D continuous. By encoding this continuity in the model as  $K$  the data model now explicitly carries information about its structure such that the implicit assumptions of

the visualization algorithms are now explicit. The explicit topology is a concise way of distinguishing visualizations that appear identical, for example heatmaps and images.

#### 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 [45],  $K$  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

$$\begin{array}{ccc} F & \hookrightarrow & E \\ & \searrow \scriptstyle \tau & \downarrow \scriptstyle \pi \\ & & K \end{array} \quad (9)$$

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

$$\tau(k) = (k, (g_{F_0}(k), \dots, g_{F_n}(k))) \quad (10)$$

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 (Equation 8, Equation 6), we can decompose  $\tau$  as

$$\tau = (\tau_0, \dots, \tau_i, \dots, \tau_n) \quad (11)$$

where each section  $\tau_i$  maps into a record on a component  $F_i \in F$ . This allows for accessing the data component wise in addition to accessing the data in terms of its location over  $K$ .

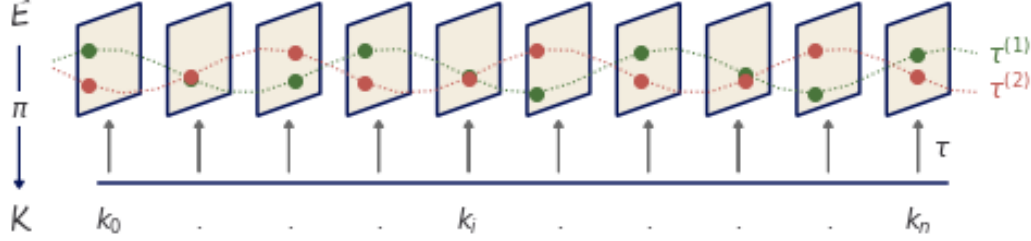


Figure 12: Fiber (time, temperature) with an interval  $K$  basespace. The sections  $\tau^{(1)}$  and  $\tau^{(2)}$  are constrained such that the time variable must be monotonic, which means each section is a timeseries of temperature values. They are included in the global set of sections  $\tau^{(1)}, \tau^{(2)} \in \Gamma(E)$

In Figure 12, the fiber is the same encoding of *(time, temperature)* illustrated in Figure 7, and the base space is the interval  $K$  shown in Figure 11. The section  $\tau^{(1)}$  is a function that for a point  $k$  returns a record in the fiber  $F$ . The section applied to a set of points in  $K$  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 Sheafs

Many types of dynamic visualizations require evaluating sections on different subspaces of  $K$ , and a sheaf, denoted  $\mathcal{O}$  provides a way to do so. A sheaf is a mathematical structure for defining collections of objects[60–62] on mathematical spaces. On the fiber bundle  $E$ , we can describe a sheaf as the collection of local sections  $\iota^* \tau$

$$\begin{array}{ccc} \iota^* E & \xleftarrow{\iota^*} & E \\ \pi \downarrow \Big) \iota^* \tau & & \pi \downarrow \Big) \tau \\ U & \xleftarrow{\iota} & K \end{array} \quad (12)$$

which are sections of  $E$  pulled back over local neighborhood  $U \subset E$  via the inclusion map  $\iota : E \rightarrow U$ . The collation of sections enabled by sheafs is necessary for navigation techniques such as pan and zoom[63] and dynamically updated visualizations such as sliding windows[64, 65].

### 3.1.6 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[66] 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[67] 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[68] 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 [69]. 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.



## 355 3.2 Graphic Space $H$

To establish that the artist is structure preserving map from data  $E$  to graphic  $H$  we construct a graphic bundle so that we can define *equivariance* in terms of maps on the fiber spaces and *continuity* in terms of maps on the base space. As with the data, we can represent the target graphic as a section  $\rho$  of a bundle  $(H, S, \pi, D)$ .

$$D \hookrightarrow H \quad \begin{array}{c} \pi \downarrow \\ S \end{array} \quad \begin{array}{c} \nearrow \rho \\ \searrow \end{array} \quad (13)$$

356 The graphic bundle  $H$  consists of a base  $S$ ([subsubsection 3.2.1](#)) that is a thickened form of  
 357  $K$  a fiber  $D$ ([subsubsection 3.2.2](#)) that is an idealized display space, and sections  $\rho$ ([??](#)) that  
 358 encode a graphic where the visual characteristics are fully specified.

### 359 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

$$(x, y, r, g, b) \in D$$

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

364 **3.2.2 Continuity of the Graphic  $S$**

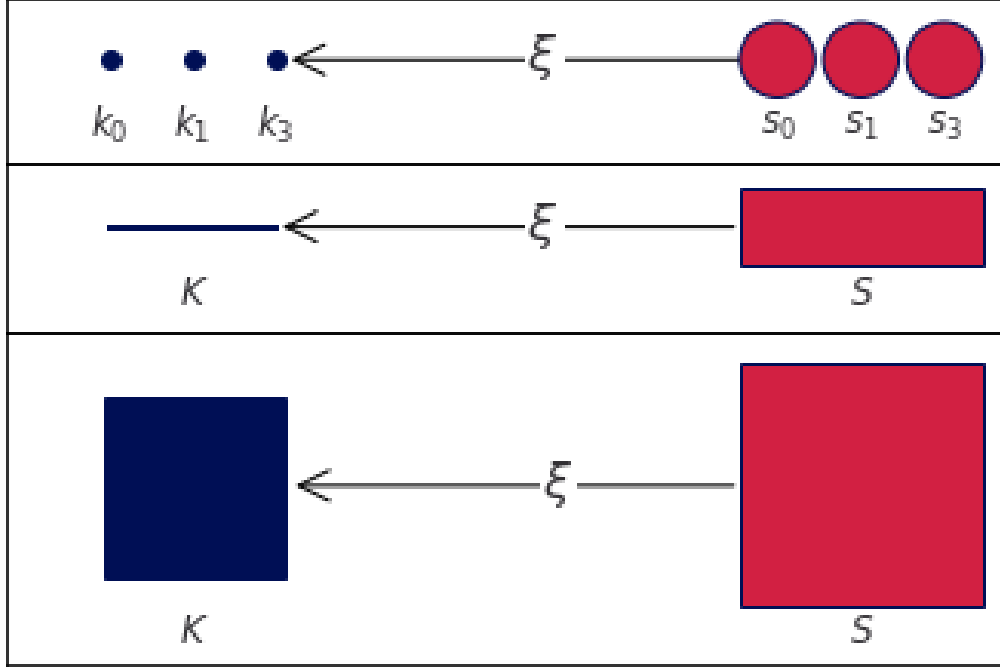


Figure 13: For a visualization component to preserve continuity, it must have a continuous surjective map  $\xi : S \rightarrow K$  from graphic continuity to data continuity. The scatter and line graphic base spaces  $S$  have one more dimension of continuity than  $K$  so that  $S$  can encode physical aspects of the glyph, such as shape (a circle) or thickness. The image has the same dimension in  $S$  as in  $K$ .

365 To establish that a visualization component preserves continuity, we propose that their must  
 366 be a continuous map  $\xi : S \rightarrow K$  from the graphic base space to the data space. For example,  
 367 consider a  $S$  that is mapped to the region of a 2D display space that represents  $K$ . For  
 368 some visualizations,  $K$  may be lower dimension than  $S$ . For example, a point that is 0D in  
 369  $K$  cannot be represented on screen unless it is thickened to 2D to encode the connectivity  
 370 of the pixels that visually represent the point. This thickening is often not necessary when  
 371 the dimensionality of  $K$  matches the dimensionality of the target space, for example if  $K$  is  
 372 2D and the display is a 2D screen. We introduce  $S$  to thicken  $K$  in a way which preserves  
 373 the structure of  $K$ .

Formally, we require that  $K$  be a deformation retract[70] of  $S$  so that  $K$  and  $S$  have the same homotopy, meaning there is a continuous map from  $S$  to  $K$ [71]. The surjective map  $\xi : S \rightarrow K$

$$\begin{array}{ccc} E & & H \\ \pi \downarrow & & \pi \downarrow \\ K & \xleftarrow{\xi} & S \end{array} \quad (14)$$

374 goes from region  $s \in S_k$  to its associated point  $s$ . This means that if  $\xi(s) = k$ , the record at  
 375  $k$  is copied over the region  $s$  such that  $\tau(k) = \xi^* \tau(s)$  where  $\xi^* \tau(s)$  is  $\tau$  pulled back over  $S$ .  
 376 When  $K$  is discrete points and the graphic is a scatter plot, each point  $k \in K$  corresponds  
 377 to a 2D disk  $S_k$  as shown in Figure 13. In the case of 1D continuous data and a line plot,  
 378 the region  $\beta$  over a point  $\alpha_i$  specifies the thickness of the line in  $S$  for the corresponding  
 379  $\tau$  on  $k$ . The image has the same dimensions in data space and graphic space such that no  
 380 extra dimensions are needed in  $S$ .

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

### 388 3.2.3 Graphic $\rho$

The section  $\rho : S \rightarrow H$  is the graphic in an idealizes prerender space and also acts as a specification for rendering the graphic to an image. It is sufficient to sketch out how an arbitrary pixel would be rendered, where a pixel  $p$  in a real display corresponds to a region  $S_p$  in the idealized display. To determine the color of the pixel, we aggregate the color values

over the region via integration:

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

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

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

389 For a 2D screen, the pixel is defined as a region  $p = [y_{top}, y_{bottom}, x_{right}, x_{left}]$  of the rendered  
 390 graphic. Since the x and y in  $p$  are in the same coordinate system as the x and y components  
 391 of  $D$  the inverse map of the bounding box  $S_p = \rho_{xy}^{-1}(p)$  is a region  $S_p \subset S$ . The color is  
 392 the result of the integration over  $S_p$ .

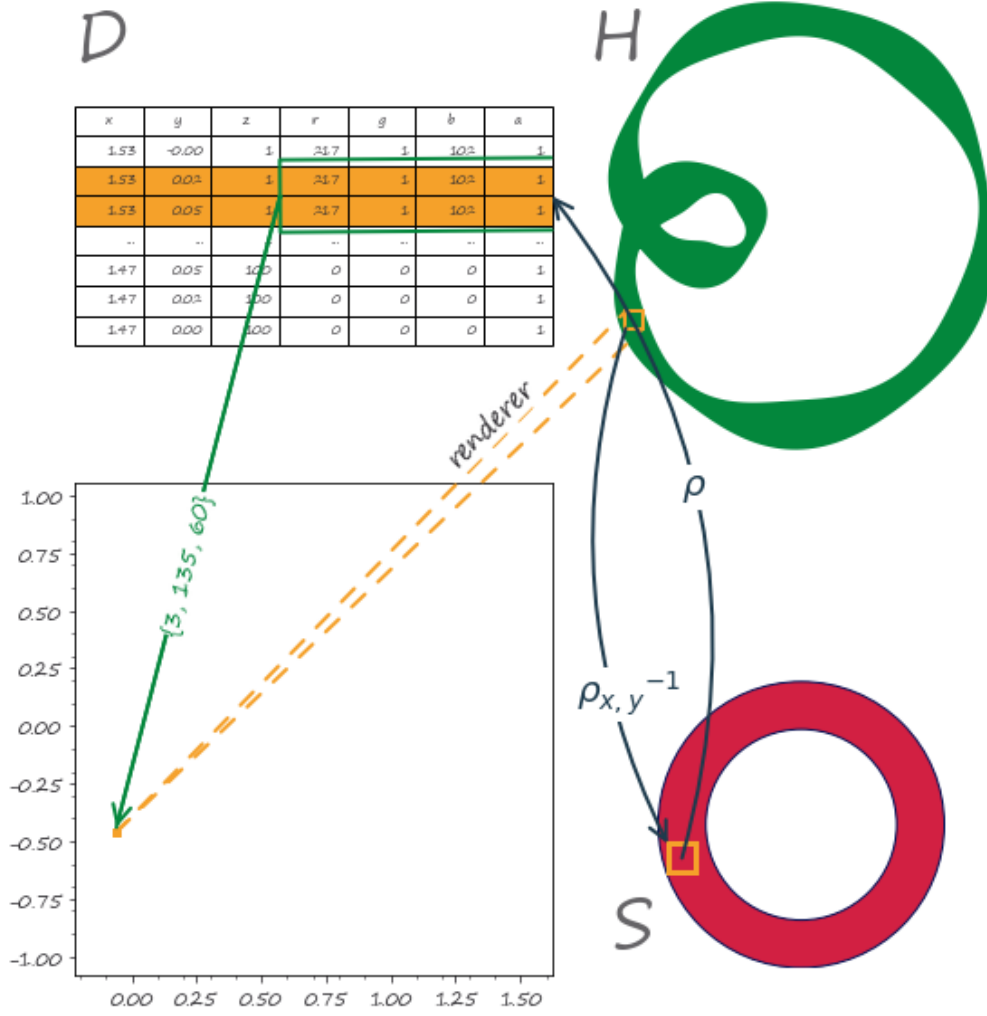


Figure 14: To render a graphic, a pixel  $p$  is selected in the display space, which is defined in the same coordinates as the  $x$  and  $y$  components in  $D$  via the renderer. In  $H$  the inverse mapping  $\rho_{xy}(p)$  returns a region  $S_p \subset S$ .  $\rho(S_p)$  returns a set of points  $(x, y, r, g, b) \in D$  that lie over  $S_p$ . The integral over the  $(r, g, b)$  pixels specifies that the pixel should be green

393 As shown in Figure 14, a pixel  $p$  in the output space, drawn in yellow, is selected and  
394 mapped, via the renderer, into a region on  $H$ . The region on  $H$  corresponds to a region  
395  $S_p \subset S$  via the inverse mapping  $\rho_{xy}(p)$ . The base space  $S$  is an annulus to match the  
396 topology of the graphic idealized in  $H$ . The section  $\rho(S_p)$  then maps into the fiber  $D$  over  
397  $S_p$  to obtain the set of points in  $D$ , here represented as a table, that correspond to that

398 section. The integral over the pixel components of this set of points in the fiber yields  
 399  $\{3, 135, 60\}$  the actual color of the pixel. In general,  $\rho$  is an abstraction of rendering.  
 400 In very broad strokes  $\rho$  can be a specification such as PDF[72], SVG[73], or an OpenGL  
 401 scene graph[74]. Alternatively,  $\rho$  can be a rendering engine such as cairo[75] or AGG[76].  
 402 Implementation of  $\rho$  is out of scope for this work,

### 403 3.3 Artist

The topological artist  $A$  is how we model the building block component that transforms data into a graphic. The artist  $A$  is a map from the sheaf on a data bundle  $E$  which is  $\mathcal{O}(E)$  to the sheaf on the graphic bundle  $H$ ,  $\mathcal{O}(H)$ .

$$A : \mathcal{O}(E) \rightarrow \mathcal{O}(H) \quad (15)$$

The artist preserves *continuity* through the  $\xi$  map discussed in [subsubsection 3.2.2](#) and is an *equivariant* map because it carries a homomorphism of monoid actions [77]

$$\varphi : M \rightarrow M' \quad (16)$$

Given  $M$  on data  $\mathcal{E}$  and  $M'$  on graphic  $\mathcal{H}$ , we propose that artists  $\mathcal{A}$  are equivariant maps

$$A(m \cdot r) = \varphi(m) \cdot A(r) \quad (17)$$

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

The monoid equivariant map has two stages: the encoders  $\nu : E' \rightarrow V$  convert the data components to visual components, and the assembly function  $Q : \xi^*V \rightarrow H$  composites the

fiber components of  $\xi^*V$  into a graphic in  $H$ .

$$\begin{array}{ccccc}
 E' & \xrightarrow{\nu} & V & \xleftarrow{\xi^*} & \xi^*V & \xrightarrow{Q} & H \\
 & \searrow \pi & \downarrow \pi & & \xi^* \downarrow \pi & \swarrow \pi & \\
 & & K & \xleftarrow{\xi} & S & & 
 \end{array} \tag{18}$$

407  $\xi^*V$  is the visual bundle  $V$  pulled back over  $S$  via the equivariant continuity map  $\xi : S \rightarrow K$   
 408 introduced in [subsubsection 3.2.2](#). The functional decomposition of the visualization artist  
 409 in [Equation 18](#) facilitates building reusable components at each stage of the transformation  
 410 because the equivariance constraints are defined on  $\nu$ ,  $Q$ , and  $\xi$ . We name this map the artist  
 411 as that is the analogous part of the Matplotlib[5] architecture that builds visual elements.

### 412 3.3.1 Visual Fiber Bundle $V$

We introduce a visual bundle  $V$  to store the mappings of the data components into components of the graphic. The visual bundle  $(V, K, \pi, P)$  is the space of possible parameters of a visualization type, such as a scatter or line plot. As with the data and graphic bundles, the visual bundle is defined by the projection map  $\pi$

$$\begin{array}{ccc}
 P & \hookrightarrow & V \\
 & & \pi \downarrow \Big)^\mu \\
 & & K
 \end{array} \tag{19}$$

413 where  $\mu$  is the visual variable encoding, as described by Bertin [9], of the data section  $\tau$ .  
 414 The visual fiber  $P$  is defined in terms of the input parameters of the visualization library's  
 415 plotting functions; by making these parameters explicit components of the fiber, we can  
 416 build consistent definitions and expectations of how these parameters behave.

$\nu_i$	$\mu_i$	$\text{codomain}(\nu_i) \subset P_i$
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

417 A section  $\mu$  is a tuple of visual values that specifies the visual characteristics of a part  
 418 of the graphic. For example, given a fiber of  $\{x, y, color\}$  one possible section could be  
 419  $\{.5, .5, (255, 20, 147)\}$ . The  $\text{codomain}(\nu_i)$  determines which monoids can act on  $P_i$ . These  
 420 fiber components are implicit in the library, as seen in [Table 1](#), and by making them explicit  
 421 as components of the fiber we can build consistent definitions and expectations of how these  
 422 parameters behave.

### 423 3.3.2 Visual Encoders $\nu$

We define the visual transformers  $\nu$

$$\{\nu_0, \dots, \nu_n\} : \{\tau_0, \dots, \tau_n\} \mapsto \{\mu_0, \dots, \mu_n\} \quad (20)$$



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 there is a monoid  $M_i'$  on  $V_i$ , then there is a monoid homomorphism from  $\varphi : M_i \rightarrow M_i'$  that  $\nu$  must preserve. As mentioned in [subsubsection 3.1.2](#), monoid actions define the structure on the fiber components and are therefore the basis for equivariance. A validly constructed  $\nu$  is one where the diagram of the monoid transform  $m$  commutes

$$\begin{array}{ccc} E_i & \xrightarrow{\nu_i} & V_i \\ m_r \downarrow & & \downarrow m_v \\ E_i & \xrightarrow{\nu_i} & V_i \end{array} \quad (21)$$

such that applying equivariant monoid actions to  $E_i$  and  $V_i$  preserves the map  $\nu_i : E_i \rightarrow V_i$ . 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  $P_i$ ; therefore an equivariant  $\nu_i$  must satisfy the constraint

$$\nu_i(m_r(E_i)) = \varphi(m_r)(\nu_i(E_i)) \quad (22)$$

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

430 The translation from weather state data to visual representation as umbrella emoji in  
 431 [Figure 15a](#) is an invalid visual encoding map  $\nu$  because it is not homomorphic. This is  
 432 because the monotonic condition  $rain \geq storm \implies \nu(rain) \geq \nu(storm)$  is not met since  
 433  $\nu(rain) \leq \nu(storm)$ . To satisfy the monotonic condition for  $rain \geq storm$ , either red arrow  
 434 in [Figure 15a](#) would have to go in a different direction. On the other hand, the mapping  
 435 from weather state to umbrella in [Figure 15b](#) is a homomorphism since  $\nu(rain) = \nu(storm)$   
 436 satisfies the monotonic condition of  $rain \geq storm$ . [Figure 15](#) is an example of how the

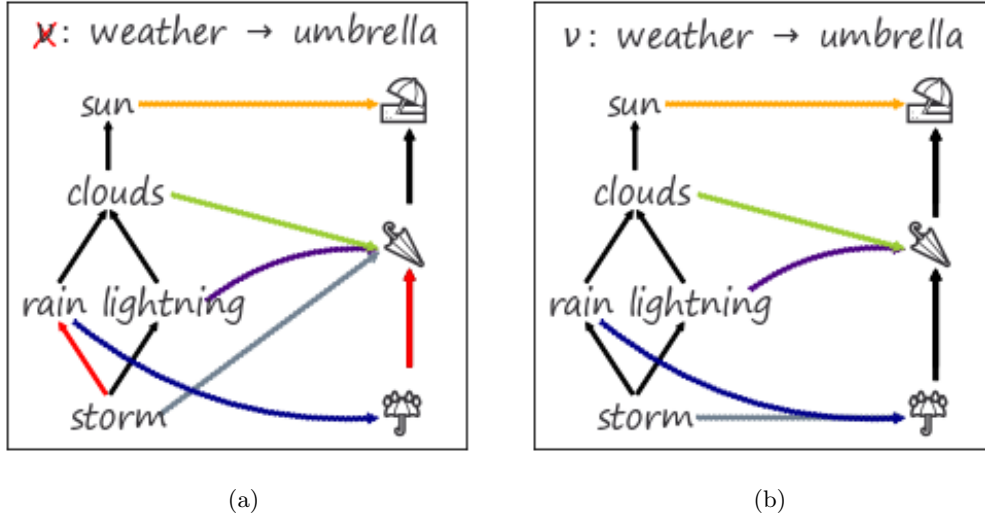


Figure 15: The map from data component to visual component in Figure 15a is not homomorphic, and therefore invalid, because  $rain \geq storm$  is mapped to elements with the reverse ordering  $v(storm) \geq v(storm)$ . In contrast, the mapping in Figure 15b is valid since  $v(storm) = v(rain)$  satisfies the condition  $v(storm) \geq v(storm)$

437 model supports partially ordered data components, which was a motivation for defining  
 438 equivariance as monoid homomorphisms.

scale	group	constraint
nominal	permutation	if $r_1 \neq r_2$ then $\nu(r_1) \neq \nu(r_2)$
ordinal	monotonic	if $r_1 \leq r_2$ then $\nu(r_1) \leq \nu(r_2)$
interval	translation	$\nu(x + c) = \nu(x) + c$
ratio	scaling	$\nu(xc) = \nu(x) * c$

Table 2

439 The Stevens measurement types[55], listed in Table 2, are specified in terms of groups,  
 440 which are monoids with invertible operations[78]. Despite critiques of the scales[79, 80], we  
 441 believe it is critical for the model to include the measurement scales since they are com-  
 442 monly used in visualization to classify components [22, 46]. By specifying the equivariance  
 443 constraints on  $\nu$  we can guarantee that the stage of the artist that transforms data compo-  
 444 nents into visual representations is equivariant. These constraints guide the implementation  
 445 of reusable component transformers  $\nu$  that are composed when generating the graphic.

### 446 3.3.3 Visualization Assembly

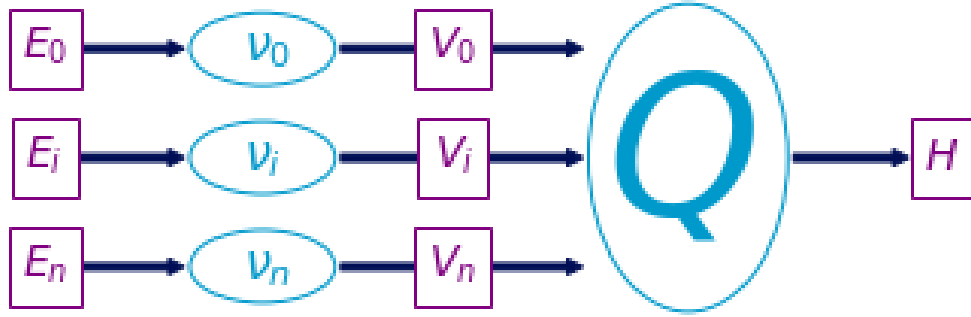


Figure 16: The transform functions  $\nu_i$  convert data  $\tau_i \in E$  to visual characteristics  $\mu_i \in V$ , then  $Q$  assembles  $\mu_i$  into a graphic  $\rho \in H$ .

447 The transformation from data into graphic is analogous to a map-reduce operation;  
 448 as illustrated in ??, data components  $E_i$  are mapped into visual components  $V_i$  that are  
 449 reduced into a graphic in  $H$ . The space of all graphics that  $Q$  can generate is the subset  
 450 of graphics reachable via applying the reduction function  $Q(\Gamma(V)) \in \Gamma(H)$  to the visual  
 451 section  $\mu \in \Gamma(V)$ . The full space of graphics is not necessarily equivariant; therefore we  
 452 formalize the constraints on  $Q$  such that it produces structure preserving graphics.

453 We formalize the expectation that visualization generation functions parameterized in  
 454 the same way should generate the same functions as the equivariant map  $Q : \mu \mapsto \rho$ . We

455 then define the constraint on  $Q$  such that if  $Q$  is applied to two visual sections  $\mu$  and  $\mu'$   
 456 that generate the same  $\rho$  then the output of  $\mu$  and  $\mu'$  acted on by the same monoid  $m$   
 457 must be the same. We do not define monoid actions on all of  $\Gamma(H)$  because there may be  
 graphics  $\rho \in \Gamma(H)$  for which we cannot construct a valid mapping from  $V$ . Lets call the

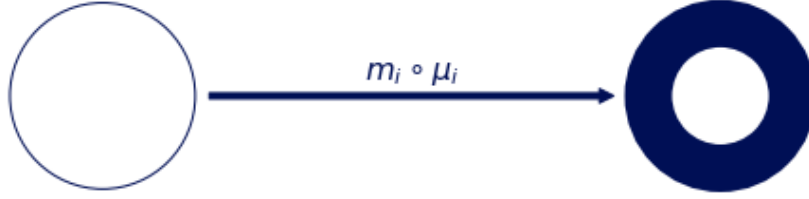


Figure 17: These two glyphs are generated by the same annulus  $Q$  function. The monoid action  $m_i$  on edge thickness  $\mu_i$  of the first glyph yields the thicker edge  $\mu_i'$  in the second glyph.

458

459 visual representations of the components  $\Gamma(V) = X$  and the graphic  $Q(\Gamma(V)) = Y$

**Proposition 1.** *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*

$$Q(\mu) = Q(\mu') \implies Q(m \circ \mu) = Q(m \circ \mu') \quad (23)$$

460 then a monoid action on  $Y$  can be defined as  $m \circ \rho = \rho'$ . If and only if  $Q$  satisfies  
 461 *Equation 23*, we can state that the transformed graphic  $\rho' = Q(m \circ \mu)$  is equivariant to a  
 462 monoid action applied on  $Q$  with input  $\mu \in Q^{-1}(\rho)$  that must generate valid  $\rho$ .

463 For example, given fiber  $P = (xpos, ypos, color, thickness)$ , then sections  $\mu = (0, 0, 0, 1)$   
 464 and  $Q(\mu) = \rho$  generates a piece of the thin hollow circle. The action  $m = (e, e, e, x + 2)$ ,  
 465 where  $e$  is identity, translates  $\mu$  to  $\mu' = (e, e, e, 3)$  and the corresponding action on  $\rho$  causes  
 466  $Q(\mu')$  to be the thicker circle in [Figure 17](#).

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

$$J = \{j \in K \text{ s. t. } \exists \gamma \text{ s.t. } \gamma(0) = k \text{ and } \gamma(1) = j\} \quad (24)$$

where the path[81]  $\gamma$  from  $k$  to  $j$  is a continuous function from the interval  $[0,1]$ . We define the glyph as the graphic generated by  $Q(S_j)$

$$H \xrightleftharpoons[\rho(S_j)]{} S_j \xrightleftharpoons[\xi^{-1}(J)]{\xi(s)} J_k \quad (25)$$

such that for every glyph there is at least one corresponding region on  $K$ , in keeping with the definition of glyph as any visually differentiable element put forth by Ziemkiewicz and Kosara[82]. The primitive point, line, and area marks[9, 83] are specially cased glyphs.

### 3.3.4 Assembly $Q$

Given the continuities described in 13, we illustrate a minimal  $Q$  that will generate the most minimal visualizations associated with those continuities: non-overlapping scatter points, a non-infinitely thin line, and an image.

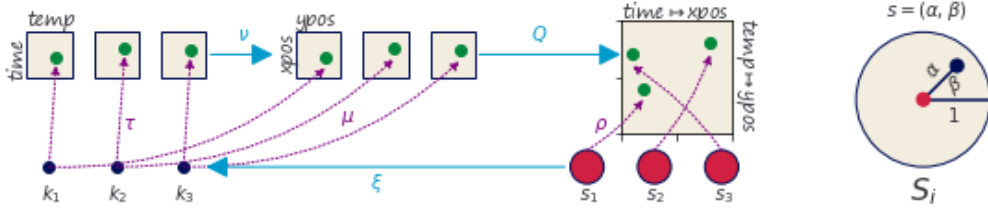


Figure 18: The data is discrete points (temperature, time). Via  $\nu$  these are converted to (xpos, ypos) and pulled over discrete  $S$ . These values are then used to parameterize  $\rho$  which returns a color based on the parameters (xpos,ypos) and position  $\alpha, \beta$  on  $S_k$  that  $\rho$  is evaluated on.

The scatter plot in Figure 18 can be defined as

$$Q(xpos, ypos)(\alpha, \beta) \quad (26)$$

with a constant  $size$  and color  $\rho_{RGB} = (0, 0, 0)$  that are defined as part of  $Q$ . The position of this swatch of color can be computed relative to the location on the disc  $(\alpha, \beta) \in S_k$  as

shown in Figure 18

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

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

474 such that  $\rho(s) = (x, y, 0, 0, 0)$  colors the point (x,y) black.

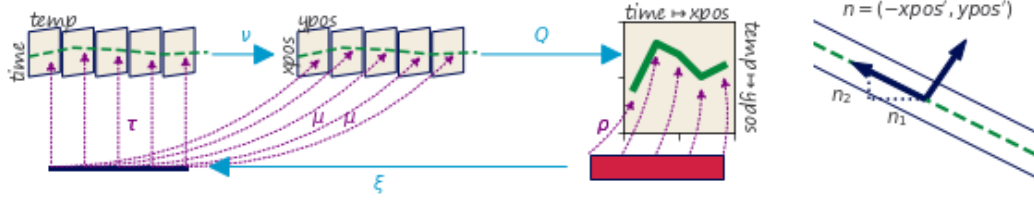


Figure 19: The line fiber  $(time, temp)$  is thickened with the derivative  $(time', temperature')$  because that information will be necessary to figure out the tangent to the point to draw a line. This is because the line needs to be pushed perpendicular to the tangent of  $(xpos, ypos)$ . The data 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, the line plot

$$Q(xpos, \hat{n}_1, ypos, \hat{n}_2)(\alpha, \beta) \quad (27)$$

in ?? has a  $\xi$  function that is not only parameterized on  $k$  but also on the  $\alpha$  distance along  $k$  and corresponding region in  $S$ . As shown in ??, 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; therefore the artist takes as input the jet bundle [84, 85]  $\mathcal{J}^2(E)$  which is the data  $E$  and the first and second derivatives of  $E$ . The magnitude of the slope is  $|n| = \sqrt{n_1^2 + n_2^2}$  such that the normal is  $\hat{n}_1 = \frac{n_1}{|n|}$ ,  $\hat{n}_2 = \frac{n_2}{|n|}$  which yields components of  $\rho$

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

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



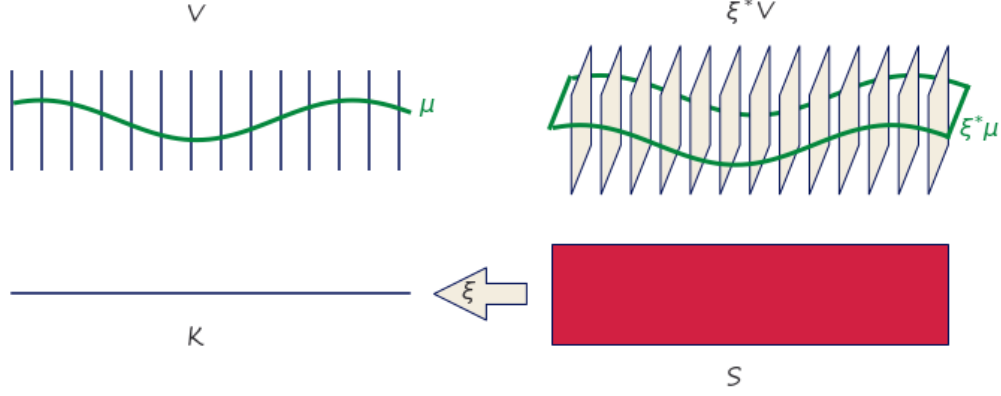


Figure 21: Because the pullback of the visual bundle  $\xi^*V$  is the replication of a  $\mu$  over all points  $s$  that map back to a single  $k$ , we can construct a  $\hat{Q}$  on  $\mu$  over  $k$  that will fabricate the  $Q$  for the equivalent region of  $s$  associated to that  $k$

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

485 As shown in **Figure 21**, given the section  $\xi^*\mu$  pulled back from  $\mu$  and the point  $s \in \xi^{-1}(k)$ ,  
 486 there is a direct map  $(k, \mu(k)) \mapsto (s, \xi^*\mu(s))$  from  $\mu$  over  $k$  to the section  $\xi^*\mu$  over  $s$ . This  
 487 means that the pulled back section  $\xi^*\mu(s) = \xi^*(\mu(k))$  is the section  $\mu$  copied over all  $s$  such  
 488 that  $\xi^*\mu$  is identical for all  $s$  where  $\xi(s) = k$ . In **Figure 21** each dot on  $P$  is equivalent to  
 489 the line on  $P^*\mu$ .

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  $\xi^*\mu$  then the graphic section  $\rho$  evaluated on a visual region  $s$

$$\rho(s) := Q(\xi^*\mu)(s) \quad (30)$$

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

$$\hat{Q}(\mu(k))(s) := Q((\xi^*\mu)(s)) \quad (31)$$



where  $\hat{Q}$  with input  $\mu$  is defined to  $Q$  that takes as input the copied section  $\xi^*\mu$  such that both functions are evaluated over the same location  $\xi^{-1}(k) = s$  in the base space  $S$ . Factoring out  $s$  from Equation 31 yields

$$\hat{Q}(\mu(k)) = Q(\xi^*\mu) \quad (32)$$

where  $Q$  is no longer bound to input but  $\hat{Q}$  is still defined in terms of  $K$ . In fact,  $\hat{Q}$  is a map from visual space to graphic space  $\hat{Q} : \Gamma(V) \rightarrow \Gamma(H)$  locally over  $k$  such that it can be evaluated on a single visual record  $\hat{Q} : \Gamma(V_k) \rightarrow \Gamma(H \mid_{\xi^{-1}(k)})$ . This allows us to construct a  $\hat{Q}$  that only depends on  $K$ , such that for each  $\mu(k)$  there is part of  $\rho \mid_{\xi^{-1}(k)}$ . The construction of  $\hat{Q}$  allows us to retain the functional map reduce benefits of  $Q$  without having to majorly restructure the existing pipeline for libraries that delegate the construction of  $\rho$  to a backend such as Matplotlib.

### 3.3.5 Composition of Artists: +

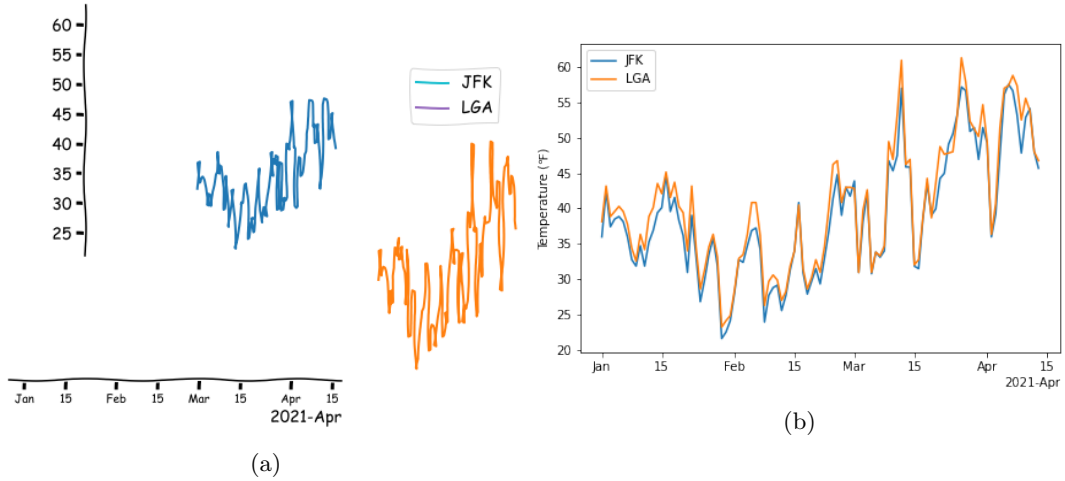


Figure 22: Each of the visual elements in ?? is generated via a unique artist  $A$ . In Figure 22a, they are added to the image independent of the other elements, creating an incoherent visualization. In Figure 22b, these artists are composited before being added to the image. Disjoint union of  $E$  aligns the two timeseries with the x and y axis so all these elements use a shared coordinate system. A more complex composition dictates that the legend is connected to the  $E$  such that it must use the same color as the data it is identifying.

Visualizations with a single artist do not provide much information, so we define addition operators for generating more complex visualizations. Given the family of artists  $(E_i : i \in I)$  on the same image, the  $+$  operator

$$+ := \sqcup_{i \in I} E_i \quad (33)$$

498 defines a simple composition of artists. For example, the components in [Figure 22a](#) are  
 499 each generated by different artists, and a visualization of solely the x axis is rarely all that  
 500 useful. In [Figure 22a](#), these artists are all added to the image independently of the other  
 501 and therefore there are no constraints on how they are generated in conjunction with each  
 502 other. In [Figure 22b](#), the data is joined via disjoint union; doing so aligns the components  
 503 in  $F$  such the  $\nu$  to the same component in  $P$  targets the same coordinate system. When  
 504 artists share a base space  $K_2 \hookrightarrow K_1$ , a composition operator can be defined such that the  
 505 artists are acting on different components of the same section. This type of composition  
 506 is important for visualizations where elements update together in a consistent way, such as  
 507 multiple views [\[86, 87\]](#) and brush-linked views[\[88, 89\]](#).

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

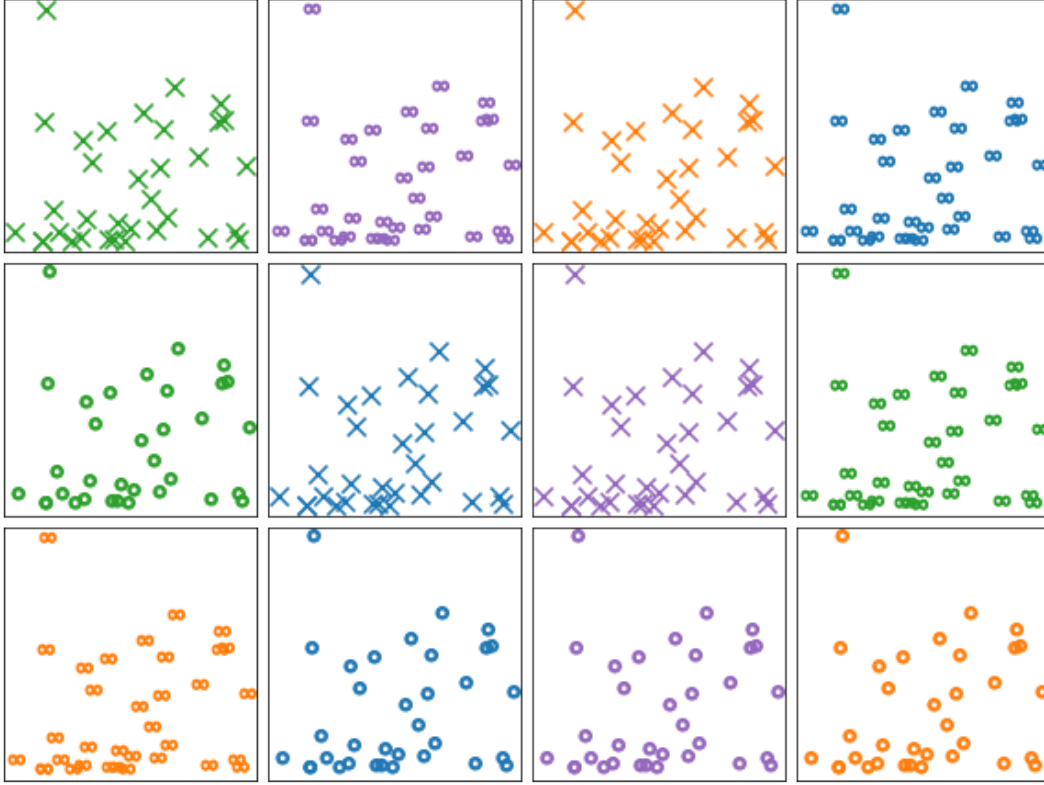


Figure 23: Each scatter plot is generated via a unique artist function  $A_i$ , but they only differ in aesthetic styling. Therefore, these artists are all members of an equivalence class  $A_i \in A'$

Representational invariance, as defined by Kindlmann and Scheidegger, is the notion that visualizations are equivalent if changing the visual representation, such as colors or shapes, does not change the meaning of the visualization[12]. We propose that visualizations are invariant if they are generated by artists that are members of an equivalence class

$$\{A \in A' : A_1 \equiv A_2\}$$

For example, every scatter plot in Figure 23 is a scatter of the same datasets mapped to the  $x$  position and  $y$  position in the same way. The scatter plots only differ in the choice of constant visual literals, differing in color and marker shape. Each scatter is generated by an

512 artist  $A_i$ , and every scatter is generated by a member of the equivalence class  $A_i \in \mathcal{A}$ . Since  
 513 it is impractical to implement a new artist for every single graphic, the equivalence class  
 514 provides a way to evaluate an implementation of a generalized artist. Given equivalent, but  
 515 no necessarily identical,  $\nu$ ,  $Q$ , and  $\xi$ , two artists are equivalent. This criteria also allows for  
 516 comparing artists across libraries.

## 517 4 Prototype Implementation: Matplottoy

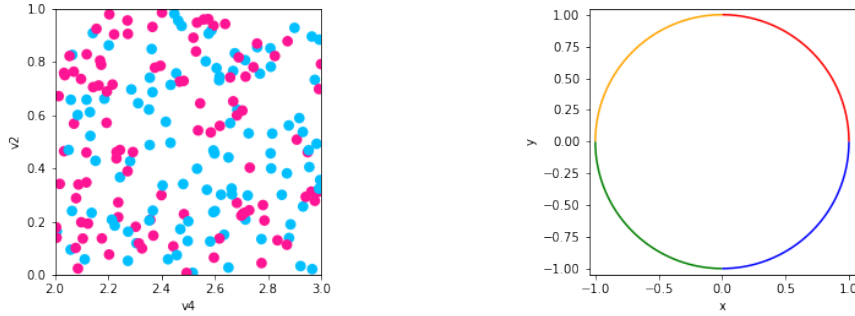


Figure 24: Scatter plot and line plot implemented using prototype artists and data models, building on Matplotlib rendering.

518 To prototype our model, we implemented the artist classes for the scatter and line plots  
 519 shown in figure 24 because they differ in every attribute: different visual channels  $\nu$  that  
 520 composite to different marks  $Q$  with different continuities  $\xi$ . We make use of the Matplotlib  
 521 figure and axes artists [4, 5] so that we can initially focus on the data to graphic transfor-  
 522 mations. We also exploit the Matplotlib transform stack to transform data coordinates into  
 523 screen coordinates. To generate the images in figure 24, we instantiate `fig`, `ax` artists that  
 524 will contain the new `Point`, `Line` primitive objects we implemented based on our topology  
 525 model.

---

```

1 fig, ax = plt.subplots()
2 artist = Point(data, transforms)
3 ax.add_artist(artist)
```

---

```
1 fig, ax = plt.subplots()
2 artist = Line(data, transforms)
3 ax.add_artist(artist)
```

---

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

#### 536 4.1 Artist Class $A'$

537 The artist is the piece of the Matplotlib architecture that constructs an internal representa-  
 538 tion of the graphic that the render then uses to draw the graphic. In the prototype artist,  
 539 `transform` is a dictionary of the form `{parameter:(variable, encoder)}` where parame-  
 540 ter is a component in  $P$ , variable is a component in  $F$ , and the  $\nu$  encoders are passed in as  
 541 functions or callable objects. The data bundle  $E$  is passed in as a `data` object. By binding  
 542 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)

```

---

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

553 The components of the visual object, denoted by the Python unpacking convention  
554 `**visual` are then passed into the `assemble` function that is  $\hat{Q}$ . This assembly function  
555 is responsible for generating a representation such that it could be serialized to recreate a  
556 static version of the graphic. Although `assemble` could be implemented outside the class  
557 such that it returns an object the artist could then parse to set attributes, the attributes are  
558 directly set here to reduce indirection. This artist is not optimized because we prioritized  
559 demonstrating the separability of  $\nu$  and  $\hat{Q}$ . The last step in the artist function is handing

560 itself off to the renderer. The extra `*arg, **kwargs` arguments in `__init__`, `draw` are  
 561 artifacts of how these objects are currently implemented in Matplotlib.

562 The `Point` artist builds on `collection` artists because collections are optimized to ef-  
 563 ficiently draw a sequence of primitive point and area marks. In this prototype, the scatter  
 564 marker shape is fixed as a circle, and the only visual fiber components are `x` and `y` position,  
 565 size, and the facecolor of the marker. We only show the `assemble` function here because  
 566 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)

```

---

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

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



578 The main difference between the `Point` and `Line` objects is in the `assemble` function  
 579 because line has different continuity from scatter and is represented by a different type of  
 580 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)

```

---

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

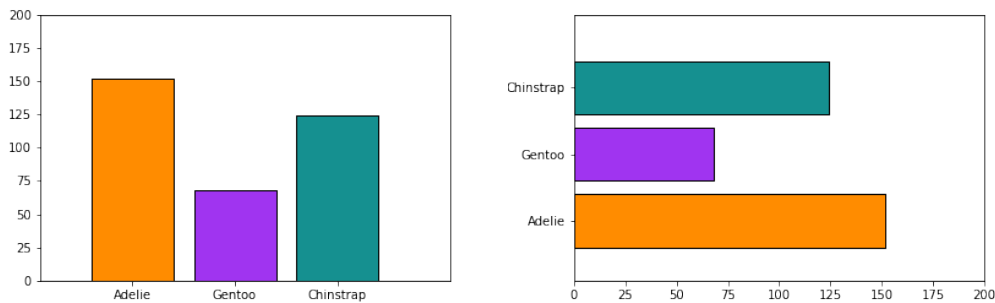


Figure 25: Frequency of Penguin types visualized as discrete bars.

588 The bar charts in figure 25 are generated with a `Bar` artist. The artist has required  
589 visual parameters  $P$  of (position, length), and an additional parameter `orientation` which  
590 controls whether the bars are arranged vertically or horizontally. This parameter only applies  
591 holistically to the graphic and never to individual data parameters, and highlights how the  
592 model encourages explicit differentiation between parameters in  $V$  and graphic parameters  
593 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)

```

---

594 The **draw** method here has a more complex unpacking of visual encodings to support passing  
 595 in visual component data directly. This is vastly simplifies building composite objects as  
 596 the alternative would be higher order functions that take as input the transforms passed in  
 597 by the user. This construction supports a constant visual parameter, an identity transform  
 598 where the value is the same in  $E$  and  $V$ , and setting the visual component directly. The  
 599 **assemble** function constructs bars and sets their face and edge colors. The **make\_bars**  
 600 function converts the input position and length to the coordinates of a rectangle of the given  
 601 width. Defaults are provided for 'width' and 'floor' to make this function more reusable.  
 602 Typically the defaults are used for the type of chart shown in figure 25, but these visual  
 603 variables are often set when building composite versions of this chart type as discussed in  
 604 section 4.4.

## 605 4.2 Encoders $\nu$

606 As mentioned above, the encoding dictionary is specified by the visual fiber component, the  
 607 corresponding data fiber component, and the mapping function. The visual parameter serves  
 608 as the dictionary key because the visual representation is constructed from the encoding  
 609 applied to the data  $\mu = \nu \circ \tau$ . For the scatter plot, the mappings for the visual fiber  
 610 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)}}

```

---

611 where the position  $(x, y)$   $\nu$  transformers are identity functions. The size  $s$  transformer is not  
 612 acting on a component of  $F$ , instead it is a  $\nu$  that returns a constant value. While size could  
 613 be embedded inside the **assemble** function, it is added to the transformers to illustrate user  
 614 configured visual parameters that could either be constant or mapped to a component in  $F$ .

615 The identity and constant  $\nu$  are explicitly implemented here to demonstrate their implicit  
616 role in the visual pipeline, but they are somewhat optimized away in `Bar`. More complex  
617 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]
```

---

618 where `__init__` can validate that the output of the  $\nu$  is a valid element of the  $P$  com-  
619 ponent the  $\nu$  function is targeting. Creating a callable class also provides a simple way to  
620 swap out the specific (data, value) mapping without having to reimplement the validation  
621 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)
```

---

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

### 625 4.3 Data $E$

626 The data input into the **Artist** will often be a wrapper class around an existing data  
627 structure. This wrapper object must specify the fiber components  $F$  and connectivity  $K$   
628 and have a **view** method that returns an atomic object that encapsulates  $\tau$ . The object  
629 returned by the view must be key valued pairs of {**component name** : **component section**}  
630 where each section is a component as defined in equation ?? . To support specifying the fiber  
631 bundle, we define a **FiberBundle** data class[90]

---

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

---

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

638 To generate the scatter plot in figure 24, we fully specify a dataset with random keys  
639 and values in a section chosen at random from the corresponding fiber component. The  
640 fiberbundle **FB** is a class level attribute since all instances of **VertexSimplex** come from the  
641 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

```

---

642 The view method returns a dictionary where the key is a fiber component name and the  
643 value is a list of values in the fiber component. The table is built one call to the section  
644 method `tau` at a time, guaranteeing that all the fiber component values are over the same  
645  $k$ . Table has a `get` method as it is a method on Python dictionaries. In contrast, the line  
646 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
7     def __init__(self, num_edges=4, num_samples=1000):
8         self.keys = range(num_edge) #edge id
9         # distance along edge
10        self.distances = np.linspace(0,1, num_samples)
11        # half generlized representation of arcs on a circle
12        self.angle_samples = np.linspace(0, 2*np.pi, len(self.keys)+1)
13
14    @staticmethod
15    def _color(edge):
16        colors = ['red','orange', 'green','blue']
17        return colors[edge%len(colors)]
18
19    @staticmethod
20    def _xy(edge, distances, start=0, end=2*np.pi):
21        # start and end are parameterizations b/c really there is
22        angles = (distances *(end-start)) + start
23        return np.cos(angles), np.sin(angles)
24
25    def tau(self, k): #will fix location on page on revision
26        x, y = self._xy(k, self.distances,
27                        self.angle_samples[k], self.angle_samples[k+1])
28        color = self._color(k)
29        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)

```

---

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

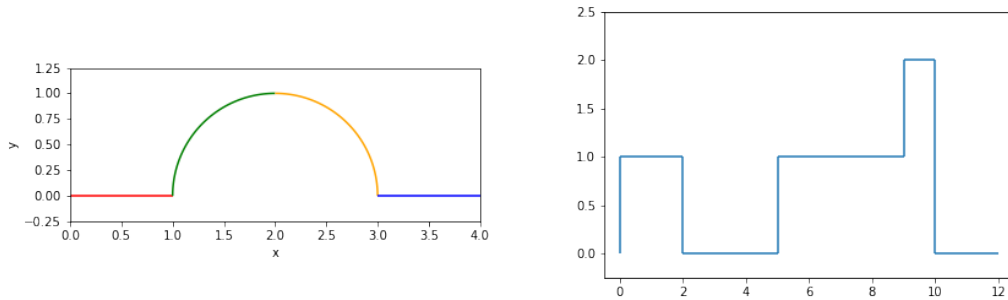


Figure 26: Continuous and discontinuous lines as defined via the same data model, and generated with the same `A'Line`

653 The graphics in figure 26 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

```

---

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

---

```

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

```

---

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

---

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

---

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

## 662 4.4 Case Study: Penguins

663 For this case study, we use the Palmer Penguins dataset[91, 92] since it is multivariate and  
664 has a varying number of penguins. We use a version of the data packaged as a pandas  
665 dataframe[93] since that is a very commonly used Python labeled data structure. The  
666 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
```

---

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

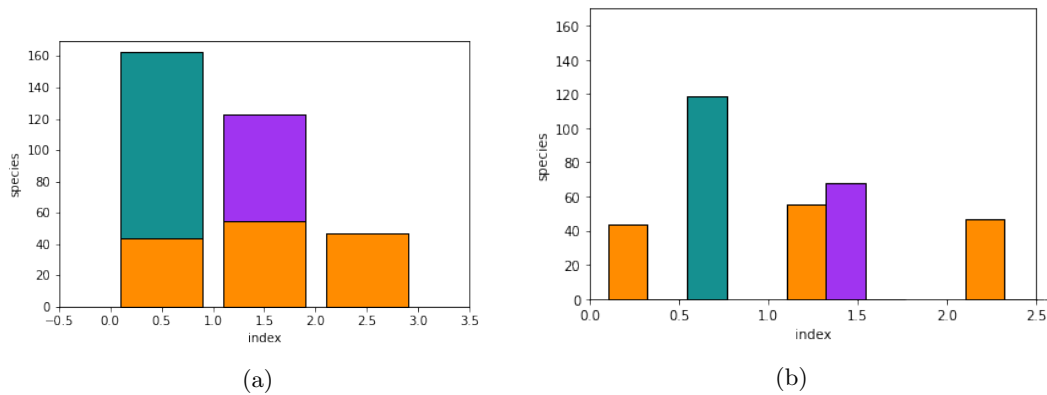


Figure 27: Penguin count disaggregated by island and species

672 The stacked and grouped bar charts in figure 27 are both out of `Bar` artists such that  
 673 the difference between `StackedBar` and `GroupedBar` is specific to the ways in which the  
 674 `Bar` are stitched together. These two artists have identical constructors and `draw` methods.  
 675 As with `Bar`, the orientation is set in the constructor. In both these artists, we separate  
 676 the transforms applied to only one component and the case `mtransforms` where the same  
 677 transform is applied to multiple components such that `V` has multiple components that map  
 678 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)

```

---

679 Since all the visual transformation is passed through to `Bar`, the `draw` method does not  
680 do any visual transformations. In `StackedBar` the `view` is used to adjust the `floor` for  
681 every subsequent bar chart since a stacked bar chart is bar chart area marks concatenated  
682 together in the `length` parameter. In contrast, `GroupedBar` does not even need the `view`, but

683 instead keeps track of the relative position of each group of bars in the visual only variable  
684 `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)
```

---

685 Since the only difference between these two glyphs is in the composition of `Bar`, they take  
686 in the exact same transform specification dictionaries. The `transform` dictionary dictates  
687 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}}]
```

---

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

---

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

---

692 but differ in assembly  $\hat{Q}$ . By decomposing the architecture into data, visual encoding,  
693 and assembly steps, we are able to build components that are more flexible and also more self  
694 contained than the existing code base. While very rough, this API demonstrates that the  
695 ideas presented in the math framework are implementable. For example, the `draw` function  
696 that maps most closely to  $A$  is functional, with state only being necessary for bookkeeping  
697 the many inputs that the function requires. In choosing a functional approach, if not  
698 implementation, we provide a framework for library developers to build reusable encoder  
699  $\nu$  assembly  $\hat{Q}$  and artists  $A$ . We argue that if these functions are built such that they  
700 are equivariant with respect to monoid actions and the graphic topology is a deformation  
701 retraction of the data topology, then the artist by definition will be a structure and property  
702 preserving map from data to graphic.

## 703 5 Discussion

704 This work contributes a mathematical description of the mapping  $A$  from data to visual rep-  
705 resentation. Combining Butler’s proposal of a fiber bundle model of visualization data with  
706 Spivak’s formalism of schema lets this model support a variety of datasets, including discrete  
707 relational tables, multivariate high resolution spatio temporal datasets, and complex net-  
708 works. Decomposing the artist into encoding  $\nu$ , assembly  $Q$ , and reindexing  $\xi$  provides the  
709 specifications for producing visualization where the structure and properties match those  
710 of the input data. These specifications are that the graphic must have continuity equiva-

711 lent to the data, and that the visual characteristics of the graphics are equivariant to their  
712 corresponding components under monoid actions. This model defines these constraints on  
713 the transformation function such that they are not specific to any one type of encoding or  
714 visual characteristic. Encoding the graphic space as a fiber bundle provides a structure rich  
715 abstraction of the target graphical design in the target display space.

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

## 727 5.1 Limitations

728 So far this model has only been worked out for a single data set tied to a primitive mark,  
729 but it should be extensible to compositing datasets and complex glyphs. The examples and  
730 prototype have so far only been implemented for the static 2D case, but nothing in the math  
731 limits to 2D and expansion to the animated case should be possible because the model is  
732 formalized in terms of the sheaf. While this model supports equivariance of figurative glyphs  
733 generated from parameters of the data[94, 95], it currently does not have a way to evaluate  
734 the semantic accuracy of the figurative representation. Effectiveness is out of scope for this  
735 model because it is not part of the structure being preserved, but potentially a developer  
736 building a domain specific library with this model could implement effectiveness criteria in  
737 the artists. Also, even though the model is designed to be backend and format independent,  
738 it has only really been tested against PNGs rendered with the AGG backend. It is especially



unknown how this framework interfaces with high performance rendering libraries such as  
OpenGL[74]. Because this model has been limited to the graphic design space, it does not  
address the critical task of laying out the graphics in the image

This model and the associated prototype is deeply tied to Matplotlib’s existing architecture. While the model is expected to generalize to other libraries, such as those built on Mackinlay’s APT framework, this has not been worked through. In particular, Mackinlay’s formulation of graphics as a language with semantic and syntax lends itself a declarative interface[96], which Heer and Bostock use to develop a domain specific visualization language that they argue makes it simpler for designers to construct graphics without sacrificing expressivity [18]. Similarly, the model presented in this work formulates visualization as equivariant maps from data space to visual space, and is designed such that developers can build software libraries with data and graphic topologies tuned to specific domains.

## 5.2 Future Work

While the model and prototype demonstrate that generation of simple marks from the data, 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[97–99] and also encoding a separate heatmap[100, 101] artist that consumes 1D discrete data. A second important proof of concept artist is a boxplot[102] 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

766 of how interactivity would work in this model. In summary, the proposed scope of work for  
767 the dissertation is

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

776 Other potential tasks for future work is implementing a data object for a non-trivial fiber  
777 bundle and exploiting the models section level formalism to build distributed data source  
778 models and concurrent artists. This could be pushed further to integrate with topologi-  
779 cal[103] and functional [104] data analysis methods. Since this model formalizes notions of  
780 structure preservation, it can serve as a good base for tools that assess quality metrics[105]  
781 or invariance [12] of visualizations with respect to graphical encoding choices. While this  
782 paper formulates visualization in terms of monoidal action homomorphisms between fiber-  
783 bundles, the model lends itself to a categorical formulation[54, 106] that could be further  
784 explored.

## 785 6 Conclusion

786 An unofficial philosophy of Matplotlib is to support making whatever kinds of plots a user  
787 may want, even if they seem nonsensical to the development team. The topological frame-  
788 work described in this work provides a way to facilitate this graph creation in a rigorous  
789 manner; any artist that meets the equivariance criteria described in this work by definition

790 generates a graphic representation that matches the structure of the data being represented.  
 791 We leave it to domain specialists to define the structure they need to preserve and the maps  
 792 they want to make, and hopefully make the process easier by untangling these components  
 793 into separate constrained maps and providing a fairly general data and display model.

## 794 References

- 795 [1] Krist Wongsuphasawat. *Navigating the Wide World of Data Visualization Libraries*  
 796 *(on the Web)*. 2021.
- 797 [2] J. Hughes. “Why Functional Programming Matters”. In: *The Computer Journal* 32.2  
 798 (Jan. 1989), pp. 98–107. ISSN: 0010-4620. DOI: [10.1093/comjnl/32.2.98](https://doi.org/10.1093/comjnl/32.2.98).
- 799 [3] Zhenjiang Hu, John Hughes, and Meng Wang. “How Functional Programming Mat-  
 800 tered”. In: *National Science Review* 2.3 (Sept. 2015), pp. 349–370. ISSN: 2095-5138.  
 801 DOI: [10.1093/nsr/nwv042](https://doi.org/10.1093/nsr/nwv042).
- 802 [4] J. D. Hunter. “Matplotlib: A 2D Graphics Environment”. In: *Computing in Science*  
 803 *Engineering* 9.3 (May 2007), pp. 90–95. ISSN: 1558-366X. DOI: [10.1109/MCSE.2007.](https://doi.org/10.1109/MCSE.2007.55)  
 804 [55](https://doi.org/10.1109/MCSE.2007.55).
- 805 [5] John Hunter and Michael Droettboom. *The Architecture of Open Source Applications*  
 806 *(Volume 2): Matplotlib*. <https://www.aosabook.org/en/matplotlib.html>.
- 807 [6] A. Sarikaya et al. “What Do We Talk About When We Talk About Dashboards?”  
 808 In: *IEEE Transactions on Visualization and Computer Graphics* 25.1 (Jan. 2019),  
 809 pp. 682–692. ISSN: 1941-0506. DOI: [10.1109/TVCG.2018.2864903](https://doi.org/10.1109/TVCG.2018.2864903).
- 810 [7] Michael Friendly. “A Brief History of Data Visualization”. en. In: *Handbook of Data*  
 811 *Visualization*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 15–56. ISBN:  
 812 978-3-540-33036-3 978-3-540-33037-0. DOI: [10.1007/978-3-540-33037-0\\_2](https://doi.org/10.1007/978-3-540-33037-0_2).
- 813 [8] L. Byrne, D. Angus, and J. Wiles. “Acquired Codes of Meaning in Data Visualization  
 814 and Infographics: Beyond Perceptual Primitives”. In: *IEEE Transactions on Visual-*

- ization and Computer Graphics 22.1 (Jan. 2016), pp. 509–518. ISSN: 1077-2626. DOI: 10.1109/TVCG.2015.2467321.
- [9] Jacques Bertin. *Semiology of Graphics : Diagrams, Networks, Maps*. English. Redlands, Calif.: ESRI Press, 2011. ISBN: 978-1-58948-261-6 1-58948-261-1.
- [10] 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.
- [11] Jock Mackinlay. “Automatic Design of Graphical Presentations”. English. PhD Thesis. Stanford, 1987.
- [12] G. Kindlmann and C. Scheidegger. “An Algebraic Process for Visualization Design”. In: *IEEE Transactions on Visualization and Computer Graphics* 20.12 (Dec. 2014), pp. 2181–2190. ISSN: 1941-0506. DOI: 10.1109/TVCG.2014.2346325.
- [13] Ricky Shadrach. *Introduction to Groups*. <https://www.mathsisfun.com/sets/groups-introduction.html>. 2017.
- [14] *Naturalness Principle - InfoVis:Wiki*. [https://infovis-wiki.net/wiki/Naturalness\\_Principle](https://infovis-wiki.net/wiki/Naturalness_Principle).
- [15] 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.
- [16] 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.
- [17] E. H. Chi. “A Taxonomy of Visualization Techniques Using the Data State Reference Model”. In: *IEEE Symposium on Information Visualization 2000. INFOVIS 2000. Proceedings*. Oct. 2000, pp. 69–75. DOI: 10.1109/INFVIS.2000.885092.
- [18] 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.

- [19] C. Stolte, D. Tang, and P. Hanrahan. “Polaris: A System for Query, Analysis, and Visualization of Multidimensional Relational Databases”. In: *IEEE Transactions on Visualization and Computer Graphics* 8.1 (Jan. 2002), pp. 52–65. ISSN: 1941-0506. DOI: [10.1109/2945.981851](https://doi.org/10.1109/2945.981851).
- [20] Pat Hanrahan. “VizQL: A Language for Query, Analysis and Visualization”. In: *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data*. SIGMOD ’06. New York, NY, USA: Association for Computing Machinery, 2006, p. 721. ISBN: 1-59593-434-0. DOI: [10.1145/1142473.1142560](https://doi.org/10.1145/1142473.1142560).
- [21] J. Mackinlay, P. Hanrahan, and C. Stolte. “Show Me: Automatic Presentation for Visual Analysis”. In: *IEEE Transactions on Visualization and Computer Graphics* 13.6 (Nov. 2007), pp. 1137–1144. ISSN: 1941-0506. DOI: [10.1109/TVCG.2007.70594](https://doi.org/10.1109/TVCG.2007.70594).
- [22] 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.
- [23] Hadley Wickham. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN: 978-3-319-24277-4.
- [24] 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).
- [25] 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).
- [26] 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).
- [27] Jeffrey D. Ullman and Jennifer. Widom. *A First Course in Database Systems*. English. Upper Saddle River, NJ: Pearson Prentice Hall, 2008. ISBN: 0-13-600637-X 978-0-13-600637-4.

- [28] 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).
- [29] Nicholas Sofroniew et al. *Napari/Napari: 0.4.5rc1*. Zenodo. Feb. 2021. DOI: [10.5281/zenodo.4533308](https://doi.org/10.5281/zenodo.4533308).
- [30] *Writing Plugins*. en. [https://imagej.net/Writing\\_plugins](https://imagej.net/Writing_plugins).
- [31] Software Studies. *Culturevis/Imageplot*. Jan. 2021.
- [32] Mathieu Bastian, Sebastien Heymann, and Mathieu Jacomy. “Gephi: An Open Source Software for Exploring and Manipulating Networks”. en. In: *Proceedings of the International AAAI Conference on Web and Social Media* 3.1 (Mar. 2009). ISSN: 2334-0770.
- [33] John Ellson et al. “Graphviz— Open Source Graph Drawing Tools”. In: *Graph Drawing*. Ed. by Petra Mutzel, Michael Jünger, and Sebastian Leipert. Berlin, Heidelberg: Springer Berlin Heidelberg, 2002, pp. 483–484. ISBN: 978-3-540-45848-7.
- [34] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. “Exploring Network Structure, Dynamics, and Function Using NetworkX”. In: *Proceedings of the 7th Python in Science Conference*. Ed. by Gaël Varoquaux, Travis Vaught, and Jarrod Millman. Pasadena, CA USA, 2008, pp. 11–15.
- [35] *Data Representation in Mayavi — Mayavi 4.7.2 Documentation*. <https://docs.entthought.com/mayavi/mayavi/d>
- [36] 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).
- [37] 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).

- [38] 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).
- [39] James Ahrens, Berk Geveci, and Charles Law. “Paraview: An End-User Tool for Large Data Visualization”. In: *The visualization handbook* 717.8 (2005).
- [40] 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).
- [41] Stephen Few and Perceptual Edge. “Dashboard Confusion Revisited”. In: *Perceptual Edge* (2007), pp. 1–6.
- [42] 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).
- [43] 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).
- [44] David I Spivak. *Databases Are Categories*. en. Slides. June 2010.
- [45] David I Spivak. “SIMPLICIAL DATABASES”. en. In: (), p. 35.
- [46] Tamara Munzner. *Visualization Analysis and Design*. AK Peters Visualization Series. CRC press, Oct. 2014. ISBN: 978-1-4665-0891-0.
- [47] Tamara Munzner. “Ch 2: Data Abstraction”. In: *CPSC547: Information Visualization, Fall 2015-2016* ().
- [48] E.H. Spanier. *Algebraic Topology*. McGraw-Hill Series in Higher Mathematics. Springer, 1989. ISBN: 978-0-387-94426-5.
- [49] *Locally Trivial Fibre Bundle - Encyclopedia of Mathematics*. [https://encyclopediaofmath.org/wiki/Locally\\_trivial\\_fibre\\_bundle](https://encyclopediaofmath.org/wiki/Locally_trivial_fibre_bundle)
- [50] “Monoid”. en. In: *Wikipedia* (Jan. 2021).
- [51] “Semigroup Action”. en. In: *Wikipedia* (Jan. 2021).
- [52] nLab authors. “Action”. In: (Mar. 2021).

- 922 [53] Eric W. Weisstein. *Group Homomorphism*. en. <https://mathworld.wolfram.com/GroupHomomorphism.html>.  
923 Text.
- 924 [54] Brendan Fong and David I. Spivak. *An Invitation to Applied Category Theory:  
925 Seven Sketches in Compositionality*. en. First. Cambridge University Press, July  
926 2019. ISBN: 978-1-108-66880-4 978-1-108-48229-5 978-1-108-71182-1. DOI: [10.1017/  
927 9781108668804](https://doi.org/10.1017/9781108668804).
- 928 [55] S. S. Stevens. “On the Theory of Scales of Measurement”. In: *Science* 103.2684 (1946),  
929 pp. 677–680. ISSN: 00368075, 10959203.
- 930 [56] W A Lea. “A Formalization of Measurement Scale Forms”. en. In: (), p. 44.
- 931 [57] Brent A Yorgey. “Monoids: Theme and Variations (Functional Pearl)”. en. In: (),  
932 p. 12.
- 933 [58] “Quotient Space (Topology)”. en. In: *Wikipedia* (Nov. 2020).
- 934 [59] Professor Denis Auroux. “Math 131: Introduction to Topology”. en. In: (), p. 113.
- 935 [60] Robert W. Ghrist. *Elementary Applied Topology*. Vol. 1. Createspace Seattle, 2014.
- 936 [61] Robert Ghrist. “Homological Algebra and Data”. In: *Math. Data* 25 (2018), p. 273.
- 937 [62] David Urbanik. “A Brief Introduction to Schemes and Sheaves”. en. In: (), p. 16.
- 938 [63] Dmitry Nekrasovski et al. “An Evaluation of Pan & Zoom and Rubber Sheet  
939 Navigation with and without an Overview”. In: *Proceedings of the SIGCHI Con-  
940 ference on Human Factors in Computing Systems*. CHI '06. New York, NY, USA:  
941 Association for Computing Machinery, 2006, pp. 11–20. ISBN: 1-59593-372-7. DOI:  
942 [10.1145/1124772.1124775](https://doi.org/10.1145/1124772.1124775).
- 943 [64] Michael S. Crouch, Andrew McGregor, and Daniel Stubbs. “Dynamic Graphs in the  
944 Sliding-Window Model”. In: *European Symposium on Algorithms*. Springer, 2013,  
945 pp. 337–348.
- 946 [65] Chia-Shang James Chu. “Time Series Segmentation: A Sliding Window Approach”.  
947 In: *Information Sciences* 85.1 (July 1995), pp. 147–173. ISSN: 0020-0255. DOI: [10.  
948 1016/0020-0255\(95\)00021-G](https://doi.org/10.1016/0020-0255(95)00021-G).



- [66] Charles R Harris et al. “Array Programming with NumPy”. In: *Nature* 585.7825 (2020), pp. 357–362.
- [67] 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).
- [68] Stephan Hoyer and Joe Hamman. “Xarray: ND Labeled Arrays and Datasets in Python”. In: *Journal of Open Research Software* 5.1 (2017).
- [69] Matthew Rocklin. “Dask: Parallel Computation with Blocked Algorithms and Task Scheduling”. In: *Proceedings of the 14th Python in Science Conference*. Vol. 126. Citeseer, 2015.
- [70] “Retraction (Topology)”. en. In: *Wikipedia* (July 2020).
- [71] Eric W. Weisstein. *Homotopy*. en. <https://mathworld.wolfram.com/Homotopy.html>. Text.
- [72] Tim Bienz, Richard Cohn, and Calif.) Adobe Systems (Mountain View. *Portable Document Format Reference Manual*. Citeseer, 1993.
- [73] 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).
- [74] 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).
- [75] *Cairographics.Org*. <https://www.cairographics.org/>.
- [76] Maxim Shemanarev. *Anti-Grain Geometry*. <https://antigrain.com/>.
- [77] A. M. Cegarra. “Cohomology of Monoids with Operators”. In: *Semigroup Forum*. Vol. 99. Springer, 2019, pp. 67–105. ISBN: 1432-2137.
- [78] Christian Remling. *Algebra (Math 5353/5363 ) Lecture Notes*. Lecture Notes. University of Oklahoma.

- [79] H. M. Johnson. “Pseudo-Mathematics in the Mental and Social Sciences”. In: *The American Journal of Psychology* 48.2 (1936), pp. 342–351. ISSN: 00029556. DOI: [10.2307/1415754](https://doi.org/10.2307/1415754).
- [80] M. A. Thomas. *Mathematization, Not Measurement: A Critique of Stevens’ Scales of Measurement*. en. SSRN Scholarly Paper ID 2412765. Rochester, NY: Social Science Research Network, Oct. 2014. DOI: [10.2139/ssrn.2412765](https://doi.org/10.2139/ssrn.2412765).
- [81] “Connected Space”. en. In: *Wikipedia* (Dec. 2020).
- [82] 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).
- [83] Sheelagh Carpendale. *Visual Representation from Semiology of Graphics by J. Bertin*. en.
- [84] “Jet Bundle”. en. In: *Wikipedia* (Dec. 2020).
- [85] 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).
- [86] 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).
- [87] 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).

- [88] 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).
- [89] 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.
- [90] *Dataclasses — Data Classes — Python 3.9.2rc1 Documentation*. <https://docs.python.org/3/library/dataclasses>.
- [91] 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).
- [92] 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).
- [93] Muhammad Chenariyan Nakhaee. *Mcnakhaee/Palmerpenguins*. Jan. 2021.
- [94] 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.
- [95] 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).
- [96] 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.
- [97] M. Tory and T. Moller. “Rethinking Visualization: A High-Level Taxonomy”. In: *IEEE Symposium on Information Visualization*. 2004, pp. 151–158. DOI: [10.1109/INFVIS.2004.59](https://doi.org/10.1109/INFVIS.2004.59).

- [98] 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.
- [99] Charles D Hansen and Chris R Johnson. *Visualization Handbook*. Elsevier, 2011.
- [100] 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).
- [101] Toussaint Loua. *Atlas Statistique de La Population de Paris*. J. Dejeu & cie, 1873.
- [102] Hadley Wickham and Lisa Stryjewski. “40 Years of Boxplots”. In: *The American Statistician* (2011).
- [103] 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).
- [104] James O Ramsay. *Functional Data Analysis*. Wiley Online Library, 2006.
- [105] Enrico Bertini, Andrada Tatu, and Daniel Keim. “Quality Metrics in High-Dimensional Data Visualization: An Overview and Systematization”. In: *IEEE Transactions on Visualization and Computer Graphics* 17.12 (2011), pp. 2203–2212.
- [106] Bartosz Milewski. “Category Theory for Programmers”. en. In: (), p. 498.