

1

TOPOLOGICAL ARTISTS

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. MARK HANWELL(CHECK SPELLING)

10

JUNE 2021

11

Abstract

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

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

Contents

| | | |
|----|--|----|
| 34 | Abstract | ii |
| 35 | 1 Introduction | 1 |
| 36 | 2 Background | 2 |
| 37 | 2.1 Tools | 3 |
| 38 | 2.2 Data | 4 |
| 39 | 2.3 Visualization | 7 |
| 40 | 2.4 Contribution | 10 |
| 41 | | |
| 42 | 3 Topological Artist Model | 10 |
| 43 | 3.1 Data Space E | 11 |
| 44 | 3.1.1 Variables: Fiber Space F | 12 |
| 45 | 3.1.2 Measurement Scales: Monoid Actions | 14 |
| 46 | 3.1.3 Continuity: Base Space K | 15 |
| 47 | 3.1.4 Data: Sections τ | 17 |
| 48 | 3.2 Graphic: H | 20 |
| 49 | 3.2.1 Idealized Display D | 20 |
| 50 | 3.2.2 Continuity of the Graphic S | 20 |
| 51 | 3.2.3 Rendering ρ | 22 |
| 52 | 3.3 Artist | 23 |
| 53 | 3.3.1 Visual Fiber Bundle V | 24 |
| 54 | 3.3.2 Visual Encoders ν | 25 |
| 55 | 3.3.3 Graphic Assembler Q | 28 |
| 56 | 3.3.4 Assembly Q | 30 |
| 57 | 3.3.5 Assembly factory \hat{Q} | 33 |
| 58 | 3.3.6 Sheaf | 35 |
| 59 | 3.3.7 Composition of Artists: $+$ | 36 |
| 60 | 3.3.8 Equivalence class of artists A' | 37 |

| | | |
|----|---|-----------|
| 61 | 4 Prototype Implementation: Matplottoy | 38 |
| 62 | 4.1 Artist Class A' | 39 |
| 63 | 4.2 Encoders ν | 44 |
| 64 | 4.3 Data E | 45 |
| 65 | 4.4 Case Study: Penguins | 50 |
| 66 | 5 Discussion | 54 |
| 67 | 5.1 Limitations | 54 |
| 68 | 5.2 Future Work | 55 |
| 69 | 6 Conclusion | 57 |

70 **1 Introduction**

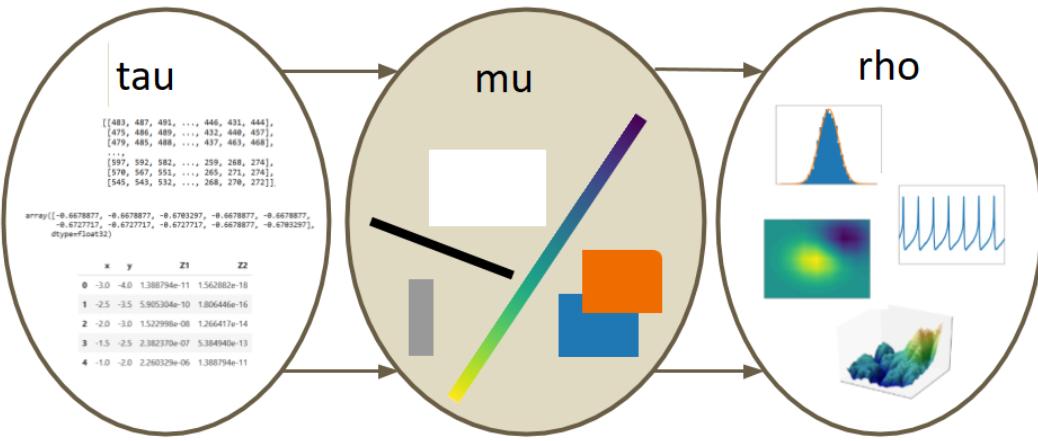


Figure 1: Visualization is equivariant maps between data and visual encoding of the variables and assembly of those encodings into a graphic. **will replace w/ overarching figure w/ same structure**

71 The work presented in this paper is motivated by a need for a library of visualization
 72 components that developers could use to build complex, domain specific tools tuned to
 73 the semantics and structure carried in domain specific data. While many researchers have
 74 identified and described important aspects of visualization, they have specialized in such
 75 different ways as to not provide a model general enough to natively support the full range
 76 of data and visualization types many general purpose modern visualization tools may need
 77 to support. The core architecture also needs to be robust to the big data needs of many
 78 visualization practitioners, and therefore support distributed and streaming data needs. To
 79 support both exploratory and confirmatory visualization[96], this tool needs to support 2D
 80 and 3D, static, dynamic and interactive visualizations.

81 Specifically, this work was driven by a rearchitecture of the Python visualization li-
 82 brary Matplotlib[53] to meet modern data visualization needs. We aim to take advantage
 83 of developments in software design, data structures, and visualization to improve the con-
 84 sistency, compositability, and discoverability of the API. To do so, this work first presents a
 85 mathematical description of how data is transformed into graphic representations, as shown
 86 in figure 1. As with other mathematical formalisms of visualization [56, 62, 92, 98], a

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

90 In other words this model is not intended to be evaluative, it is intended to be a reference
91 specification for visualization library API. To make this model as implementation indepen-
92 dent as possible, we propose fairly general mathematical abstractions of the data container
93 such that we do not need to assume the data has any specific structure, such as a relational
94 database. We reuse this structure for the graphic as that allows us to specifically discuss
95 how structure is preserved. We take a functional approach because functional paradigms
96 encourage writing APIs that are flexible, concise and predictable due to the lack of side
97 effects [61]. Furthermore, by structuring the API in terms of composition of the smallest
98 units of transformation for which we can define correctness, a functional paradigm naturally
99 leads to a library of highly modular components that are composable in such a way that
100 by definition the composition is also correct. This allows us to ensure that domain specific
101 visualizations built on top of these components are also correct without needing knowl-
102 edge of the domain. As with the other mathematical formalisms of visualization, we factor
103 out the rendering into a separate stage; but, our framework describes how these rendering
104 instructions are generated.

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

109 **2 Background**

110 One of the reasons we developed a new formalism rather than adopting the architecture of
111 an existing library is that most information visualization software design patterns, as cate-
112 gorized by Heer and Agrawala[47], are tuned to very specific data structures. This in turn
113 restricts the design space of visual algorithms that display information (the visualization

114 types the library supports) since the algorithms are designed such that the structure of data
115 is assumed, as described in Tory and Möller’s taxonomy [94]. In proposing a new architec-
116 ture, we contrast the trade offs libraries make, describe different types of data continuity,
117 and discuss metrics by which a visualization library is traditionally evaluated.

118 **2.1 Tools**

119 One extensive family of relational table based libraries are those based on Wilkenson’s Gram-
120 mar of Graphics (GoG) [103], including ggplot[101], protovis[14] and D3 [15], vega[82] and
121 altair[97]. The restriction to tables in turn restricts the native design space to visualizations
122 suited to tables. Since the data space and graphic space is very well defined in this gram-
123 mar, it lends itself to a declarative interface [48]. This grammar oriented approach allows
124 users to describe how to compose visual elements into a graphical design [105], while we
125 are proposing a framework for building those elements. An example of this distinction is
126 that the GoG grammar includes computation and aggregation of the table as part of the
127 grammar, while we propose that most computations are specific to domains and only try to
128 describe them when they are specifically part of the visual encoding - for example mapping
129 data to a color. Disentangling the computation from the visual transforms allows us to
130 determine whether the visualization library needs to handle them or if they can be more
131 efficiently computed by the data container.

132 A different class of user facing tools are those that support images, such as ImageJ[83]
133 or Napari[85]. These tools often have some support for visualizing non image components
134 of a complex data set, but mostly in service to the image being visualized. These tools
135 are ill suited for general purpose libraries that need to support data other than images
136 because the architecture is oriented towards building plugins into the existing system [106]
137 where the image is the core data structure. Even the digital humanities oriented ImageJ
138 macro ImagePlot[91], which supports some non-image aggregate reporting charts, is still
139 built around image data as the primary input.

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

ture, as does VTK [39, 45] and its derivatives such as MayaVi[77] and extensions such as ParaView[5] and the infoviz themed Titan[16]. Where GoG and ImageJ type libraries have very consistent APIs for their visualization tools because the data structure is the same, the APIs for visualizations in VTK and Matplotlib are significantly dependent on the structure of the data it expects. This in turn means that every new type of visualization must carry implicit assumptions about data structure in how it interfaces with the input data. This has lead to poor API consistency and brittle code as every visualization type has a very different point of view on how the data is structured. This API choice particularly breaks down when the same dataset is fed into visualizations with different assumptions about structure or into a dashboard consisting of different types of visualization[1, 36] because there is no consistent way to update the data and therefore no consistent way of guaranteeing that the views stay in sync. Our model is a structure dependent formalism, but then also provides a core representation of that structure that is abstract enough to provide a common interface for many different types of visualization.

2.2 Data

Discrete and continuous data and their attributes form a discipline independent design space [73], so one of the drivers of this work was to facilitate building libraries that could natively support domain specific data containers that do not make assumptions about data continuity. As shown in figure 2, there are many types of connectivity. A database typically consists of unconnected records, while an image is an implicit 2D grid and a network is some sort of explicitly connected graph. These data structures typically contain not only the measurements or values of the data, but also domain specific semantic information such as that the data is a map or an image that a modern visualization library could exploit if this information was exposed to the API.

Datasets

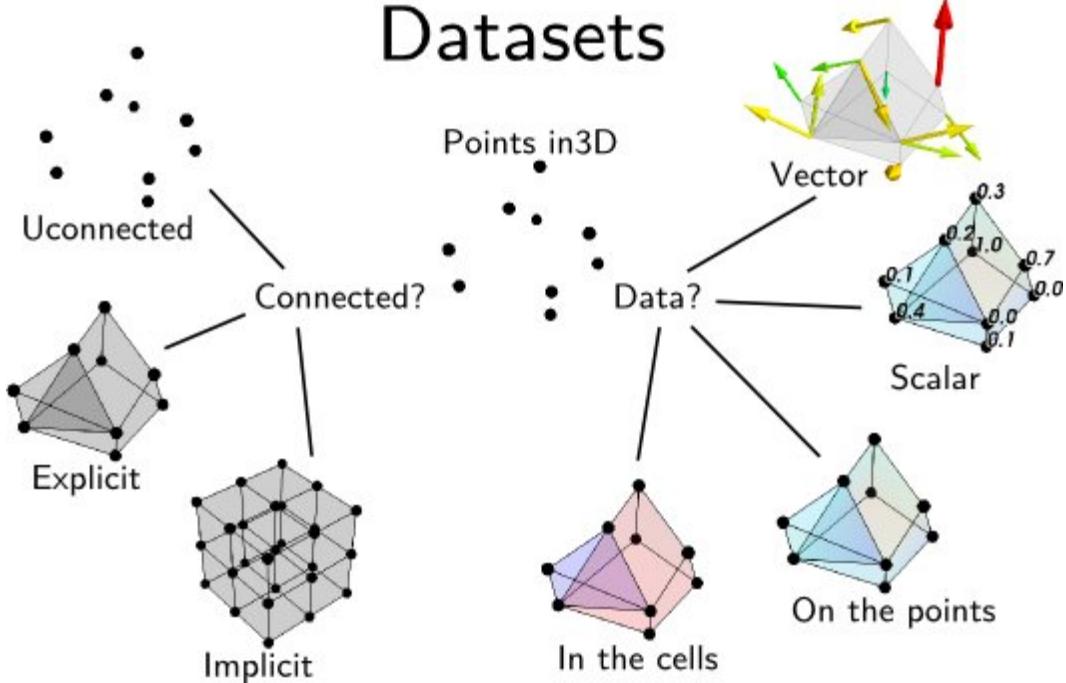


Figure 2: One way to describe data is by the connectivity of the points in the dataset. A database for example is often discrete unconnected points, while an image is an implicitly connected 2D grid. This image is from the Data Representation chapter of the MayaVi 4.7.2 documentation.[32]

166 As shown in figure 2, there are many distinct ways of encoding each specific type of
 167 structure, while as mentioned in section 2.1 APIs are clearer when structured around a
 168 common data representation. Fiber bundles were proposed by Butler as one such represen-
 169 tation because they encode the continuity of the data separately from the types of variables
 170 and are flexible enough to support discrete and ND continuous datasets [18, 19]. Since
 171 Butler’s model lacks a robust way of describing variables, we fold in Spivak’s Simplicial
 172 formulation of databases [87, 88] so that we can encode a schema like description of the
 173 data in the fiber bundle. In this work we will refer to the points of the dataset as *records*
 174 to indicate that a point can be a vector of heterogenous elements. Each *component* of the
 175 record is a single object, such as a temperature measurement, a color value, or an image.
 176 We also generalize *component* to mean all objects in the dataset of a given type, such as

₁₇₇ all temperatures or colors or images. The way in which these records are connected is the
₁₇₈ *connectivity, continuity, or more generally topology.*

definitions

records points, observations, entries

components variables, attributes, fields

connectivity how the records are connected to each other

₁₇₉ Often this topology has metadata associated with it, describing for example when or
₁₈₀ where the measurement was taken. Building on the idea of metadata as *keys* and their
₁₈₁ associated *value* proposed by Munzner [67], we propose that information rich metadata are
₁₈₂ part of the components and instead the values are keyed on coordinate free structural ids.
₁₈₃ In contrast to Munzner's model where the semantic meaning of the key is tightly coupled
₁₈₄ to the position of the value in the dataset, our model allows for renaming all the metadata,
₁₈₅ for example changing the coordinate systems or time resolution, without imposing new
₁₈₆ semantics on the underlying structure.

¹⁸⁷ 2.3 Visualization

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

Figure 3: Retinal variables are a codification of how position, size, shape, color and texture are used to illustrate variations in the components of a visualization. The best to show column describes which types of information can be expressed in the corresponding visual encoding. This tabular form of Bertin's retinal variables is from Understanding Graphics [64] who reproduced it from Krygier and Wood's *Making Maps: A Visual Guide to Map Design for GIS* [57]

¹⁸⁸ Visual representations of data, by definition, reflect something of the underlying structure
¹⁸⁹ and semantics[38], whether through direct mappings from data into visual elements or via
¹⁹⁰ figurative representations that have meaning due to their similarity in shape to external
¹⁹¹ concepts [20]. The components of a visual representation were first codified by Bertin[11].
¹⁹² As illustrated in figure 3, Bertin proposes that there are classes of visual encodings such as
¹⁹³ shape, color, and texture that when mapped to from specific types of measurement, quan-

titative or qualitative, will preserve the properties of that measurement type. For example, that nominal data mapped to hue preserves the selectivity of the nominal measurements. Furthermore he proposes that the visual encodings be composited into graphical marks that match the connectivity of the data - for example discrete data is a point, 1D continuous is the line, and 2D data is the area mark. A general form of marks are glyphs, which are graphical objects that convey one or more attributes of the data entity mapped to it[68, 99] and minimally need to be differentiable from other visual elements [108]. The set of encoding relations from data to visual representation is termed the graphical design by Mackinlay [62, 63] and the design rendered in an idealized abstract space is what throughout this paper we will refer to as a graphic.

The measure of how much of the structure of the data the graphic encodes is a concept Mackinlay termed expressiveness, while the graphic's effectiveness describes how much design choices are made in deference to perceptual saliency [26, 28, 29, 68]. When the properties of the representation match the properties of the data, then the visualization is easier to understand according to Norman's Naturalness Principal[71]. These ideas are combined into Tufte's notion of graphical integrity, which is that a visual representation of quantitative data must be directly proportional to the numerical quantities it represents (Lie Principal), must have the same number of visual dimensions as the data, and should be well labeled and contextualized, and not have any extraneous visual elements [95]. This notion of matching is explicitly formalized by Mackinlay as a structure preserving mapping of a binary operator from one domain to another [63]. A functional dependency framework for evaluating visualizations was proposed by Sugibuchi et al [92], and an algebraic basis for visualization design and evaluation was proposed by Kindlmann and Scheidegger[56]. Vickers et al. propose a category theory framework[98] that extends structural preservation to layout, but is focused strictly on the design layer like the other mathematical frameworks.

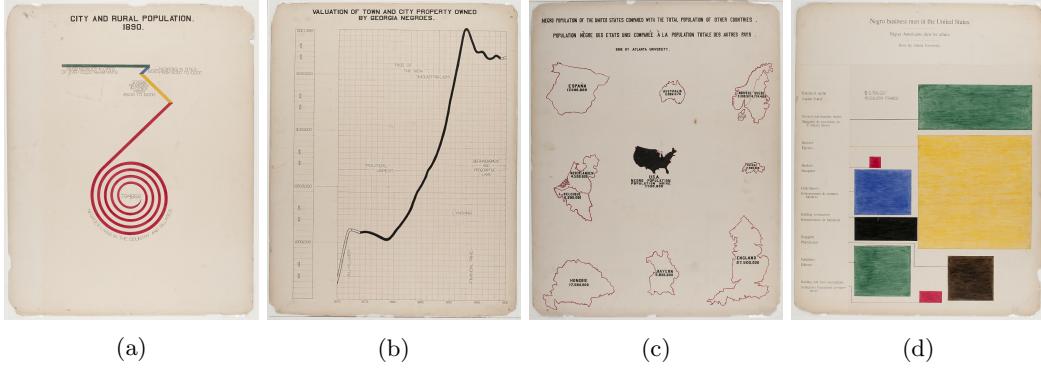


Figure 4: Du Bois' data portraits[35] of post reconstruction Black American life exemplify that the fundamental characteristics of data visualization is that the visual elements vary in proportion to the source data. In figure 4a, the length of each segment maps to population; in figure 4b, the line changes color to indicate a shift in the political environment; in figure 4c the countries are scaled to population size; and figure 4d is a treemap where the area of the rectangle is representative of the number of businesses in each field. The images here are from the Prints and Photographs collection of the Library of Congress [2, 3, 34, 93]

219 One example of highly expressive visualizations are the data portraits by Du Bois shown
 220 in figure 4. While the Du Bois charts are different from the usual scatter, line, and plot
 221 charts, they conform to the constraint that a graphic is a structure preserving map from
 222 data to visual representation. Figure 4a is semantically similar to a bar chart in that the
 223 lengths of the segments are mapped to the values, but in this chart the segments are stacked
 224 together. Figure 4b is a multicolored line chart where the color shifts are at periods of
 225 political significance. In figure 4c, Du Bois combines a graphical representation where glyph
 226 size varies by population with a figurative representation of those glyphs as the countries
 227 the data is from, which means that the semantic and numerical properties of the data are
 228 preserved in the graph. Figure 4b is simply a treemap[49] with space between the marks.
 229 Since the Du Bois data portraits meet the criteria of a faithful visual representation, we
 230 propose a mathematical framework and implementation that allows us to express the Du
 231 Bois charts and common chart types with equal fidelity.

232 **2.4 Contribution**

233 This work presents a mathematical model of the transformation from data to graphic rep-
234 resentation and a proof of concept implementation. Specifically, this work contributes

- 235 1. a functional oriented visualization tool architecture
236 2. topology-preserving maps from data to graphic
237 3. monoidal action equivariant maps from component to visual variable
238 4. algebraic sum such that more complex visualizations can be built from simple ones
239 5. prototype built on Matplotlib’s infrastructure

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

244 **3 Topological Artist Model**

As discussed in the introduction, visualization is generally defined as structure preserving maps from data to graphic representation. In order to formalize this statement, we describe the connectivity of the records using topology and define the structure on the components in terms of the monoid actions on the component types. By formalizing structure in this way, we can evaluate the extent to which a visualization preserves the structure of the data it is representing and build structure preserving visualization tools. We introduce the notion of an artist \mathcal{A} as an equivariant map from data to graphic

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

that carries a homomorphism of monoid actions $\varphi : M \rightarrow M'$ [25], which are discussed in detail in section 3.1.2. Given M on data \mathcal{E} and M' on graphic \mathcal{H} , we propose that artists

\mathcal{A} are symmetric

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

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

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

3.1 Data Space E

Building on Butler’s proposal of using fiber bundles as a common data representation format for visualization data[18, 19], a fiber bundle is a tuple (E, K, π, F) defined by the projection map π

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

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

By definition fiber bundles are locally trivial[59, 86], meaning that over a localized neighborhood we can dispense with extra structure on E and focus on the components and continuity. We use fiber bundles as the data model because they are inclusive enough to express all the types of data described in section 2.2.

²⁶⁴ **3.1.1 Variables: Fiber Space F**

To formalize the structure of the data components, we use notation introduced by Spivak [88] that binds the components of the fiber to variable names and types. 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}_σ

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

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

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

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

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

which is equivalent to

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

²⁶⁵ which allows us to decouple F into components F_i .

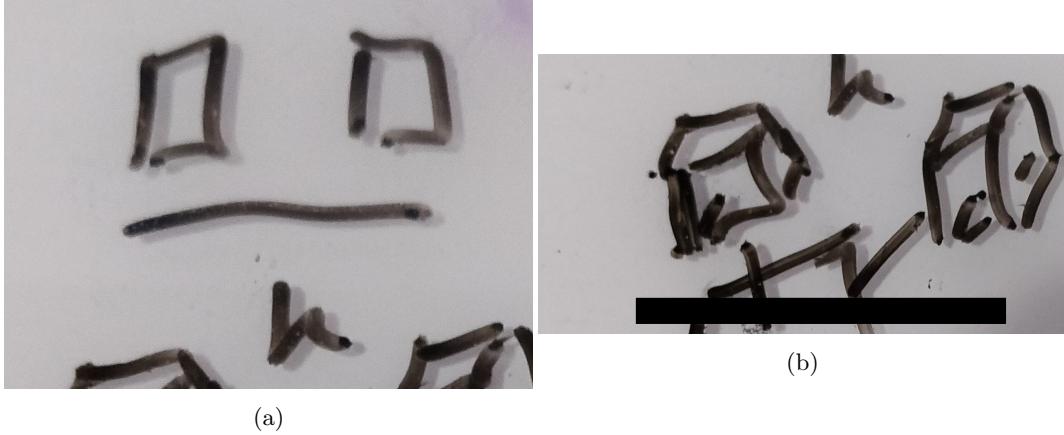


Figure 5: These two datasets have the same base space K but figure 5a has fiber $F = \mathbb{R} \times \mathbb{R}$ which is (time, temperature) while figure 5b has fiber $\mathbb{R}^+ \times \mathbb{R}^2$ which is (time, wind=(speed, direction))

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

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

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

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

²⁶⁶ such that F_1 is specified by ($c = \text{wind}$, $T = \text{wind}$, $\mathbb{U}_\sigma = \mathbb{R}^2$). As illustrated in figure 5,
²⁶⁷ Spivak's framework provides a consistent way to describe potentially complex components
²⁶⁸ of the input data.

269 **3.1.2 Measurement Scales: Monoid Actions**

270 Implementing expressive visual encodings requires formally describing the structure on the
271 components of the fiber, which we define by the action of a monoid on the component.
272 While structure on a set of values is often described algebraically as operations or through the
273 actions of a group, for example Steven’s scales [89], we generalize to monoids to support more
274 component types. Monoids are also commonly found in functional programming because
275 they specify compositions of transformations [90, 107].

A monoid [66] M is a set with an associative binary operator $* : M \times M \rightarrow M$. A monoid has an identity element $e \in M$ such that $e * a = a * e = a$ for all $a \in M$. As defined on a component of F , a left monoid action [4, 84] of M_i is a set F_i with an action $\bullet : M \times F_i \rightarrow F_i$ with the properties:

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

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

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

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

276 of each monoid M_i on F_i . The monoid is also added to the specification of the fiber
277 $(c_i, T_i, \mathbb{U}_\sigma M_i)$

278 Steven’s described the measurement scales[58, 89] in terms of the monoid actions on
279 the measurements: nominal data is permutable, ordinal data is monotonic, interval data is
280 translatable, and ratio data is scalable [100]. For example, given an arbitrary interval scale
281 fiber component ($c = \text{temperature}$, $T = \text{float}$, $\mathbb{U}_\sigma = \mathbb{R}$) with arbitrarily chosen monoid
282 translation actions actions

283 • monoid operator addition $* = +$

284 • monoid operations: $f : x \mapsto x + 1$, $g : x \mapsto x + 2$

285 • monoid action operator composition $\bullet = \circ$

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

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

286 where 1° and 2° are valid distances between two temperatures x . What this diagram means
287 is that either the fiber could be shifted by 1 (vertical line) then by 2 (horizontal), or the
288 two shifts could be combined such that in this case the fiber is shifted by 3 (diagonal) and
289 these two paths yield the same temperature.

290 While many component types will be one of the measurement scale types, we gen-
291 eralize to monoids specifically for the case of partially ordered set. Given a set $W =$
292 $\{mist, drizzle, rain\}$, then the map $f : W \rightarrow W$ defined by

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

294 2. $f(drizzle) = mist$

295 3. $f(mist) = mist$

296 is order preserving such that $mist \leq drizzle \leq rain$ but has no inverse since $drizzle$ and
297 $mist$ go to the same value $mist$. Therefore order preserving maps do not form a group, and
298 instead we generalize to monoids to support partial order component types. Defining the
299 monoid actions on the components serves as the basis for identifying the invariance[56] that
300 must be preserved in the visual representation of the component. We propose equivariance
301 of monoid actions individually on the fiber to visual component maps and on the graphic
302 as a whole.

303 **3.1.3 Continuity: Base Space K**

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

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

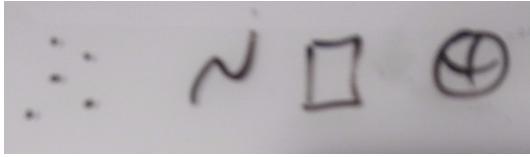


Figure 6: The topological base space K encodes the connectivity of the data space, for example if the data is independent points or a map or on a sphere

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

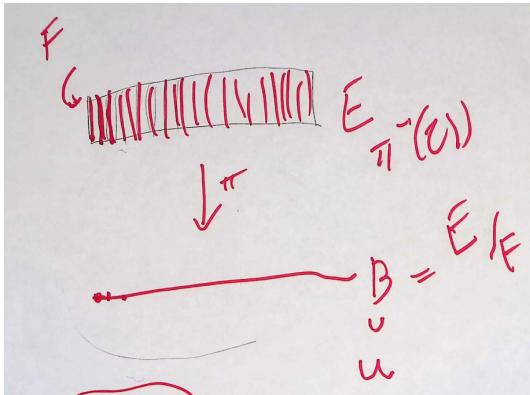


Figure 7: The base space E is divided into fiber segments F . The base space K acts as an index into the records in the fibers. **this figure might be good all the way up top to lay out the components of fb**

Formally K is the quotient space [76] of E meaning it is the finest space[8] such that every $k \in K$ has a corresponding fiber F_k [76]. In figure 7, E is a rectangle divided by vertical fibers F , so the minimal K for which there is always a mapping $\pi : E \rightarrow K$ is the closed interval $[0, 1]$. As with fibers and monoids, we can decompose the total space into components $\pi : E_i \rightarrow K$ where

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

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

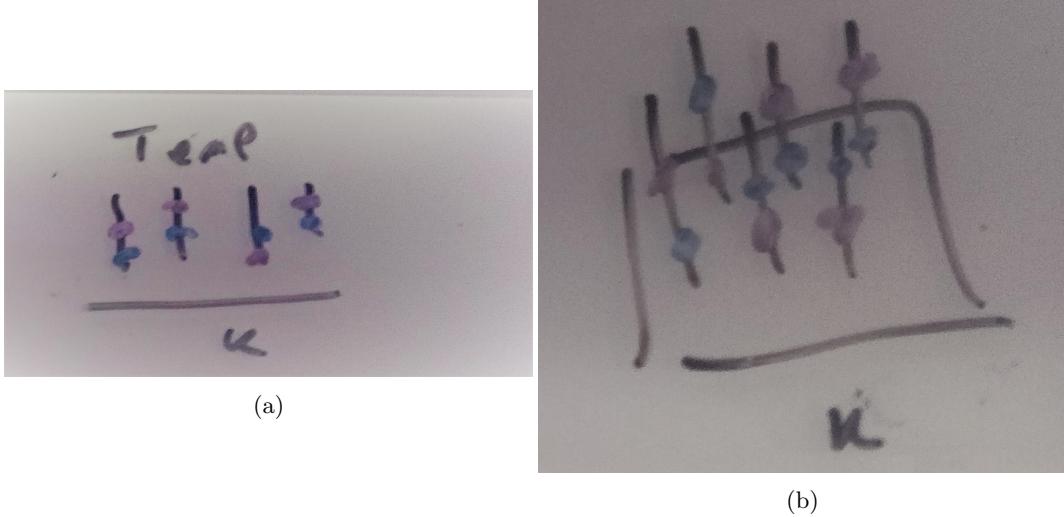


Figure 8: These two datasets have the same (time, temperature) fiber. In figure ?? the total space E is discrete over points $k \in K$, meaning the records in the fiber are also discrete. In figure ?? E lies over the continuous interval K , meaning the records in the fiber are sampled from a continuous space. *revamp figure: F=Plane, k1 = dots, k2=line*

312

313 The datasets in figure 8 have the same fiber of (temperature, time). In figure 8a the
 314 fibers lie over discrete K such that the records in the datasets in the fiber bundles are
 315 discrete. The same fiber in figure 8b lies over a continuous interval K such that the records
 316 are samples from a continuous function defined on K . By encoding this continuity in the
 317 model as K the data model now explicitly carries information about its structure such
 318 that the implicit assumptions of the visualization algorithms are now explicit. The explicit
 319 topology is a concise way of distinguishing visualizations that appear identical, for example
 320 heatmaps and images.

321 **3.1.4 Data: Sections τ**

The section $\tau : K \rightarrow E$ is what ties together the base space K with the fiber F . A section
 is a function that takes as input location $k \in K$ and returns a record $r \in E$. For example,
 in the special case of a table [88], K is a set of row ids, F is the columns, and the section τ

returns the record r at a given key in K . For any fiber bundle, there exists a map

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

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

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

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

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

322 where each section τ_i is a variable or set of variables. This allows for accessing the data
323 component wise rather than just as sections on K .

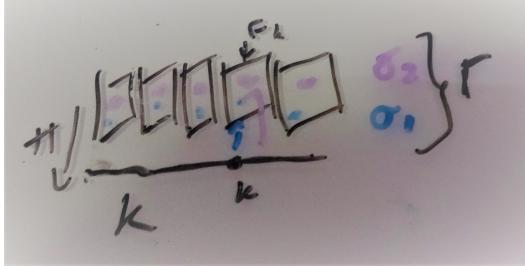


Figure 9: Fiber (time, temperature) with an interval K basespace. The sections τ_i and τ_j are constrained such that the time variable must be monotonic, which means each section is a timeseries of temperature values. They are included in the global set of sections $\tau_1, \tau_2 \in \Gamma(E)$

324 In the example in figure 9, the fiber is $(\text{time}, \text{temperature})$ as described in figure 5 and
325 the base space is the interval K . The section $\tau^{(1)}$ resolves to a series of monotonically
326 increasing in time records of (time, temperature) values. Section $\tau^{(2)}$ returns a different

327 timeseries of (time, temperature) values. Both sections are included in the global set of
328 sections $\tau^{(1)}, \tau^{(2)} \in \Gamma(E)$.

329 This model provides a common interface to widely used data containers without sacri-
330 ficing the semantic structure embedded in each container. For example, the section can be
331 any instance of a univariate numpy array[46] that stores an image. This could be a section
332 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
333 color, and the other two components are the x and y positions of the sampled data in the
334 image. Instead of an image, the numpy array could also store a 2D discrete table. The
335 fiber would not change, but the K would now be 0D discrete points. These different choices
336 in topology indicate, for example, what sorts of interpolation would be appropriate when
337 visualizing the data.

338 There are also many types of labeled containers that can richly be described in this
339 framework because of the fiber. For example, a pandas series which stores a labeled list,
340 or a dataframe[79] which stores a relational table. A series could store the values of $\tau^{(1)}$
341 and a second series could be $\tau^{(2)}$. We could also flatten the fiber to hold two temperature
342 series, such that a section would be an instance of a dataframe with a time column and two
343 temperature columns. While the series and dataframe explicitly have a time index column,
344 they are components in our model and the index is always assumed to be random keys.

345 Where this model particularly shines are N dimensional labeled data structures. For
346 example, an xarray[52] data that stores temperature field could have a K that is a continuous
347 volume and the components would be the temperature and the time, latitude, and longitude
348 the measurements were sampled at. A section can also be an instance of a distributed data
349 container, such as a dask array [81]. As with the other containers, K and F are defined in
350 terms of the index and dtypes of the components of the array. Because our framework is
351 defined in terms of the fiber, continuity, and sections, rather than the exact values of the
352 data, our model does not need to know what the exact values are until the renderer needs
353 to fill in the image.

354 **3.2 Graphic: H**

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

360 **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

$$(x, y, r, g, b) \in D \quad (16)$$

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

365 **3.2.2 Continuity of the Graphic S**

366 Just as the K encodes the connectivity of the records in the data, we propose an equivalent
367 S that encodes the connectivity of the rendered elements of the graphic. For some visual-
368 izations, K may be lower dimension than S . For example, in a typical 2D display (ignoring
369 depth), a point that is 0D in K cannot be represented on screen unless it is thickened to 2D
370 to encode the connectivity of the pixels that visually represent the point. This thickening is
371 often not necessary when the dimensionality of K matches the dimensionality of the target
372 space, for example if K is 2D and the display is a 2D screen. We introduce S to thicken K
373 in a way which preserves the structure of K .



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

Formally, we require that K be a deformation retract[80] of S so that K and S have the same homotopy. The surjective map $\xi : S \rightarrow K$

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

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 τ 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 10. In the case of 1D continuous data and a line plot,
378 the region β over a point α_i specifies the thickness of the line in S for the corresponding
379 τ 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 ξ 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 ξ 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 .

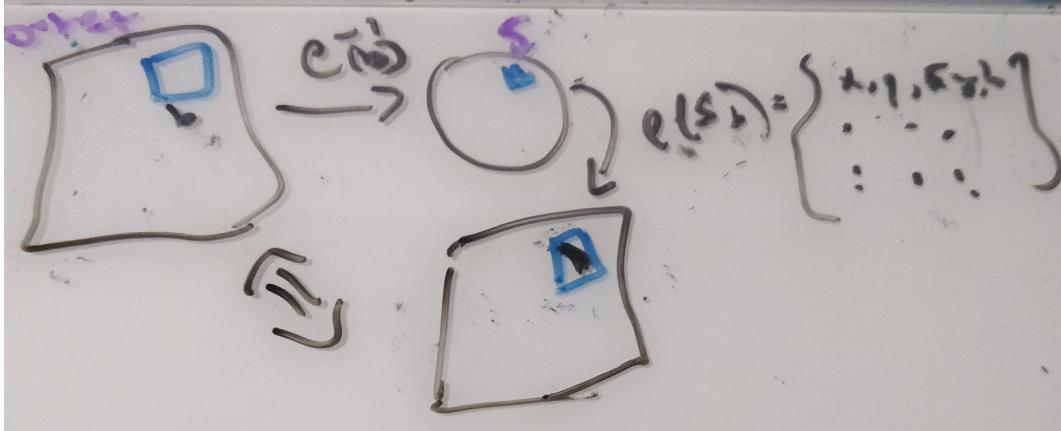


Figure 11: To render a graphic, a pixel p is selected in the display space, which is defined in the same coordinates as the x and y components in D . The inverse mapping $\rho_{xy}^{-1}(p)$ returns a region $S_p \subset S$. $\rho(S_p)$ returns the list of elements $(x, y, r, g, b) \in D$ that lie over S_p . The integral over the (r, g, b) elements is the color of the pixel.

388 3.2.3 Rendering ρ

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

394 For a 2D screen, the pixel is defined as a region $p = [y_{top}, y_{bottom}, x_{right}, x_{left}]$ of the
395 rendered graphic. Since the x and y in p are in the same coordinate system as the x and y
396 components of D the inverse map of the bounding box $S_p = \rho_{xy}^{-1}(p)$ is a region $S_p \subset S$.
397 To compute the color, we integrate on S_p

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

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

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

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

402 In general, ρ is an abstraction of rendering. In very broad strokes ρ can be a specification
 403 such as PDF[13], SVG[75], or an OpenGL scene graph[24]. Alternatively, ρ can be a rendering
 404 engine such as cairo[22] or AGG[7]. Implementation of ρ is out of scope for this work,

405 **3.3 Artist**

We propose that the transformation from data to visual representation can be described as a structure preserving map from one topological space to another. We name this map the artist as that is the analogous part of the Matplotlib[54] architecture that builds visual elements to pass off to the renderer. The topological artist A is a monoid equivariant sheaf map from the sheaf on a data bundle E which is $\mathcal{O}(E)$ to the sheaf on the graphic bundle H , $\mathcal{O}(H)$.

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

406 Sheafs are a mathematical object with restriction maps that define how to glue τ over local
 407 neighborhoods $U \subseteq K$, discussed in section ??, such that the A maps are consistent over
 408 continuous regions of K . While A can usually construct graphical elements solely with the
 409 data in τ , some visualizations, such as line, may also need some finite number n of derivatives,
 410 which is captured by the jet bundle \mathcal{J}^n [55, 69] with $\mathcal{J}^0(E) = E$. In this work, we at most

411 need $\mathcal{J}^2(E)$ which is the value at τ and its first and second derivatives; therefore the artist
412 takes as input the data bundle padded with the jet bundle $E' = E + \mathcal{J}^2(E)$.

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

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

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

420 3.3.1 Visual Fiber Bundle V

421 We introduce a visual bundle V to store the visual representations the artist needs to
422 assemble into a graphic. The visual bundle (V, K, π, P) has section $\mu : V \rightarrow K$ that
423 resolves to a visual variable in the fiber P . The visual bundle V is the latent space of
424 possible parameters of a visualization type, such as a scatter or line plot. We define P
425 in terms of the parameters of a visualization libraries compositing functions; for example
426 table 1 is a sample of the fiber space for Matplotlib [53].

| ν_i | μ_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

427 A section μ is a tuple of visual values that specifies the visual characteristics of a part of
 428 the graphic. For example, given a fiber of $\{xpos, ypos, color\}$ one possible section could be
 429 $\{.5, .5, (255, 20, 147)\}$. The $\text{codomain}(\nu_i)$ determines the monoid actions on P_i . These fiber
 430 components are implicit in the library, by making them explicit as components of the fiber
 431 we can build consistent definitions and expectations of how these parameters behave.

432 **3.3.2 Visual Encoders ν**

As introduced in section 2.3, there are many ways to visually represent data components.

We define the visual transformers ν

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

```

[2]: nu = {'confused': ':(', 'woozy': '=(', 'shruggy': '=@')
[3]: nu.keys()
[3]: dict_keys(['confused', 'woozy', 'shruggy'])
[4]: nu.values()
[4]: dict_values([(':(', '=(', '@=')])
[14]: values
[14]: ['woozy', 'shruggy', 'confused']
[15]: [nu[v] for v in values]
[15]: ['=((', '@=)', ':(']

```

Figure 12: In this artis, ν maps the strings to the emojis. For ν to be equivariant, a shuffle in the words should have an equivalent shuffle in the emojis, and a shuffle in the emojis should have an equivalent shuffle in the words.

⁴³³ 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
⁴³⁴ there is a monoid M'_i on V_i , then there is a monoid homomorphism from $M_i \rightarrow M'_i$ that
⁴³⁵ ν must preserve. As mentioned in section 3.1.2, we choose monoid actions as the basis for
⁴³⁶ equivariance because they define the structure on the fiber components.

A validly constructed ν is one where the diagram of the monoid transform m

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

commutes such that

$$\nu_i(m_r(E_i)) = \varphi(m_v)(\nu_i(E_i)) \tag{25}$$

This equivariance constraint yields guidance on what makes for an invalid transform. For example, given a visual fiber of the same type as the data fiber $F_i = P_i$, the transform $\nu_i(x) = .5$ does not commute under translation $t(x) = x + 2$

$$\nu(t(r + 2)) \stackrel{?}{=} \nu(r) + \nu(2) \tag{26}$$

$$.5 \neq .5 + .5 \tag{27}$$

On the other hand figure 12 illustrates a valid ν mapping from **Strings** to symbols. The group action on these sets is permutation, so shuffling the words must have an equivalent shuffle of the symbols they are mapped to. Continuing with the assumption that $F = P$, we can describe the constraints on the other Steven's measurement group types. To preserve ordinal and partial order monoid actions, ν must be a monotonic function such that given $r_1, r_2 \in E_i$

$$\text{if } delement_1 \leq r_2 \text{ then } \nu(r_1) \leq \nu(r_2) \tag{28}$$

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

$$\nu(x + c) = \nu(x) + \nu(c) \tag{29}$$

while for ratio data, there must be equivalent scaling

$$\nu(xc) = \nu(x)\nu(c) \quad (30)$$

437 The assumption $F_i = P_i$ is made here so that we could omit $\varphi : M \rightarrow M'$ maps that convert
438 the monoid action on F_i to the equivalent action on P_i . As shown in table 1, it is not
439 uncommon for data and visual variables to resolve to the same types. The constraints on ν
440 can be embedded into our artist such that the ν functions are equivariant and also provide
441 guidance on constructing new equivariant ν functions.

442 **3.3.3 Graphic Assembler Q**

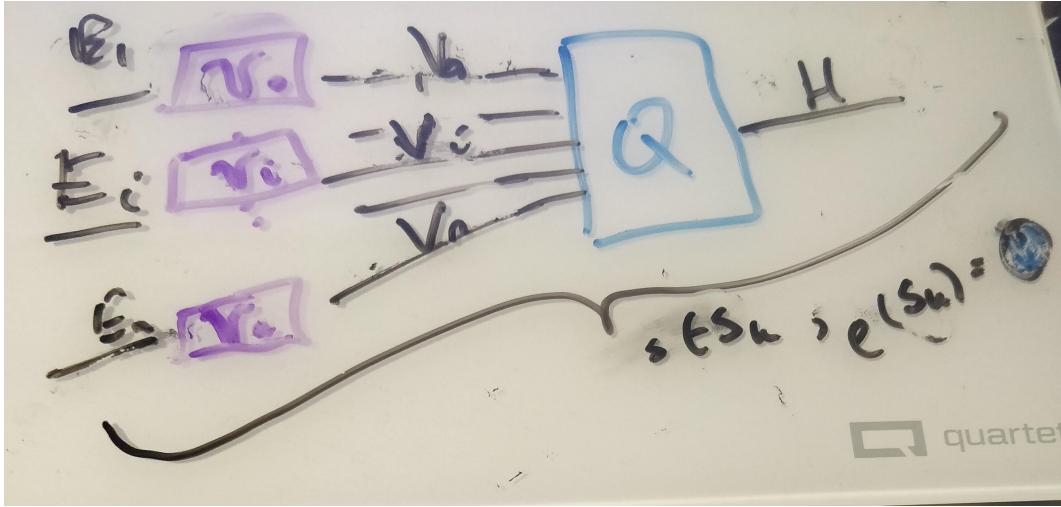


Figure 13: ν functions convert data τ_i to visual characteristics μ_i , then Q assembles μ_i into a graphic ρ such that there is a map ξ preserving the continuity of the data. ρ applied to a region of connected components S_j generates a part of a graphic, for example the point graphical mark.

443 As shown in figure 13, the assembly function Q combines the fiber F_i wise ν transforms into
444 a graphic in H . Together, ν and Q are a map-reduce operation: map the data into their
445 visual encodings, reduce the encodings into a graphic. As with ν the constraint on Q is
446 that for every monoid action on the input μ there is corresponding monoid action on the
447 output ρ .

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

$$J = \{j \in K \text{ s. t. } \exists \gamma \text{ s.t. } \gamma(0) = k \text{ and } \gamma(1) = j\} \quad (31)$$

where the path[30] γ from k to j is a continuous function from the interval [0,1]. We define the glyph as the graphic generated by $Q(S_j)$

$$H \xrightleftharpoons[\rho(S_j)]{} S_j \xrightleftharpoons[\xi^{-1}(J)]{} J_k \quad (32)$$

such that for every glyph there is at least one corresponding section on K . This is in keeping with the definition of glyph as any differentiable element put forth by Ziemkiewicz and Kosara[108]. The primitive point, line, and area marks[11, 23] are specially cased glyphs.

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

Lets call the visual representations of the components $\Gamma(V) = X$ and the graphic $Q(\Gamma(V)) = Y$. If for all monoids $m \in M$ and for all $\mu, \mu' \in X$, the output is equivalent

$$Q(\mu) = Q(\mu') \implies Q(m \circ \mu) = Q(m \circ \mu') \quad (33)$$

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

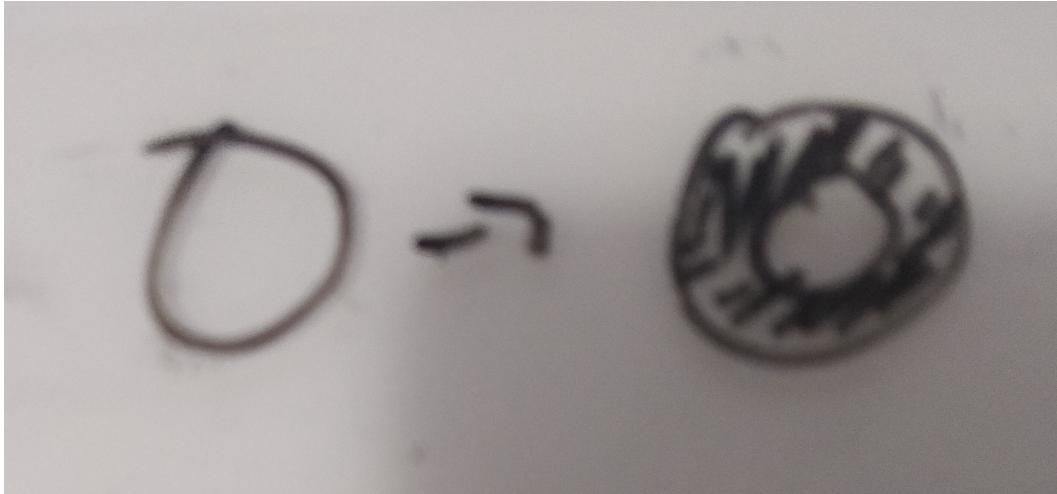


Figure 14: These two glyphs are generated by the same Q function, but differ in the value of the edge thickness parameter μ_i . A valid Q is one where a shift in μ_i is reflected in the glyph generated by ρ .

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

466 **3.3.4 Assembly Q**

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

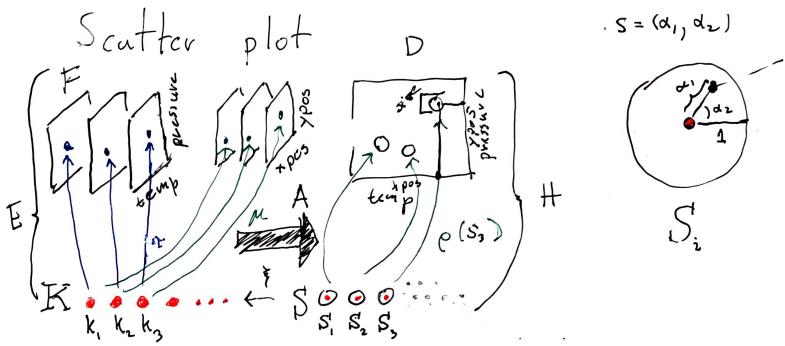


Figure 15: The data is discrete points (temperature, time). Via ν these are converted to (xpos, ypos) and pulled over discrete S . These values are then used to parameterize ρ which returns a color based on the parameters (xpos,ypos) and position α, β on S_k that ρ is evaluated on.

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

$$x = size \bullet \alpha \bullet \cos(\beta) + xpos \quad (34)$$

$$y = size \bullet \alpha \bullet \sin(\beta) + ypos \quad (35)$$

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

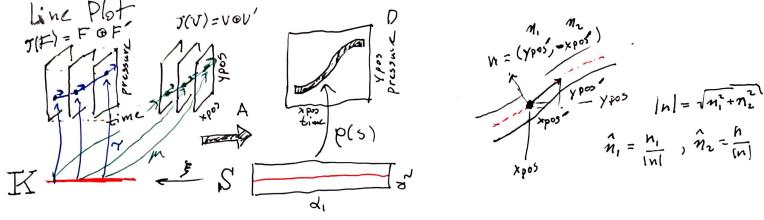


Figure 16: The line fiber (*time, temp*) is thickened with the derivative (*time', temperature'*) because that information will be necessary to figure out the tangent to the point to draw a thick line. This is because the line needs to be pushed perpendicular to the tangent of (*xpos, ypos*). *this is gonna move once this gets regenerated w/ labels* The data is converted to visual characteristics (*xpos, ypos*). The α coordinates on S specifies the position of the line, the β coordinate specifies thickness.

The line plot $Q(xpos, \hat{n}_1, ypos, \hat{n}_2)(\alpha, \beta)$ shown in fig 15 exemplifies the need for the jet. The line needs to know the tangent of the data to draw an envelope above and below each (*xpos,ypos*) such that the line appears to have a thickness. The magnitude of the thickness is

$$|n| = \sqrt{n_1^2 + n_2^2} \quad (36)$$

such that the normal is

$$\hat{n}_1 = \frac{n_1}{|n|}, \quad \hat{n}_2 = \frac{n_2}{|n|} \quad (37)$$

which yields components of ρ

$$x = xpos(\xi(\alpha)) + \beta \hat{n}_1(\xi(\alpha)) \quad (38)$$

$$y = ypos(\xi(\alpha)) + \beta \hat{n}_2(\xi(\alpha)) \quad (39)$$

470 where (x,y) look up the position $\xi(\alpha)$ on the data. At that point, we also look up the
471 the derivatives \hat{n}_1, \hat{n}_2 which are then multiplied by data to specify the thickness.

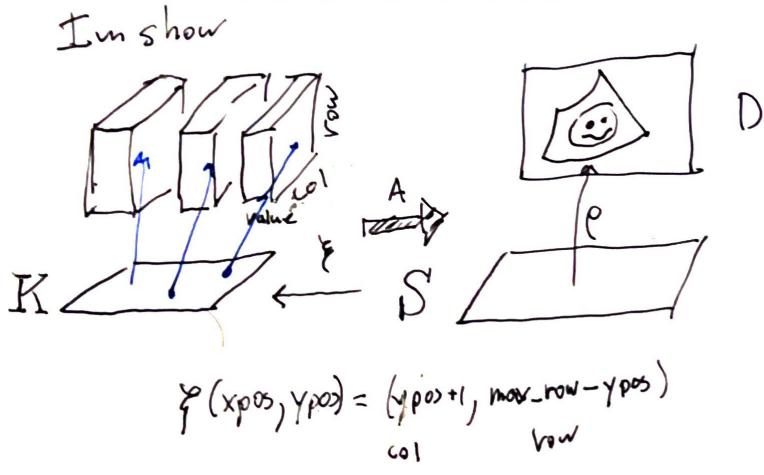


Figure 17: The only visual parameter a image requires is color since ξ encodes the mapping between position in data and position in graphic.

⁴⁷² The image $Q(\text{color})$ in figure 17 is a direct lookup $\xi : S \rightarrow K$ such that

$$R = R(\xi(\alpha, \beta)) \quad (40)$$

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

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

⁴⁷³ where ξ may do some translating to a convention expected by Q for example reorientng the
⁴⁷⁴ array such that the first row in the data is at the bottom of the graphic.

⁴⁷⁵ **3.3.5 Assembly factory \hat{Q}**

⁴⁷⁶ The graphic base space S is not accessible in many architectures, including Matplotlib;
⁴⁷⁷ instead we can construct a factory function \hat{Q} over K that can build a Q . As shown in
⁴⁷⁸ eq 22, Q is a bundle map $Q : \xi^*V \rightarrow H$ where ξ^*V and H are both bundles over S .

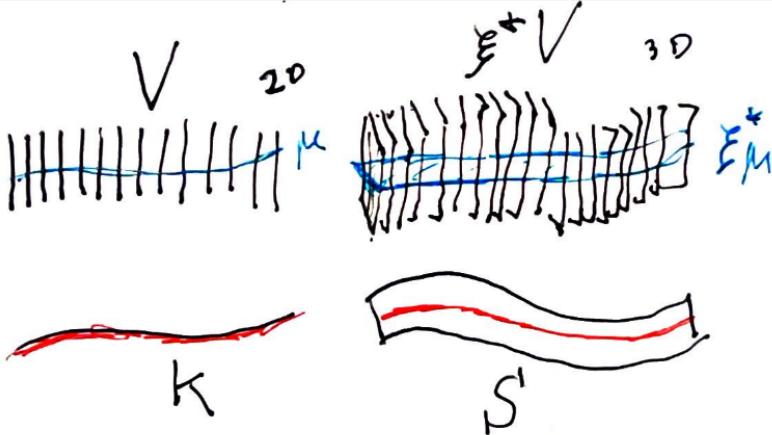


Figure 18: The pullback of the visual bundle ξ^*V is the replication of a μ over all points s that map back to a single k . Because the μ is the same, we can construct a \hat{Q} on μ over k that will fabricate the Q for the equivalent set of s associated to that k

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

$$\xi^*V|_{\xi^{-1}(k)} = \xi^{-1}(k) \times P \quad (43)$$

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

481 This copying is illustrated in figure 18, where the 1D fiber $P \hookrightarrow V$ over K is copied
 482 repeatedly to become the 2D fiber $P^*\mu \hookrightarrow \xi^*V$ with identical components over S . Given
 483 the section $\xi^*\mu$ pulled back from μ and the point $s \in \xi^{-1}(k)$, there is a direct map from μ on
 484 a point k , there is a direct map from the visual section over data base space $(k, \mu(k)) \mapsto$
 485 $(s, \xi^*\mu(s))$ to the visual section $\xi^*\mu$ over graphic base space. This map means that the
 486 pulled back section $\xi^*\mu(s) = \xi^*(\mu(k))$ is the section μ copied over all s . This means that
 487 $\xi^*\mu$ is identical for all s where $\xi(s) = k$, which is illustrated in figure 18 as each dot on P is
 488 equivalent to the line intersection $P^*\mu$.

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

to the pulled back visual section $\xi^*\mu$ then

$$\rho(s) := Q(\xi^*\mu)(s) \quad (44)$$

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

$$\hat{Q}(\mu(k))(s) := Q((\xi^*\mu)(s)) \quad (45)$$

where the assembly function \hat{Q} that takes as input the visual section on data μ is defined to be the assembly function 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 $\hat{Q}(\mu(k)) = Q(\xi^*\mu)$ generates a curried Q . 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 $\hat{Q} : \Gamma(V_k) \rightarrow \Gamma(H|_{\xi^{-1}(k)})$. This allows us to construct a \hat{Q} that only depends on K , such that for each $\mu(k)$ there is part of $\rho|_{\xi^{-1}(k)}$. The construction of \hat{Q} allows us to retain the functional map reduce benefits of Q without having to majorly restructure the existing rendering pipeline.

3.3.6 Sheaf

The restriction maps of a sheaf describe how local τ can be glued into larger sections [40, 41]. As part of the definition of local triviality, there is an open neighborhood $U \subset K$ for every $k \in K$. We can define the inclusion map $\iota : U \rightarrow K$ which pulls E over U

$$\begin{array}{ccc} \iota^* E & \xhookrightarrow{\iota^*} & E \\ \pi \downarrow \lrcorner \iota^* \tau & & \pi \downarrow \lrcorner \tau \\ U & \xhookrightarrow{\iota} & K \end{array} \quad (46)$$

such that the pulled back $\iota^*\tau$ only contains records over $U \subset K$. By gluing $\iota^*\tau$ together, the sheaf is putting a continuous structure on local sections which allows for defining a section

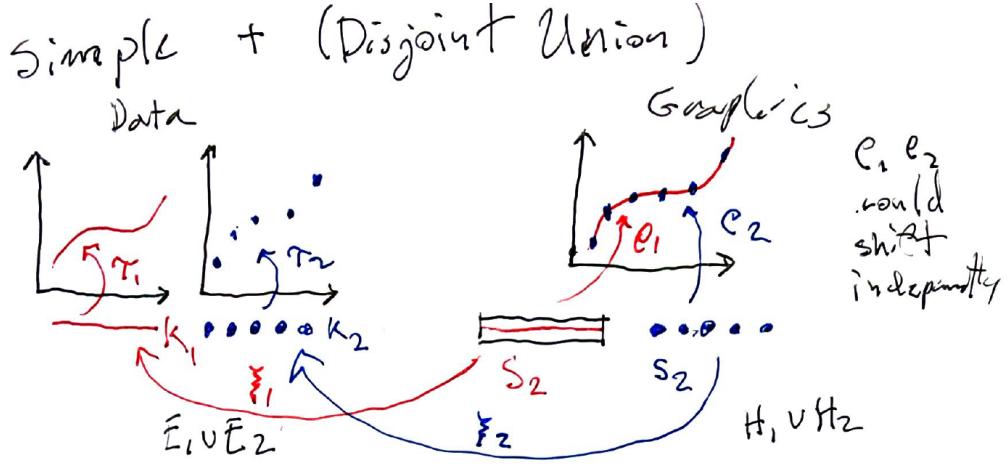


Figure 19: τ_1 and τ_2 are distinct datasets passed through artists A_1 and A_2 to generate graphics ρ_1 and ρ_2 . These graphics happen to be rendered to the same image, but otherwise have no intrinsic link.

501 over a subset in K . That section over subset K maps to the graphic generated by A for
 502 visualizations such as sliding windows[27, 31] streaming data, or navigation techniques such
 503 as pan and zoom[72].

504 **3.3.7 Composition of Artists: +**

505 To build graphics that are the composites of multiple artists, we define a simple addition
 506 operator that is the disjoint union of fiber bundles E . For example, in figure 19 the scatter
 507 plot E_1 and the line plot E_2 have different K that are mapped to separate S . To fully
 508 display both graphics, the composite graphic $A_1 + A_2$ needs to include all records on both
 509 K_1 and K_2 , which are the sections on the disjoint union $K_1 \sqcup K_2$. This in turn yields disjoint
 510 graphics $S_1 \sqcup S_2$ rendered to the same image. Constraints can be placed on the disjoint
 511 union such as that the fiber components need to have the same v position encodings or that
 512 the position μ need to be in a specified range. There is a second type of composition where
 513 E_1 and E_2 share a base space $K_2 \hookrightarrow K_1$ such that the the artists can be considered to be
 514 acting on different components of the same section. This type of composition is important
 515 for creating visualizations where elements need to update together in a consistent way, such
 516 as multiple views [6, 74] and brush-linked views[10, 17].

517 **3.3.8 Equivalence class of artists A'**

518 It is impractical to implement an artist for every single graphic; instead we implement an
519 approximation of an the equivalence class of artists $\{A \in A' : A_1 \equiv A_2\}$. Roughly, equivalent
520 artists have the same fiber bundle V and same assembly function Q but act on different
521 sections μ , but we will formalize the definition of the equivalence class in future work.

522 As a first pass for implementation purposes, we identify a minimal P associated with
523 each A' that defines what visual characteristics of the graphic must originate in the data
524 such that the graphic is identifiable as a given chart type.

Figure 20: Each of these graphics is generated by a different artist A which is the equivalence
class of scatter plots A'
this is gonna be a whole bunch of scatter plots

525 For example, a scatter plot of red circles is the output of one artist, a scatter plot of
526 green squares the output of another, as shown in figure 20. These two artists are equivalent
527 since their only difference is in the literal visual encodings (color, shape). Shape and color
528 could also be defined in Q but the position must come from the fiber $P = (xpos, ypos)$ since
529 fundamentally a scatter plot is the plotting of one position against another[38]. We also use
530 this criteria to identify derivative types, for example the bubble chart[95] is a type of scatter
531 where by definition the glyph size is mapped from the data. The criteria for equivalence
532 class membership serves as the basis for evaluating invariance[56].

533 4 Prototype Implementation: Matplottoy

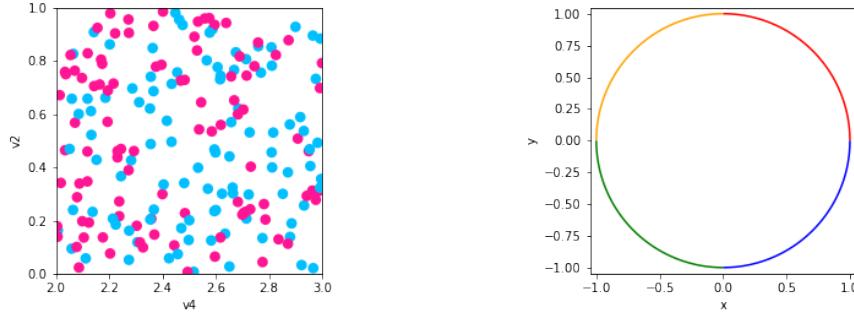


Figure 21: Scatter plot and line plot implemented using prototype artists and data models, building on Matplotlib rendering.

534 To prototype our model, we implemented the artist classes for the scatter and line plots
 535 shown in figure 21 because they differ in every attribute: different visual channels ν that
 536 composite to different marks Q with different continuities ξ . We make use of the Matplotlib
 537 figure and axes artists [53, 54] so that we can initially focus on the data to graphic trans-
 538 formations. To generate the images in figure 21, we instantiate `fig`, `ax` artists that will
 539 contain the new `Point`, `Line` primitive objects we implemented based on our topology
 540 model.

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

541 We then add the `Point` and `Line` artist that construct the scatter and line graphics.
 542 These artists are implemented as the equivalence class A' with the aesthetic configurations
 543 factored out into a `transforms` dictionary that specifies the visual bundle V . The equivalence
 544 classes A' map well to Python classes since the functional aspects- ν , \hat{Q} , and ξ - are completely
 545 reusable in a consistent composition, while the visual values in V are what change between
 546 different artists belonging to the same class A' . The `data` object is an abstraction of a
 547 data bundle E with a specified section τ . Implementing H and ρ are out of scope for this
 548 prototype because they are part of the rendering process. We also did not implement any
 549 form of ξ because the scatter, line, and bar plots prototyped here directly broadcast from k
 550 to s , unlike for example an image which may need to be rotated.

551 4.1 Artist Class A'

552 The artist is the piece of the matplotlib architecture that constructs an internal representa-
 553 tion of the graphic that the render then uses to draw the graphic. In the prototype artist,
 554 `transform` is a dictionary of the form `{parameter:(variable, encoder)}` where param-
 555 eter is a component in P , variable is a component in F , and the ν encoders are passed in as
 556 functions or callable objects. The data bundle E is passed in as a `data` object. By binding
 557 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)

```

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

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

577 The `Point` artist builds on `collection` artists because collections are optimized to ef-
 578 ficiently draw a sequence of primitive point and area marks. In this prototype, the scatter
 579 marker shape is fixed as a circle, and the only visual fiber components are x and y position,
 580 size, and the facecolor of the marker. We only show the `assemble` function here because
 581 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)

```

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

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

593 The main difference between the `Point` and `Line` objects is in the `assemble` function
 594 because line has different continuity from scatter and is represented by a different type of
 595 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)

```

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

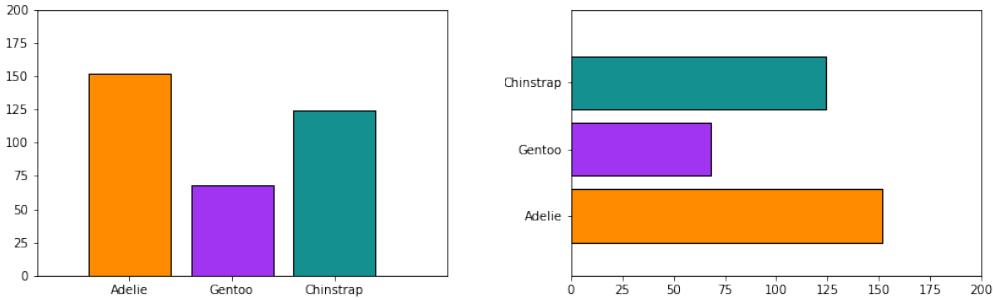


Figure 22: Frequency of Penguin types visualized as discrete bars.

603 The bar charts in figure 22 are generated with a `Bar` artist. The have the same required
 604 P components of (position, length), but in `__init__` of `Bar` an additional parameter is set,
 605 `orientation` which controls whether the bars are arranged vertically or horizontally. This
 606 parameter only applies holistically to the graphic and never to individual data parameters,
 607 and highlights how the model encourages explicit differentiation between parameters in V
 608 and graphic parameters 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, facecolors='C0', edgecolors='k', offset=0):
14         #set some defaults
15         width = itertools.repeat(width) if np.isscalar(width) else width
16         floor = itertools.repeat(floor) if np.isscalar(floor) else (floor, )
17
18         # offset is passed through via assemblers such as multigroup, not supposed to be directly
19         position = position + offset
20
21     def make_bars(xval, xoff, yval, yoff):
22         return [[(x, y), (x, y+y0), (x+x0, y+y0), (x+x0, y), (x, y)]
23                 for (x, xo, y, yo) in zip(xval, xoff, yval, yoff)]
24         #build bar glyphs based on graphic parameter
25     if self.orientation in {'vertical', 'v'}:
26         verts = make_bars(position, width, floor, length)
27     elif self.orientation in {'horizontal', 'h'}:
28         verts = make_bars(floor, length, position, width)
29
30     self._paths = [mpath.Path(xy, closed=True) for xy in verts]
31     self.set_edgecolors(edgecolors)
32     self.set_facecolors(facecolors)
33
34     def draw(self, renderer, *args, **kwargs):
35         view = self.data.view(self.axes)
36         visual = {}
37         for (p, t) in self.transforms.items():
38             if isinstance(t, dict):
39                 if t['name'] in self.data.FB.F and 'encoder' in t:
40                     visual[p] = t['encoder'](view[t['name']])
41                 elif 'encoder' in t: # constant value
42                     visual[p] = t['encoder'](t['name'])
43                 elif t['name'] in self.data.FB.F: # identity
44                     visual[p] = view[t['name']]
45             else: # no transform
46                 visual[p] = t
47         self.assemble(**visual)
48         super().draw(renderer, *args, **kwargs)

```

609 The `draw` method here has a more complex unpacking of visual encodings to support passing
 610 in visual component data directly. This is vastly simplifies building composite objects as
 611 the alternative would be higher order functions that take as input the transforms passed in
 612 by the user. This construction supports a constant visual parameter, an identity transform
 613 where the value is the same in E and V , and setting the visual component directly. The
 614 `assemble` function constructs bars and sets their face and edge colors. The `make_bars`
 615 function converts the input position and length to the coordinates of a rectangle of the given
 616 width. Defaults are provided for 'width' and 'floor' to make this function more reusable.
 617 Typically the defaults are used for the type of chart shown in figure 22, but these visual
 618 variables are often set when building composite versions of this chart type as discussed in
 619 section 4.4.

620 4.2 Encoders ν

621 As mentioned above, the encoding dictionary is specified by the visual fiber component, the
 622 corresponding data fiber component, and the mapping function. The visual parameter serves
 623 as the dictionary key because the visual representation is constructed from the encoding
 624 applied to the data $\mu = \nu \circ \tau$. For the scatter plot, the mappings for the visual fiber
 625 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)}}

```

626 where the position (x, y) ν transformers are identity functions. The size s transformer is not
 627 acting on a component of F , instead it is a ν that returns a constant value. While size could
 628 be embedded inside the `assemble` function, it is added to the transformers to illustrate user
 629 configured visual parameters that could either be constant or mapped to a component in F .
 630 The identity and constant ν are explicitly implemented here to demonstrate their implicit

631 role in the visual pipeline, but they are somewhat optimized away in `Bar`. More complex
632 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]
```

633 where `__init__` can validate that the output of the ν is a valid element of the P com-
634 ponent the ν function is targeting. Creating a callable class also provides a simple way to
635 swap out the specific (data, value) mapping without having to reimplement the validation
636 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)
```

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

640 4.3 Data E

641 The data input into the `Artist` will often be a wrapper class around an existing data
642 structure. This wrapper object must specify the fiber components F and connectivity K
643 and have a `view` method that returns an atomic object that encapsulates τ . The object
644 returned by the view must be key valued pairs of `{component name : component section}`

645 where each section is a component as defined in equation 15. To support specifying the fiber
 646 bundle, we define a `FiberBundle` data class[33]

```

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

```

647 that asks the user to specify how K is triangulated and the attributes of F . Python
 648 dataclasses are a good abstraction for the fiber bundle class because the `FiberBundle` class
 649 only stores data. The K is specified as tables because the `assemble` functions expect
 650 tables that match the continuity of the graphic; scatter expects a vertex table because it
 651 is discontinuous, line expects an edge table because it is 1D continuous. The fiber informs
 652 appropriate choice of ν therefore it is a dictionary of attributes of the fiber components.

653 To generate the scatter plot in figure 21, we fully specify a dataset with random keys
 654 and values in a section chosen at random from the corresponding fiber component. The
 655 fiberbundle `FB` is a class level attribute since all instances of `VertexSimplex` come from the
 656 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': {'type': float, 'monoid': 'interval', 'range': [0,1]},
6                          'v2': {'type': str, 'monoid': 'nominal', 'range': ['true', 'false']},
7                          'v3': {'type': float, 'monoid': 'interval', 'range': [2,3]}})
8
9     def __init__(self, sid = 45, size=1000, max_key=10**10):
10         # create random list of keys

```

```

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

```

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

```

1  class EdgeSimplex:
2      # assign a class level FB attribute
3      def __init__(self, num_edges=4, num_samples=1000):
4          self.keys = range(num_edge) #edge id
5          # distance along edge
6          self.distances = np.linspace(0,1, num_samples)
7          # half generalized representation of arcs on a circle
8          self.angle_samples = np.linspace(0, 2*np.pi, len(self.keys)+1)
9
10     @staticmethod
11     def _color(edge):
12         colors = ['red', 'orange', 'green', 'blue']
13         return colors[edge%len(colors)]
14
15     @staticmethod
16     def _xy(edge, distances, start=0, end=2*np.pi):
17         # start and end are parameterizations b/c really there is
18         angles = (distances *(end-start)) + start
19         return np.cos(angles), np.sin(angles)

```

```

20
21     def tau(self, k): #will fix location on page on revision
22         x, y = self._xy(k, self.distances,
23                            self.angle_samples[k], self.angle_samples[k+1])
24         color = self._color(k)
25         return (k, (x, y, color))
26
27     def view(self, axes):
28         table = defaultdict(list)
29         for k in self.keys():
30             table['index'].append(k)
31             # (name, value) pair, value is [x0, ..., xn] for x, y
32             for (name, value) in zip(self.FB.fiber.keys(), self.tau(k, simplex)[1]):
33                 table[name].append(value)

```

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

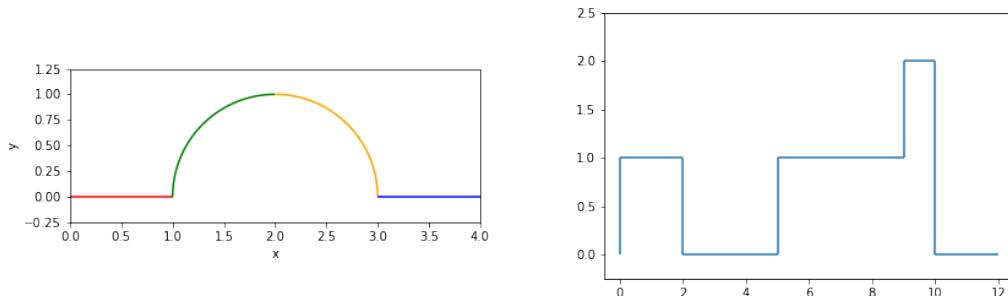


Figure 23: Continuous and discontinuous lines as defined by different data models, but generated with the same `A'Line`

668 The graphics in figure 23 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          # 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
23         return table

```

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

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

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

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

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

677 **4.4 Case Study: Penguins**

678 For this case study, we use the Palmer Penguins dataset[42, 51] since it is multivariate and
679 has a varying number of penguins. We use a version of the data packaged as a pandas
680 dataframe[70] since that is a very commonly used Python labeled data structure. The
681 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
```

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

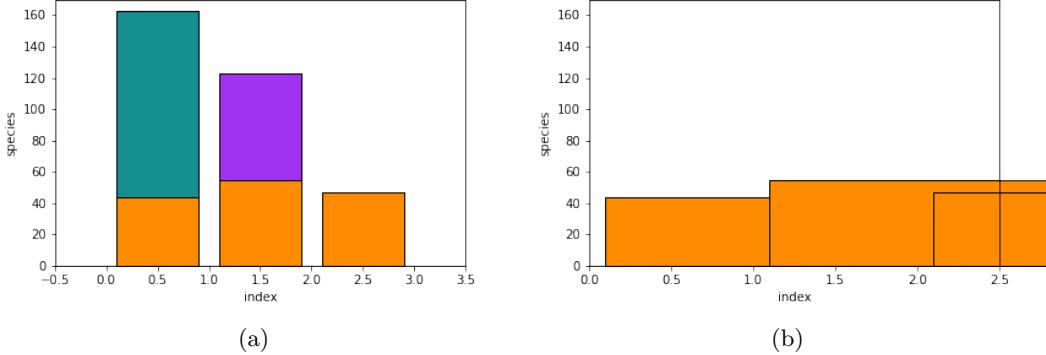


Figure 24: Penguin count disaggregated by island and species

687 The stacked and grouped bar charts in figure 24 are both out of `Bar` artists such that
 688 the difference between `StackedBar` and `GroupedBar` is specific to the ways in which the
 689 `Bar` are stitched together. These two artists have identical constructors and `draw` methods.
 690 As with `Bar`, the orientation is set in the constructor. In both these artists, we separate
 691 the transforms applied to only one component and the case `mtransforms` where the same
 692 transform is applied to multiple components such that V has multiple components that map
 693 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
12          super().__init__(*args, **kwargs)
13          self.data = data
14          self.orientation = orientation
15          self.transforms = copy.deepcopy(transforms)
16          self.mtransforms = copy.deepcopy(mtransforms)
```

```

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)

```

694 Since all the visual transformation is passed through to `Bar`, the `draw` method does not
695 do any visual transformations. In `StackedBar` the `view` is used to adjust the `floor` for
696 every subsequent bar chart since a stacked bar chart is bar chart area marks concatenated
697 together in the `length` parameter. In contrast, `GroupedBar` does not even need the `view`, but
698 instead keeps track of the relative position of each group of bars in the visual only variable
699 `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)
```

700 Since the only difference between these two glyphs is in the composition of `Bar`, they take
701 in the exact same transform specification dictionaries. The `transform` dictionary dictates
702 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}}]
```

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

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

707 but differ in assembly \hat{Q} . By decomposing the architecture into data, visual encoding,
708 and assembly steps, we are able to build components that are more flexible and also more
709 self contained than the existing code base. While very rough, this API demonstrates that
710 the ideas presented in the math framework are implementable. In choosing a functional
711 approach, if not implementation, we provide a framework for library developers to build
712 reusable encoder ν assembly \hat{Q} and artists A . We argue that if these functions are built
713 such that they are equivariant with respect to monoid actions and the graphic topology is a
714 deformation retraction of the data topology, then the artist by definition will be a structure
715 and property preserving map from data to graphic.

716 **5 Discussion**

717 This work contributes a mathematical description of the mapping A from data to visual
718 representation. Combining Butler’s proposal of a fiber bundle model of visualization data
719 with Spivak’s formalism of schema lets this mode; support a variety of datasets, includ-
720 ing discrete relational tables,, multivariate high resolution spatio temporal datasets, and
721 complex networks. Decomposing the artist into encoding ν , assembly Q , and reindexing ξ
722 provides the specifications for producing visualization where the structure and properties
723 match those of the input data. These specifications are that the graphic must have continu-
724 ity equivalent to the data, and that the visual characteristics of the graphics are equivariant
725 to their corresponding components under monoid actions. This model defines these con-
726 straints on the transformation function such that they are not specific to any one type of
727 encoding or visual characteristic. Encoding the graphic space as a fiber bundle provides a
728 structure rich abstraction of the target graphical design in the target display space.

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

740 **5.1 Limitations**

741 So far this model has only been worked out for a single data set tied to a primitive mark,
742 but it should be extensible to compositing datasets and complex glyphs. The examples and

743 prototype have so far only been implemented for the static 2D case, but nothing in the math
744 limits to 2D and expansion to the animated case should be possible because the model is
745 formalized in terms of the sheaf. While this model supports equivariance of figurative glyphs
746 generated from parameters of the data[9, 21], it currently does not have a way to evaluate
747 the semantic accuracy of the figurative representation. Effectiveness is out of scope for this
748 model because it is not part of the structure being preserved, but potentially a developer
749 building a domain specific library with this model could implement effectiveness criteria in
750 the artists. Also, even though the model is designed to be backend independent, it has only
751 really been tested against the AGG backend. It is especially unknown how this framework
752 interfaces with high performance rendering libraries such as OpenGL[24]. Because this model
753 has been limited to the graphic design space, it does not address the critical task of laying
754 out the graphics in the image

755 This model and the associated prototype is deeply tied to Matplotlib’s existing archi-
756 tecture. While the model is expected to generalize to other libraries, such as those built on
757 Mackinlay’s APT framework, this has not been worked through. In particular, Mackinlay’s
758 formulation of graphics as a language with semantic and syntax lends itself a declarative in-
759 terface[61], which Heer and Bostock use to develop a domain specific visualization language
760 that they argue makes it simpler for designers to construct graphics without sacrificing
761 expressivity [48]. Similarly, the model presented in this work formulates visualization as
762 equivariant maps from data space to visual space, and is designed such that developers can
763 build software libraries with data and graphic topologies tuned to specific domains.

764 5.2 Future Work

765 While the model and prototype demonstrate that generation of simple marks from the data,
766 there is a lot of work left to develop a model that underpins a minimally viable library.
767 Foremost is implementing a data object that encodes data with a 2D continuous topology
768 and an artist that can consume data with a 2D topology to visualize the image[43, 44, 94] and
769 also encoding a separate heatmap[60, 104] artist that consumes 1D discrete data. A second
770 important proof of concept artist is a boxplot[102] because it is a graphic that assumes

771 computation on the data side and the glyph is built from semantically defined components
772 and a list of outliers. The model supports simple composition of glyphs by overlaying glyphs
773 at the same position, but more work is needed to define an operator where the fiber bundles
774 have shared $S_2 \hookrightarrow S_1$ such that fibers could be pulled back over the subset. While the
775 model's simple addition supports axes as standalone artists with overlapping visual position
776 encoding, the complex operator would allow for binding together data that needs to be
777 updated together. Additionally, implementing the complex addition operator and explicit
778 graphic to data maps would allow for developing a mathematical formalism and prototype
779 of how interactivity would work in this model. In summary, the proposed scope of work for
780 the dissertation is

- 781 • expansion of the mathematical framework to include complex addition
 - 782 • formalization of definition of equivalence class A'
 - 783 • implementation of artist with explicit ξ
 - 784 • specification of interactive visualization
 - 785 • mathematical formulation of a graphic with axes labeling
 - 786 • implementation of new prototype artists that do not inherit from Matplotlib artists
 - 787 • provisional mathematics and implementation of user level composite artists
 - 788 • proof of concept domain specific user facing library
- 789 Other potential tasks for future work is implementing a data object for a non-trivial fiber
790 bundle and exploiting the models section level formalism to build distributed data source
791 models and concurrent artists. This could be pushed further to integrate with topologi-
792 cal[50] and functional [78] data analysis methods. Since this model formalizes notions of
793 structure preservation, it can serve as a good base for tools that assess quality metrics[12]
794 or invariance [56] of visualizations with respect to graphical encoding choices. While this
795 paper formulates visualization in terms of monoidal action homomorphisms between fiber-

796 bundles, the model lends itself to a categorical formulation[37, 65] that could be further
797 explored.

798 6 Conclusion

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

807 References

- 808 [1] A. Sarikaya et al. “What Do We Talk About When We Talk About Dashboards?”
809 In: *IEEE Transactions on Visualization and Computer Graphics* 25.1 (Jan. 2019),
810 pp. 682–692. ISSN: 1941-0506. DOI: 10.1109/TVCG.2018.2864903.
- 811 [2] [*A Series of Statistical Charts Illustrating the Condition of the Descendants of Former African Slaves Now in Residence in the United States of America*] Negro Business Men in the United States. eng. <https://www.loc.gov/item/2014645363/>. Image.
- 814 [3] [*A Series of Statistical Charts Illustrating the Condition of the Descendants of Former African Slaves Now in Residence in the United States of America*] Negro Population of the United States Compared with the Total Population of Other Countries /. eng.
815 <https://www.loc.gov/item/2013650368/>. Image.
- 817 [4] *Action in nLab*. https://ncatlab.org/nlab/show/action#actions_of_a_monoid.
- 819 [5] James Ahrens, Berk Geveci, and Charles Law. “Paraview: An End-User Tool for Large Data Visualization”. In: *The visualization handbook* 717.8 (2005).

- 821 [6] Yael Albo et al. “Off the Radar: Comparative Evaluation of Radial Visualization
822 Solutions for Composite Indicators”. In: *IEEE Transactions on Visualization and*
823 *Computer Graphics* 22.1 (Jan. 2016), pp. 569–578. ISSN: 1077-2626. DOI: 10.1109/
824 *TVCG*.2015.2467322.
- 825 [7] *Anti-Grain Geometry* -. <http://agg.sourceforge.net/antigrain.com/index.html>.
- 826 [8] Professor Denis Auroux. “Math 131: Introduction to Topology”. en. In: (), p. 113.
- 827 [9] F. Beck. “Software Feathers Figurative Visualization of Software Metrics”. In: *2014*
828 *International Conference on Information Visualization Theory and Applications*
829 (*IVAPP*). Jan. 2014, pp. 5–16.
- 830 [10] Richard A. Becker and William S. Cleveland. “Brushing Scatterplots”. In: *Techno-*
831 *metrics* 29.2 (May 1987), pp. 127–142. ISSN: 0040-1706. DOI: 10.1080/00401706.
832 1987.10488204.
- 833 [11] Jacques Bertin. *Semiology of Graphics : Diagrams, Networks, Maps*. English. Red-
834 lands, Calif.: ESRI Press, 2011. ISBN: 978-1-58948-261-6 1-58948-261-1.
- 835 [12] Enrico Bertini, Andrada Tat, and Daniel Keim. “Quality Metrics in High-
836 Dimensional Data Visualization: An Overview and Systematization”. In: *IEEE*
837 *Transactions on Visualization and Computer Graphics* 17.12 (2011), pp. 2203–2212.
- 838 [13] Tim Bienz, Richard Cohn, and Calif.) Adobe Systems (Mountain View. *Portable*
839 *Document Format Reference Manual*. Citeseer, 1993.
- 840 [14] M. Bostock and J. Heer. “Protopis: A Graphical Toolkit for Visualization”. In: *IEEE*
841 *Transactions on Visualization and Computer Graphics* 15.6 (Nov. 2009), pp. 1121–
842 1128. ISSN: 1941-0506. DOI: 10.1109/TVCG.2009.174.
- 843 [15] M. Bostock, V. Ogievetsky, and J. Heer. “D³ Data-Driven Documents”. In: *IEEE*
844 *Transactions on Visualization and Computer Graphics* 17.12 (Dec. 2011), pp. 2301–
845 2309. ISSN: 1941-0506. DOI: 10.1109/TVCG.2011.185.
- 846 [16] Brian Wylie and Jeffrey Baumes. “A Unified Toolkit for Information and Scientific
847 Visualization”. In: *Proc.SPIE*. Vol. 7243. Jan. 2009. DOI: 10.1117/12.805589.

- 848 [17] Andreas Buja et al. “Interactive Data Visualization Using Focusing and Linking”. In:
849 *Proceedings of the 2nd Conference on Visualization '91*. VIS '91. Washington, DC,
850 USA: IEEE Computer Society Press, 1991, pp. 156–163. ISBN: 0-8186-2245-8.
- 851 [18] D. M. Butler and M. H. Pendley. “A Visualization Model Based on the Mathematics
852 of Fiber Bundles”. en. In: *Computers in Physics* 3.5 (1989), p. 45. ISSN: 08941866.
853 DOI: 10.1063/1.168345.
- 854 [19] David M. Butler and Steve Bryson. “Vector-Bundle Classes Form Powerful Tool
855 for Scientific Visualization”. en. In: *Computers in Physics* 6.6 (1992), p. 576. ISSN:
856 08941866. DOI: 10.1063/1.4823118.
- 857 [20] L. Byrne, D. Angus, and J. Wiles. “Acquired Codes of Meaning in Data Visualization
858 and Infographics: Beyond Perceptual Primitives”. In: *IEEE Transactions on Visual-*
859 *ization and Computer Graphics* 22.1 (Jan. 2016), pp. 509–518. ISSN: 1077-2626. DOI:
860 10.1109/TVCG.2015.2467321.
- 861 [21] Lydia Byrne, Daniel Angus, and Janet Wiles. “Figurative Frames: A Critical Vocab-
862 ulary for Images in Information Visualization”. In: *Information Visualization* 18.1
863 (Aug. 2017), pp. 45–67. ISSN: 1473-8716. DOI: 10.1177/1473871617724212.
- 864 [22] *Cairoglyphics.Org*. <https://www.cairoglyphics.org/>.
- 865 [23] Sheelagh Carpendale. *Visual Representation from Semiology of Graphics by J. Bertin*.
866 en.
- 867 [24] George S. Carson. “Standards Pipeline: The OpenGL Specification”. In: *SIGGRAPH*
868 *Comput. Graph.* 31.2 (May 1997), pp. 17–18. ISSN: 0097-8930. DOI: 10.1145/271283.
869 271292.
- 870 [25] A. M. Cegarra. “Cohomology of Monoids with Operators”. In: *Semigroup Forum*.
871 Vol. 99. Springer, 2019, pp. 67–105. ISBN: 1432-2137.
- 872 [26] John M Chambers et al. *Graphical Methods for Data Analysis*. Vol. 5. Wadsworth
873 Belmont, CA, 1983.

- 874 [27] Chia-Shang James Chu. “Time Series Segmentation: A Sliding Window Approach”.
875 In: *Information Sciences* 85.1 (July 1995), pp. 147–173. ISSN: 0020-0255. DOI: 10.
876 1016/0020-0255(95)00021-G.
- 877 [28] William S. Cleveland. “Research in Statistical Graphics”. In: *Journal of the American
878 Statistical Association* 82.398 (June 1987), p. 419. ISSN: 01621459. DOI: 10.2307/
879 2289443.
- 880 [29] William S. Cleveland and Robert McGill. “Graphical Perception: Theory, Experi-
881 mentation, and Application to the Development of Graphical Methods”. In: *Journal
882 of the American Statistical Association* 79.387 (Sept. 1984), pp. 531–554. ISSN: 0162-
883 1459. DOI: 10.1080/01621459.1984.10478080.
- 884 [30] “Connected Space”. en. In: *Wikipedia* (Dec. 2020).
- 885 [31] Michael S. Crouch, Andrew McGregor, and Daniel Stubbs. “Dynamic Graphs in the
886 Sliding-Window Model”. In: *European Symposium on Algorithms*. Springer, 2013,
887 pp. 337–348.
- 888 [32] *Data Representation in Mayavi — Mayavi 4.7.2 Documentation*. <https://docs.enthought.com/mayavi/mayavi/d>
- 889 [33] *Dataclasses — Data Classes — Python 3.9.2rc1 Documentation*. <https://docs.python.org/3/library/dataclasses>.
- 890 [34] W. E. B. (William Edward Burghardt) Du Bois. *[The Georgia Negro] Valuation of
891 Town and City Property Owned by Georgia Negroes*. en. <https://www.loc.gov/pictures/item/2013650441/>.
892 Still Image. 1900.
- 893 [35] T. W. E. B. Du Bois Center at the University of Massachusetts, W. Battle-Baptiste,
894 and B. Rusert. *W. E. B. Du Bois’s Data Portraits: Visualizing Black America*.
895 Princeton Architectural Press, 2018. ISBN: 978-1-61689-706-2.
- 896 [36] Stephen Few and Perceptual Edge. “Dashboard Confusion Revisited”. In: *Perceptual
897 Edge* (2007), pp. 1–6.
- 898 [37] Brendan Fong and David I. Spivak. *An Invitation to Applied Category Theory:
899 Seven Sketches in Compositionality*. en. First. Cambridge University Press, July

- 900 2019. ISBN: 978-1-108-66880-4 978-1-108-48229-5 978-1-108-71182-1. DOI: 10.1017/
 901 9781108668804.
- 902 [38] Michael Friendly. “A Brief History of Data Visualization”. en. In: *Handbook of Data*
 903 *Visualization*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 15–56. ISBN:
 904 978-3-540-33036-3 978-3-540-33037-0. DOI: 10.1007/978-3-540-33037-0_2.
- 905 [39] Berk Geveci et al. “VTK”. In: *The Architecture of Open Source Applications* 1 (2012),
 906 pp. 387–402.
- 907 [40] Robert Ghrist. “Homological Algebra and Data”. In: *Math. Data* 25 (2018), p. 273.
- 908 [41] Robert W. Ghrist. *Elementary Applied Topology*. Vol. 1. Createspace Seattle, 2014.
- 909 [42] Kristen B. Gorman, Tony D. Williams, and William R. Fraser. “Ecological Sexual
 910 Dimorphism and Environmental Variability within a Community of Antarctic Pen-
 911 guins (Genus Pygoscelis)”. In: *PLOS ONE* 9.3 (Mar. 2014), e90081. DOI: 10.1371/journal.pone.0090081.
- 912 [43] Robert B Haber and David A McNabb. “Visualization Idioms: A Conceptual Model
 913 for Scientific Visualization Systems”. In: *Visualization in scientific computing* 74
 914 (1990), p. 93.
- 915 [44] Charles D Hansen and Chris R Johnson. *Visualization Handbook*. Elsevier, 2011.
- 916 [45] Marcus D. Hanwell et al. “The Visualization Toolkit (VTK): Rewriting the Rendering
 917 Code for Modern Graphics Cards”. en. In: *SoftwareX* 1-2 (Sept. 2015), pp. 9–12. ISSN:
 918 23527110. DOI: 10.1016/j.softx.2015.04.001.
- 919 [46] Charles R Harris et al. “Array Programming with NumPy”. In: *Nature* 585.7825
 920 (2020), pp. 357–362.
- 921 [47] J. Heer and M. Agrawala. “Software Design Patterns for Information Visualization”.
 922 In: *IEEE Transactions on Visualization and Computer Graphics* 12.5 (2006), pp. 853–
 923 860. DOI: 10.1109/TVCG.2006.178.
- 924

- 925 [48] Jeffrey Heer and Michael Bostock. “Declarative Language Design for Interactive Vi-
926 sualization”. In: *IEEE Transactions on Visualization and Computer Graphics* 16.6
927 (Nov. 2010), pp. 1149–1156. ISSN: 1077-2626. DOI: 10.1109/TVCG.2010.144.
- 928 [49] Jeffrey Heer, Michael Bostock, and Vadim Ogievetsky. “A Tour Through the Visu-
929 alization Zoo”. In: *Communications of the ACM* 53.6 (June 2010), pp. 59–67. ISSN:
930 0001-0782. DOI: 10.1145/1743546.1743567.
- 931 [50] C. Heine et al. “A Survey of Topology-Based Methods in Visualization”. In: *Computer*
932 *Graphics Forum* 35.3 (June 2016), pp. 643–667. ISSN: 0167-7055. DOI: 10.1111/cgf.
933 12933.
- 934 [51] Allison Marie Horst, Alison Presmanes Hill, and Kristen B Gorman. *Palmerpen-*
935 *guins: Palmer Archipelago (Antarctica) Penguin Data*. Manual. 2020. DOI: 10.5281/
936 zenodo.3960218.
- 937 [52] Stephan Hoyer and Joe Hamman. “Xarray: ND Labeled Arrays and Datasets in
938 Python”. In: *Journal of Open Research Software* 5.1 (2017).
- 939 [53] J. D. Hunter. “Matplotlib: A 2D Graphics Environment”. In: *Computing in Science*
940 *Engineering* 9.3 (May 2007), pp. 90–95. ISSN: 1558-366X. DOI: 10.1109/MCSE.2007.
941 55.
- 942 [54] John Hunter and Michael Droettboom. *The Architecture of Open Source Applications*
943 (*Volume 2*): Matplotlib. <https://www.aosabook.org/en/matplotlib.html>.
- 944 [55] “Jet Bundle”. en. In: *Wikipedia* (Dec. 2020).
- 945 [56] Gordon Kindlmann and Carlos Scheidegger. “An Algebraic Process for Visualization
946 Design”. In: *IEEE transactions on visualization and computer graphics* 20.12 (2014),
947 pp. 2181–2190.
- 948 [57] John Krygier and Denis Wood. *Making Maps: A Visual Guide to Map Design for*
949 *GIS*. English. 1 edition. New York: The Guilford Press, Aug. 2005. ISBN: 978-1-59385-
950 200-9.
- 951 [58] W A Lea. “A Formalization of Measurement Scale Forms”. en. In: (), p. 44.

- 952 [59] *Locally Trivial Fibre Bundle* - Encyclopedia of Mathematics. https://encyclopediaofmath.org/wiki/Locally_trivial_fibre_bundle
- 953 [60] Toussaint Loua. *Atlas Statistique de La Population de Paris*. J. Dejey & cie, 1873.
- 954 [61] 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.
- 955 [62] Jock Mackinlay. “Automating the Design of Graphical Presentations of Relational Information”. In: *ACM Transactions on Graphics* 5.2 (Apr. 1986), pp. 110–141. ISSN: 0730-0301. DOI: [10.1145/22949.22950](https://doi.org/10.1145/22949.22950).
- 956 [63] JOCK D. MACKINLAY. “AUTOMATIC DESIGN OF GRAPHICAL PRESENTATIONS (DATABASE, USER INTERFACE, ARTIFICIAL INTELLIGENCE, INFORMATION TECHNOLOGY)”. English. PhD Thesis. 1987.
- 957 [64] Connie Malamed. *Information Display Tips*. <https://understandinggraphics.com/visualizations/information-display-tips/>. Blog. Jan. 2010.
- 958 [65] Bartosz Milewski. “Category Theory for Programmers”. en. In: (), p. 498.
- 959 [66] “Monoid”. en. In: *Wikipedia* (Jan. 2021).
- 960 [67] Tamara Munzner. “Ch 2: Data Abstraction”. In: *CPSC547: Information Visualization, Fall 2015-2016* ().
- 961 [68] Tamara Munzner. *Visualization Analysis and Design*. AK Peters Visualization Series. CRC press, Oct. 2014. ISBN: 978-1-4665-0891-0.
- 962 [69] 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).
- 963 [70] Muhammad Chenariyan Nakhaee. *Mcnakhaee/Palmerpenguins*. Jan. 2021.
- 964 [71] *Naturalness Principle - InfoVis:Wiki*. https://infovis-wiki.net/wiki/Naturalness_Principle.

- 976 [72] Dmitry Nekrasovski et al. “An Evaluation of Pan & Zoom and Rubber Sheet
 977 Navigation with and without an Overview”. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. CHI ’06. New York, NY, USA:
 978 Association for Computing Machinery, 2006, pp. 11–20. ISBN: 1-59593-372-7. DOI:
 979 10.1145/1124772.1124775.
- 981 [73] Z. Pousman, J. Stasko, and M. Mateas. “Casual Information Visualization: Depic-
 982 tions of Data in Everyday Life”. In: *IEEE Transactions on Visualization and Com-*
 983 *puter Graphics* 13.6 (Nov. 2007), pp. 1145–1152. ISSN: 1941-0506. DOI: 10.1109/TVCG.2007.70541.
- 985 [74] Z. Qu and J. Hullman. “Keeping Multiple Views Consistent: Constraints, Validations,
 986 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.
- 989 [75] A. Quint. “Scalable Vector Graphics”. In: *IEEE MultiMedia* 10.3 (July 2003), pp. 99–
 990 102. ISSN: 1941-0166. DOI: 10.1109/MMUL.2003.1218261.
- 991 [76] “Quotient Space (Topology)”. en. In: *Wikipedia* (Nov. 2020).
- 992 [77] P. Ramachandran and G. Varoquaux. “Mayavi: 3D Visualization of Scientific Data”.
 993 In: *Computing in Science Engineering* 13.2 (Mar. 2011), pp. 40–51. ISSN: 1558-366X.
 994 DOI: 10.1109/MCSE.2011.35.
- 995 [78] James O Ramsay. *Functional Data Analysis*. Wiley Online Library, 2006.
- 996 [79] Jeff Reback et al. *Pandas-Dev/Pandas: Pandas 1.0.3*. Zenodo. Mar. 2020. DOI: 10.5281/zenodo.3715232.
- 998 [80] “Retraction (Topology)”. en. In: *Wikipedia* (July 2020).
- 999 [81] Matthew Rocklin. “Dask: Parallel Computation with Blocked Algorithms and Task
 1000 Scheduling”. In: *Proceedings of the 14th Python in Science Conference*. Vol. 126.
 1001 Citeseer, 2015.

- 1002 [82] Arvind Satyanarayan, Kanit Wongsuphasawat, and Jeffrey Heer. “Declarative Inter-
 1003 action Design for Data Visualization”. en. In: *Proceedings of the 27th Annual ACM*
 1004 *Symposium on User Interface Software and Technology*. Honolulu Hawaii USA: ACM,
 1005 Oct. 2014, pp. 669–678. ISBN: 978-1-4503-3069-5. DOI: 10.1145/2642918.2647360.
- 1006 [83] Caroline A Schneider, Wayne S Rasband, and Kevin W Eliceiri. “NIH Image to
 1007 ImageJ: 25 Years of Image Analysis”. In: *Nature Methods* 9.7 (July 2012), pp. 671–
 1008 675. ISSN: 1548-7105. DOI: 10.1038/nmeth.2089.
- 1009 [84] “Semigroup Action”. en. In: *Wikipedia* (Jan. 2021).
- 1010 [85] Nicholas Sofroniew et al. *Napari/Napari: 0.4.5rc1*. Zenodo. Feb. 2021. DOI: 10.5281/
 1011 zenodo.4533308.
- 1012 [86] E.H. Spanier. *Algebraic Topology*. McGraw-Hill Series in Higher Mathematics.
 1013 Springer, 1989. ISBN: 978-0-387-94426-5.
- 1014 [87] David I Spivak. *Databases Are Categories*. en. Slides. June 2010.
- 1015 [88] David I Spivak. “SIMPLICIAL DATABASES”. en. In: (), p. 35.
- 1016 [89] S. S. Stevens. “On the Theory of Scales of Measurement”. In: *Science* 103.2684 (1946),
 1017 pp. 677–680. ISSN: 00368075, 10959203.
- 1018 [90] Michele Stieven. *A Monad Is Just a Monoid...* en. <https://medium.com/@michelestieven/a-monad-is-just-a-monoid-a02bd2524f66>. Apr. 2020.
- 1019 [91] Software Studies. *Culturevis/Imageplot*. Jan. 2021.
- 1020 [92] T. Sugibuchi, N. Spyros, and E. Siminenko. “A Framework to Analyze Information
 1021 Visualization Based on the Functional Data Model”. In: *2009 13th International*
 1022 *Conference Information Visualisation*. 2009, pp. 18–24. DOI: 10.1109/IV.2009.56.
- 1023 [93] *[The Georgia Negro] City and Rural Population. 1890*. eng. <https://www.loc.gov/item/2013650430/>.
 1024 Image. 1900.
- 1025 [94] M. Tory and T. Moller. “Rethinking Visualization: A High-Level Taxonomy”. In:
 1026 *IEEE Symposium on Information Visualization*. Oct. 2004, pp. 151–158. DOI: 10.
 1027 1109/INFVIS.2004.59.

- 1029 [95] Edward R. Tufte. *The Visual Display of Quantitative Information*. English. Cheshire,
1030 Conn.: Graphics Press, 2001. ISBN: 0-9613921-4-2 978-0-9613921-4-7 978-1-930824-13-
1031 3 1-930824-13-0.
- 1032 [96] John W. Tukey. “We Need Both Exploratory and Confirmatory”. In: *The American
1033 Statistician* 34.1 (Feb. 1980), pp. 23–25. ISSN: 0003-1305. DOI: 10.1080/00031305.
1034 1980.10482706.
- 1035 [97] Jacob VanderPlas et al. “Altair: Interactive Statistical Visualizations for Python”.
1036 en. In: *Journal of Open Source Software* 3.32 (Dec. 2018), p. 1057. ISSN: 2475-9066.
1037 DOI: 10.21105/joss.01057.
- 1038 [98] P. Vickers, J. Faith, and N. Rossiter. “Understanding Visualization: A Formal Ap-
1039 proach Using Category Theory and Semiotics”. In: *IEEE Transactions on Visualiza-
1040 tion and Computer Graphics* 19.6 (June 2013), pp. 1048–1061. ISSN: 1941-0506. DOI:
1041 10.1109/TVCG.2012.294.
- 1042 [99] C. Ware. *Information Visualization: Perception for Design*. Interactive Technologies.
1043 Elsevier Science, 2019. ISBN: 978-0-12-812876-3.
- 1044 [100] Eric W. Weisstein. *Similarity Transformation*. en. <https://mathworld.wolfram.com/SimilarityTransformation.html>
1045 Text.
- 1046 [101] Hadley Wickham. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New
1047 York, 2016. ISBN: 978-3-319-24277-4.
- 1048 [102] Hadley Wickham and Lisa Stryjewski. “40 Years of Boxplots”. In: *The American
1049 Statistician* (2011).
- 1050 [103] Leland Wilkinson. *The Grammar of Graphics*. en. 2nd ed. Statistics and Computing.
1051 New York: Springer-Verlag New York, Inc., 2005. ISBN: 978-0-387-24544-7.
- 1052 [104] Leland Wilkinson and Michael Friendly. “The History of the Cluster Heat Map”.
1053 In: *The American Statistician* 63.2 (May 2009), pp. 179–184. ISSN: 0003-1305. DOI:
1054 10.1198/tas.2009.0033.

- 1055 [105] Krist Wongsuphasawat. *Navigating the Wide World of Data Visualization Libraries*
1056 (*on the Web*). 2021.
- 1057 [106] *Writing Plugins*. en. https://imagej.net/Writing_plugins.
- 1058 [107] Brent A Yorgey. “Monoids: Theme and Variations (Functional Pearl)”. en. In: (),
1059 p. 12.
- 1060 [108] Caroline Ziemkiewicz and Robert Kosara. “Embedding Information Visualization
1061 within Visual Representation”. In: *Advances in Information and Intelligent Systems*.
1062 Ed. by Zbigniew W. Ras and William Ribarsky. Berlin, Heidelberg: Springer Berlin
1063 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).