

# Colorgorical: Creating discriminable and preferable color palettes for information visualization

Connor C. Gramazio, *Student Member, IEEE*, David H. Laidlaw, *Fellow, IEEE*, Karen B. Schloss

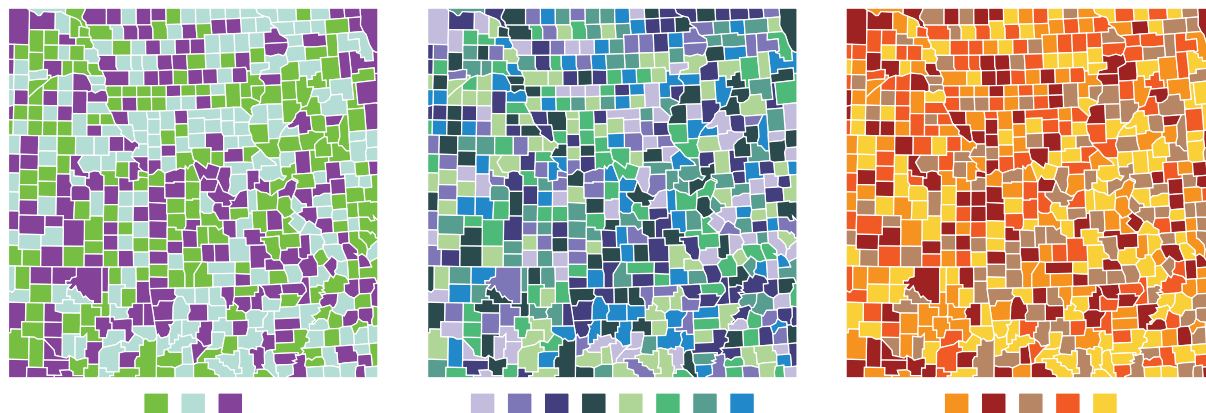


Fig. 1. Colorgorical palettes: (left) a 3-color palette generated with high name difference; (center) an 8-color palette generated with high pair preference; (right) a 5-color palette generated with high pair preference, medium perceptual distance, and a hue filter.

**Abstract**— We present an evaluation of Colorgorical, a web-based tool for creating discriminable and aesthetically preferable categorical color palettes. Colorgorical uses iterative semi-random sampling to pick colors from CIELAB space based on user-defined discriminability and preference importances. Colors are selected by assigning each a weighted sum score that applies the user-defined importances to Perceptual Distance, Name Difference, Name Uniqueness, and Pair Preference scoring functions, which compare a potential sample to already-picked palette colors. After, a color is added to the palette by randomly sampling from the highest scoring palettes. Users can also specify hue ranges or build off their own starting palettes. This procedure differs from previous approaches that do not allow customization (e.g., pre-made ColorBrewer palettes) or do not consider visualization design constraints (e.g., Adobe Color and ACE). In a Palette Score Evaluation, we verified that each scoring function measured different color information. Experiment 1 demonstrated that slider manipulation generates palettes that are consistent with the expected balance of discriminability and aesthetic preference for 3-, 5-, and 8-color palettes, and also shows that the number of colors may change the effectiveness of pair-based discriminability and preference scores. For instance, if the Pair Preference slider were upweighted, users would judge the palettes as more preferable on average. Experiment 2 compared Colorgorical palettes to benchmark palettes (ColorBrewer, Microsoft, Tableau, Random). Colorgorical palettes are as discriminable and are at least as preferable or more preferable than the alternative palette sets. In sum, Colorgorical allows users to make customized color palettes that are, on average, as effective as current industry standards by balancing the importance of discriminability and aesthetic preference.

**Index Terms**—Aesthetics in Visualization, Color Perception, Metrics & Benchmarks, Visual Design, Visualization

## 1 INTRODUCTION

We present an evaluation of Colorgorical (Fig. 2), a model and tool for creating arbitrarily sized, preferable, and discriminable color palettes for categorical information visualization (Fig. 1). As in other areas of design, it is important that a visualization color palette is aesthetically pleasing; but, unlike many other areas of design, visualization color palettes must also be highly discriminable. Balancing discriminability and aesthetic preference is challenging because they can be inversely related (i.e., preference increases with hue similarity [30], whereas discriminability decreases). Navigating this tradeoff requires design

skill and experience, both beyond those of many visualization creators.

Colorgorical addresses this problem by operationalizing effective color palette selection with three color-scoring functions to balance discriminability and aesthetic preference: Perceptual Distance (CIEDE2000) [32], Name Difference [8], and a quantified model of color Pair Preference [30] (Sec. 3). A fourth, Name Uniqueness, was originally included, but was later removed because it had little effect on behavior (Sec. 6). With Colorgorical, color palette creation is simplified so that users need only specify the number of desired colors and drag sliders controlling color-scoring function importance to (1) create custom palettes that the average individual would find preferable while maintaining discriminability, and (2) explore how relative weights on discriminability vs. preference affect palette appearance. Users can further customize palettes by specifying desired hues and by building onto existing palettes (Sec. 4).

We evaluated Colorgorical's effectiveness in four ways: (1) runtime benchmarks (Sec. 4), (2) discriminability and preference score analysis (Sec. 5), (3) human-subject evaluation of different model settings (Sec. 6), and (4) human-subject evaluation of Colorgorical compared to industry standards (Sec. 7). We make the following contributions:

- We provide a technique to generate custom color palettes via

• C.C. Gramazio and D.H. Laidlaw are with the Dept. of Computer Science at Brown University. E-mail: {connor,dhl}@cs.brown.edu.

• K.B. Schloss is with the Dept. of Cognitive, Linguistic, and Psychological Sciences at Brown University. E-mail: karen Schloss@gmail.com.

Manuscript received 31 Mar. 2016; accepted 1 Aug. 2016. Date of publication 15 Aug. 2016; date of current version 23 Oct. 2016.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TVCG.2016.2598918

user-defined importance of discriminability and preference

- We detail the relations between Perceptual Distance, Name Difference, Name Uniqueness, and Pair Preference scoring functions
- We show how varying the relative weights of discriminability and preference sliders affects human discrimination performance and preference ratings
- We present evidence that Colorgical palettes are as discriminable and often more preferable than industry standard, professionally hand-made color palettes

Colorgical combines three features, making it a novel approach to palette design. First, it is designed specifically for visualization rather than for general art and design applications. Second, it uses empirically derived color preference data to inform categorical palette generation [30]. Third, it approaches visualization palette design by balancing categorical palette discriminability and preference.

## 2 RELATED WORK

Current color palette tools are typically designed based on three types of strategies: discriminability optimization, color-term association mapping, or harmonic template application. We describe these approaches and discuss how Colorgical targets limitations of past research.

### 2.1 Palette discriminability methods

A key issue in palette discriminability is whether a graphical mark can be quickly and accurately identified. Healey demonstrated that this problem can be addressed by using palettes whose colors are named with the 10 Munsell hues and that maximize perceptual distance between colors (CIEDE1976, Sec. 3) [7]. Maxwell also developed a discriminability-based technique to create categorical color palettes for multidimensional datasets based on classification dissimilarity of categories [16]. These approaches created discriminable palettes, but each has multiple limitations for design more broadly: (1) they do not address aesthetics, (2) Healey's technique is constrained to 10 or fewer color terms, and (3) they define perceptual distance using Euclidean distance (Healey) or maximum scaled difference (Maxwell) in CIELAB color space, which can be problematic due to perceptual uniformity limitations [13] (i.e., the same distance can have different perceptual consequences depending on the sampled region).

Colorgical addresses these issues by (1) considering aesthetics in addition to discriminability (Sec. 3) [30], (2) using 153 crowdsourced color terms compared to the 10 Munsell hues in Healey's method, and (3) using an updated perceptual distance function (CIEDE2000) that improves perceptual uniformity in the distance metric [32].

The color-name associations in Colorgical are based on Heer and Stone's color-name statistics (Sec. 3) [8], which are derived from color-name association frequencies from the 153 most commonly-used names from the XKCD color-name crowdsourcing survey [19]. Name Difference measures the difference in color-name association frequency distributions between two colors. For example, green and red colors have large name differences because green colors have few associations with red names and vice versa. Presumably Name Difference is related to Perceptual Distance, but it is possible that they differ systematically, which we test in Sections 5 and 6. Name Salience, which we call Name Uniqueness to avoid confusion with color salience, captures the degree to which a color is specifically named (highly associated with only a few colors) vs. broadly named (moderately associated with many colors) (Fig. 1 in Supp. Mat.).

Another approach to designing discriminable palettes is for color experts to make pre-defined palettes (e.g., ColorBrewer [6]). Typically made through iterative design, experts construct these palettes by selecting colors that are discriminable under a variety of viewing conditions (e.g., after photocopying) and that support specialized tasks (e.g., ColorBrewer's "Accent" palettes emphasize certain colors). Although pre-made palettes are easy to use, they do not give visualization creators design flexibility or customizability. And although guidelines for hand-designing palettes exist [37], a visualization creator might not

want to spend time or effort to learn about palette design. Colorgical addresses this problem by allowing customization while building in constraints on aesthetics and discriminability; however, we leave support for specialized palettes (e.g., accent colors) for future research.

### 2.2 Color-term tools

Another way to create categorical palettes is through color-term associations. Crowdsourcing and linguistics-based approaches can produce color-term associations that create semantically meaningful palettes (e.g., a "mango ice cream" category might produce a light orange) [11, 31]. Setlur and Stone show that various natural language processing techniques can be used to mine color-semantic pairings from large text datasets [31]. Colorgical does not currently support semantic mappings, but it is an exciting future direction.

### 2.3 Harmonic template tools

Many harmony-based categorical color palette tools are targeted for general-purpose design and do not focus on visualization design constraints (e.g., discriminability). These tools create palettes based on harmony principles in color theory [20, 22]. A common implementation of harmony is through *harmonic templates* based on hue relations [15], such as the two-color complementary relation that stems from Itten's version of harmony (e.g., blue and orange) [10]. For example, Adobe Color creates 5-color palettes based on harmonic templates and optional image color analysis [21]. Similarly, Dial-a-color, uses harmonic templates as a starting point and allows users to alter color properties like lightness and saturation [18]. ACE lets users manipulate discrimination and harmony importance for interface design by answering a series of questions in a text interface about each colored interface component [17] (unlike ACE, Colorgical is not limited to interface coloration and uses sliders to balance discrimination and aesthetic preference rather than a text interface). Finally, the Harmonious Color Scheme Generator constructs color palettes through *familial factors* (promoting similarity along hue, saturation, or lightness dimensions) and *rhythmic spans* (sampling colors using a fixed uniform interval along a color dimension) [9].

Harmonic templates were generated from color theory in art without empirical validation [10], and do not necessarily correspond to human judgments of harmony. For example, the notion that complementary colors are harmonious is key to the notion of harmonic templates. Yet humans judge complementary hues as among the least harmonious and instead judge more similar hues as more harmonious [23, 24, 30, 34].

Although the term "harmony" is often used interchangeably with aesthetic preference [2], the two are not the same [24, 30]. Schloss and Palmer demonstrated how they differ, where harmony was defined as "how strongly an observer experiences the colors in the combination as going or belonging together, regardless of whether the observer likes the combination or not," and preference is "how much an observer likes a given pair of colors as a Gestalt, or whole" [30]. Although both increased with hue similarity, pair preference relied more on preference ratings for individual colors and on lightness contrast, whereas harmony relied more on desaturation (i.e., pairs with less saturated colors were more harmonious).

Colorgical uses Schloss and Palmer's pair preference model (Sec. 3) [30] rather than harmony because we reasoned that how much people like visualization palette colors is more central to the present aims than how well they feel the colors go together.

## 3 BACKGROUND: MODEL SCORING FUNCTIONS

Colorgical iteratively samples colors using three color discriminability scores (*Perceptual Distance*, *Name Difference*, *Name Uniqueness*) and a color preference score (*Pair Preference*). Colorgical assumes that discriminability and preference for large combinations of colors can be predicted by these lower-order scores. Name Uniqueness was ultimately removed from the model because it had little effect on discriminability performance or preference (Sec. 6).

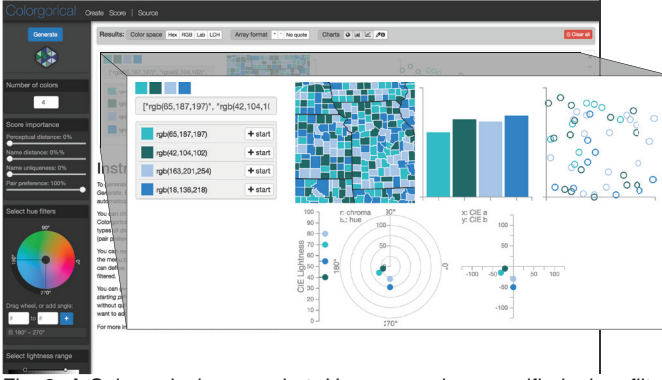


Fig. 2. A Colorgorical screenshot. Here a user has specified a hue filter (left) and has generated a 4-color palette (detail). Users can list colors in many color spaces and render colors in a variety of charts.

Each score operates in CIELAB. The  $L^*$  axis of CIELAB approximates a color's lightness, the  $a^*$  axis approximates its redness-to-greenness, and the  $b^*$  axis approximates its blueness-to-yellowness. To support Name Difference and Name Uniqueness, we use a modification of CIELAB that quantizes the space into 8,325 discrete colors by sampling every 5 units along each axis starting at the origin [8]. Some scores also depend on CIE LCh, which is a polar representation of the Euclidean CIELAB space. In CIE LCh,  $L^*$  is the same as in CIELAB, but the  $a^*$  and  $b^*$  axis are converted to chroma ( $C$ , radius) and hue ( $h$ , angle).

### 3.1 Color discriminability scores

We used multiple discriminability scores because perceptual difference might differ from name difference. For instance, a chartreuse (yellow-green) might be perceptually distinct from a green or yellow but might be called green or yellow, making it easy to confuse with other greens or yellows in a visualization when referenced by name.

#### 3.1.1 CIEDE2000: Perceptual Distance

To calculate Perceptual Distance between two colors we use CIEDE2000 ( $DE_{00}$ ) [32]. It is similar to the original CIEDE,  $DE_{76}$  (Euclidean CIELAB), but  $DE_{00}$  calculates distance in CIE LCh with a hue rotation term ( $R_T$ ) and corrections for lightness ( $S_L$ ), chroma ( $S_C$ ), and hue ( $S_H$ ) to improve perceptually uniformity [14].

$$DE_{76} = \sqrt{\Delta L^2 + \Delta a^2 + \Delta b^2} \quad (1)$$

$$DE_{00} = \sqrt{\left(\frac{\Delta L}{S_L}\right)^2 + \left(\frac{\Delta C}{S_C}\right)^2 + \left(\frac{\Delta H}{S_H}\right)^2 + R_T \frac{\Delta C}{S_C} \frac{\Delta H}{S_H}} \quad (2)$$

#### 3.1.2 Name Difference

Name Difference (ND) captures the degree to which two colors have distinct color-name association frequency distributions [8]. Color-name associations are mappings between colors and names (e.g.,  $\text{rgb}(255, 0, 0) \rightarrow \text{"bright red"}$ ). The name data are composed of the discretized CIELAB color space ( $C$ ) described earlier, a list of 153 popular color names ( $W$ ), and a color-name association frequency matrix ( $T$ ) that has  $C$  rows and  $W$  columns. The scores also rely on the conditional probability of a color name  $w$  given any color in  $C$ :

$$p(w|c) = T_{c,w} / \sum_w T_{c,w} \quad (3)$$

We calculate Name Difference using Hellinger distance [8]:

$$ND(c_1, c_2) = \sqrt{1 - \sum_{w \in W} \sqrt{p(w|c_1)p(w|c_2)}} \quad (4)$$

#### 3.1.3 Name Uniqueness

Name Uniqueness (NU) captures the degree to which colors have uniform distributions of color-name association frequencies. Colors that have few strongly associated names (i.e., a focal distribution) result in lower scores, whereas colors that have many weakly associated names (i.e., a more-uniform distribution) result in higher scores. Name Uniqueness is calculated by using the negative entropy of a color's name-association frequency distribution from the color-name-association matrix ( $T$ ) and the list of color names ( $W$ ):

$$NU(c) = -H(p(W|c)) = -\sum_{w \in W} p(w|c) \log p(w|c) \quad (5)$$

Unlike the other two discriminability measures, Name Uniqueness relies on individual colors rather than relations between other colors within the palette. We believe this is a key reason why it was not useful in the Colorgorical model (Sec. 6).

### 3.2 Aesthetic preference score: Pair Preference

Pair Preference (PP) is based on a linear regression model used to predict pair preferences from three color-appearance and color-relation factors, which was previously operationalized in Munsell space [30]. The best-fit model explained 53.5% of the variance in pair preference judgments with three factors: coolness ( $\kappa$ ), hue similarity  $\Delta H$ , and lightness contrast  $\Delta L$ . We have altered the original equation to use CIE LCh rather than Munsell color space coordinates, as is reflected in the hue similarity and lightness contrast terms<sup>1</sup>. Coolness scores are calculated in CIE LCh using a linear interpolation of the original 32 color-coolness mappings, which approximates the number of hue-steps a color is from Munsell 10R, such that greenish blues are cool and orangish reds are not cool (Supp. Mat.). The Pair Preference scoring function reflects people's preference for color combinations that contain cool colors that differ in lightness and are similar in hue.

$$PP(c_1, c_2) = 75.15(\kappa_1 + \kappa_2) + 47.61|\Delta L| - 46.42|\Delta H| \quad (6)$$

## 4 COLORGORICAL MODEL

Colorgorical generates color palettes using iterative semi-random sampling. Users specify the number of desired colors and use sliders to set the relative balance of aesthetic preference and discriminability (Sec. 3). Generated palettes are displayed to the user as a swatch, map, bar chart, and scatterplot, which highlights how the discriminability may shift with different types and sizes of graphical marks [4, 33].

### 4.1 Minimum discriminability & preference assertions

Each palette is built from an 8,325-color discretized D65 CIELAB space (Sec. 3) and is additionally filtered in three ways to help increase discriminability and preference, which we describe below: (1) noticeable difference; (2) lightness clamping (from  $L^* = 25$  to  $L^* = 85$ ) and (3) filtering the dark yellow (generally disliked) region of color space. Although the same RGB coordinates can result in different CIELAB colors on different monitors if monitors are uncalibrated, Stone et al. show that using a fixed correspondence between D65 CIELAB and RGB can be used effectively for online tools in practice [33].

**Discriminability** The model enforces a lower discriminability bound by sampling *noticeably different colors* using Stone et al.'s noticeable difference function, which provides a minimum CIELAB interval required to discriminate the colors of two graphical marks more than 50% of the time (based on their physical size) [33]. We use a small, conservative visual angle in our calculations ( $1/3^\circ$ ) and multiply the function's suggested interval by three for extra caution.

To ensure discriminability we also exclude colors that are lighter than  $L = 85$  and darker than  $L = 25$  so that all colors are visible on black or white backgrounds ( $L_{\text{black}} = 0$ ,  $L_{\text{white}} = 100$ ). Colorgorical only includes RGB-valid colors.

<sup>1</sup>The CIE LCh model explains 51.8% of the variance in Schloss and Palmer's preference data (their Munsell-based model explains 53.5%).



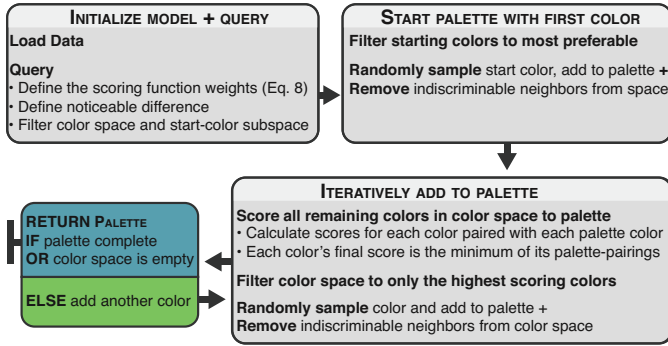


Fig. 3. Diagram of Colorgorical palette construction procedure.

**Preference** Colorgorical excludes the dark yellowish-green region of CIE LCh, which has strongly disliked colors, on average, across many cultures [26, 35, 36]. We define this region as  $L \in [35, 75]$  and  $H \in [85^\circ, 114^\circ]$ . While there are individual differences in preference [25, 30] and some observers may like these colors [29], the goal is to cater to the average observer. This filter was especially important for generating aesthetically preferable discriminable palettes because of the way Pair Preference and discriminability functions interact. In the Pair Preference equation, the coolness term biases selection toward bluish hues and the lightness term biases selection of contrasting lightness. The discriminability functions bias selection for colors that are far apart in CIELAB color space (i.e., contrasting hue and lightness). Once bluish hues are selected, discriminability would be promoted in subsequent color selections by selecting opposite, yellowish hues of a different lightness level (opposite ends of the  $b^*$  and  $L^*$  axes). If the blues are remotely light, then selected yellows will be the dark yellows that people generally dislike. The removal of this region still retained a large region of color space that was sufficiently discriminable to pair with blues, while increasing typical aesthetic palette preference.

To maximize preference within a defined balance, the model generates 10 palettes and returns the palette with the highest minimum-Pair-Preference given all color pairings in each palette.

## 4.2 User-defined model parameters

In addition to specifying the number of colors and manipulating discriminability and preference sliders, users can also configure two optional parameters. First, they can limit color sampling to certain hue ranges (e.g., reds only, or reds and blues), which supports tasks such as designing around brand colors. Second, users can supply an existing palette for Colorgorical to build on. If users provide a palette, Colorgorical rounds the input to the nearest quantized CIELAB color and adds new colors until the palette reaches the desired size.

## 4.3 Palette construction process

Palettes are generated in three steps: (1) initialize, (2) start a palette with the first color, and (3) iteratively add new colors (Fig. 3). Colorgorical can typically generate palettes with up to 22 colors before exhausting color space. However, it is inadvisable to use that many colors due to perceptual limitations [5]. If no more colors can be sampled, Colorgorical returns a partial palette and an error message.

### 4.3.1 Step 1: Initialize

Initialization starts by loading CIELAB space, color coolness scores, and color-name associations into memory. A CIELAB subspace is also loaded into memory, which samples every 15 units along each CIELAB axis and is used along with a precomputed Pair Preference score matrix to pick the first palette color. We use a coarser subspace to select the first color because using precomputed Pair-Preference scores for all pairs of 8,325-colors takes too long for interactivity due to combinatorial explosion. Color space can be filtered based on parameters provided by the user (e.g., hue filters). After applying optional filters, the model limits the subsampled space colors ( $c$ ) and the color

pair preference matrix ( $\Phi$ ) to highly preferable colors (i.e., no dark yellows) using a standard deviation (SD) preference threshold (Eq. 7). The threshold removes any color-pair row from  $\Phi$  whose pair preference score is less than the standard deviation-based limit. Then, the starting color is sampled from the unique colors remaining in  $\Phi$ 's color-pair rows.

$$\text{threshold}(c) = \Phi_c > \max(\Phi) - 0.75 * \text{SD}(\Phi) \quad (7)$$

The last initialization step also defines a noticeable difference with Stone et al.'s CIELAB intervals described above, which removes colors that are too similar to each sampled color. Sampled color differences must have at least one axis above the following intervals:  $\Delta_L = 22.747, \Delta_a = 31.427, \Delta_b = 44.757$ .

### 4.3.2 Step 2: Start palette

The first color of a palette is selected by randomly sampling a seed color from the remaining colors after Step 1. Next, all colors that are not noticeably different from the seed are removed from color space using the CIELAB intervals defined in Step 1. Sampling is skipped if users provide their own seed color(s), but indiscriminable neighboring colors are still eliminated.

### 4.3.3 Step 3: Add to palette

To add a new color, the model computes scores for all remaining colors using a weighted sum (Eq. 8). This function sums each of the four minimum palette scores ( $\tilde{\Psi}$ ) with user-defined weights ( $\tilde{w}$ ), given all possible scores between a potential new color ( $c$ ) and the already picked colors ( $P$ ). The model uses minimum palette scores assuming that a palette is only as discriminable or preferable as its lowest score. There is also a hue-dependent penalty term ( $\tau$ ) to reduce the likelihood of sampling a color bordering the dark yellow filter region. The new color is then randomly sampled from colors that fall above a score threshold (Eq. 7, where  $\Phi$  is now weighted-sum scores). Non-discriminable colors are removed after sampling.

$$\text{score}(c, P) = \tau(\tilde{w} \cdot \tilde{\Psi})$$

$$\tau = \begin{cases} 0.75, & \text{if } 115^\circ < c_{\text{hue}} < 138^\circ \wedge c_L \leq 45 \\ 0.8, & \text{if } 70^\circ \leq c_{\text{hue}} \leq 115^\circ \wedge 45 < c_L \leq 75 \\ 0.85, & \text{if } 70^\circ \leq c_{\text{hue}} \leq 115^\circ \wedge c_L > 75 \\ 1, & \text{otherwise} \end{cases} \quad (8)$$

## 4.4 Implementation and Performance

Colorgorical is implemented in C-accelerated Python. To evaluate average model runtime (50 runs) as a function of palette size (1 to 20 colors), we profiled single-palette generation on a Mid 2012 MacBook Pro Retina with a 2.6 GHz Intel Core i7 CPU and 16GB 1600MHz DDR3 RAM. Average initialization time was 140ms (SEM = 0.004). If a palette was returned before reaching the required number of colors, it was discarded and the test was run again. Runtime performance increased linearly in the number of colors such that adding a color increased runtime by 17.6ms on average (Supp. Mat.).

## 5 PALETTE SCORE EVALUATION

Before conducting human-subject testing, we first tested whether any of Colorgorical's scoring functions (i.e., Perceptual Distance, Name Difference, Name Uniqueness, and Pair Preference) could be removed from the model to simplify its design without significantly affecting palette output. For instance, if Perceptual Distance were to explain most of the variance in Name Difference scores, then the Name Difference scoring function could be removed from the model with little effect on palette output.

We examined the similarity among the four Colorgorical scoring functions using multiple linear regressions to predict 39,600 *palette scores* for each *palette scoring function* (Sec. 3) from the three remaining functions (e.g., predicting Perceptual Distance from Name Difference, Name Uniqueness, and Pair Preference). *Palette scores* are the minimum palette scoring function output given all color pairs in a

palette. We use the minimum score because we assumed that a palette is only as preferable or discriminable as its lowest pair. The number 39,600 stems from the full range of possible Colorgorical slider settings and 3 palette sizes (66 settings,  $\{3,5,8\}$ -colors, 200 repeats). The 66 settings were made from the different unique combinations from dragging each of the four sliders to 0%, 50%, or 100%, which ignore duplicate settings encountered when moving one or more sliders to 0%.

We also examined how the four palette scores changed with palette size. Below we highlight results and implications from our analyses, and the methods and full results are in Supplementary Material.

Both Perceptual Distance and Name Difference were strong positive predictors of one another. Name Uniqueness was always a weak negative predictor of the other scores. Pair Preference was always a strong negative predictor of Perceptual Distance and Name Difference. Further, Pair Preference was more strongly related to Name Difference than to Perceptual Distance. Given that palette scores in each palette scoring function were significantly predicted by all three of the other palette scoring functions, we concluded that each scoring function measured sufficiently different color information to justify keeping them all in the model for Experiment 1.

## 6 EXP. 1: MODEL HUMAN-SUBJECT EVALUATION

Experiment 1 tested how palette discriminability performance and preference ratings varied as the relative weights on the Colorgorical sliders varied (i.e., the relative importance of each scoring function; Sec. 4). We also identified which slider settings produced the most discriminable or preferable palettes to prepare for a comparison between Colorgorical and current industry standards in Experiment 2 (Sec. 7).

Experiment 1 used the same representative palettes as in Section 5, which were analogous to the slider settings produced by moving each to either 0%, 50%, or 100% for 3-, 5-, and 8-color palettes.

Discrimination performance and preference were assessed using two difference tasks (Fig. 4). In the discrimination task, participants reported which side of a map had more counties of a target color, providing data on number of errors and response time (RT). In the aesthetic preference task, participants rated how much they liked the color combinations in each palette. We predicted that:

- P1** Palettes with fewer color would be more discriminable
- P2** Discrimination RT and error would correlate in a strong negative direction with Perceptual Distance and in a strong positive direction with Pair Preference, whereas preference ratings would show the opposite pattern
- P3** Palette size would modulate the discriminability and preference ratings associated with each slider setting.
- P4** Slider settings would significantly predict discrimination performance and preference ratings

**P1** is based on previous evidence that visualizations with more colors are harder to process [5]. **P2** extends Palette Score Evaluation findings that Perceptual Distance and Name Difference negatively predicted Pair Preference. **P3** builds on the first two predictions: based on **P1** we expect that palette size will modulate the discriminability of slider settings, and based on **P2** we expect that preference will be negatively correlated with discriminability. **P4** makes two strings of assumptions based on the Palette Score Evaluation: (1) the trade-off between discrimination and preference palette scores will extend to behavior (**P2**) and (2) the relative importance of scoring functions (i.e., slider settings) would affect behavior in the same manner as palette scores (e.g., a higher relative importance of Pair Preference will produce higher Pair Preference palette scores). By transitivity, we predict that slider settings will be indicative of behavior.

### 6.1 Methods

#### 6.1.1 Participants

77 participants completed the discrimination task and 60 completed the preference rating task (recruited through Amazon Mechanical

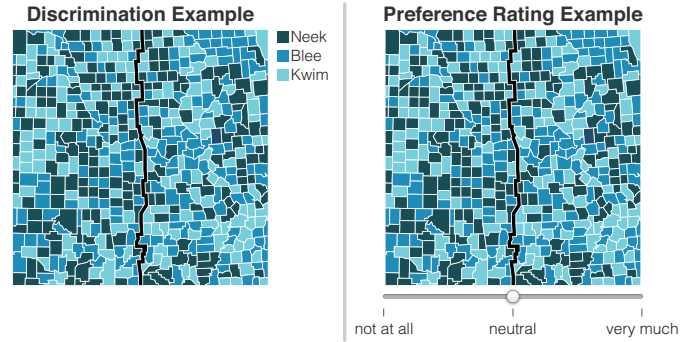


Fig. 4. Discrimination and preference rating task stimuli. The discrimination task asked users which side had more “Neek” counties (← and → keys). The preference rating task asked users to click on the slider.

Turk, \$3 compensation). Palette size (3-, 5-, 8-colors) was a between-subjects factor. For quality control, we determined *a priori* to discard participants who were < 60% accurate across all trials in the discriminability task (3-color:  $n = 3$ ; 5-color:  $n = 6$ ; 8-color:  $n = 8$ ). No participants were discarded in the preference task. In the final datasets there were 20 participants per palette size in each task, and discard frequency did not significantly differ between palette size conditions ( $\chi^2(2) = 1.793, p = 0.408$ ). All self-reported having normal color vision and gave informed consent. The Brown University IRB approved the experiment protocol.

#### 6.1.2 Design & Displays

The experimental designs for the discrimination and preference tasks were similar. In both, each participant saw 660 palettes from 66 slider settings (see Sec. 5 for setting information) with 10 different color palettes within each slider setting (treated as repetitions). The specific colors in each palette varied across participants (simulating different runs of Colorgorical), but were generated with the same experimental design. Palette size varied between-subjects (3, 5, or 8 colors).

The palettes that comprised the displays for the discrimination task were also used for the preference task, such that each discrimination participant was yoked to a preference participant (i.e., both saw the same palettes). Palettes were displayed on a predefined map of 554 counties in the U.S. ( $300 \times 300$  pixels). The map itself differed slightly based on the task (Fig. 4).

For the discrimination task, a 5-pixel-wide contour bisected the map (adhering to county borders). The contour was black and the county borders were white so that both would fall outside of Colorgorical’s default lightness sampling range ( $L \in [25, 85]$ ;  $L_{\text{Black}} = 0$ ;  $L_{\text{White}} = 100$ ). The size of the counties on each side were slightly altered so they were approximately equal (left: 165 px; right: 163 px). A legend rendered to the right of each map assigned each palette color to a nonsense word category. The target “Neek” color was always at the top of the legend to prevent participants from having to search for the target color. One side of the map had over-represented target color (“Neek”;  $1.5 \times$  more frequent on one side than the base rate) and the opposite side had an over-represented distractor color ( $1.3 \times$  more frequent). The target side was left/right balanced across trials. Based on our assumption that a palette is only as effective as its least discriminable pair of colors, the target and distractor colors were always the palette colors with the lowest and second-lowest Perceptual Distance scores compared with all other colors in the palette, respectively.

In the preference task, there was no dividing contour and no legend, the colors were roughly equal in proportion, and they were randomly assigned to positions across the map (no left/right asymmetry). Below the map there was a 300-pixel-wide continuous response slider scale ranging from -100 to 100 with labeled extrema and midpoint (left: “not at all”; right: “very much”; midpoint: “neutral”) [30]. The scale was initialized with the slider set to “neutral” to avoid biasing participants.

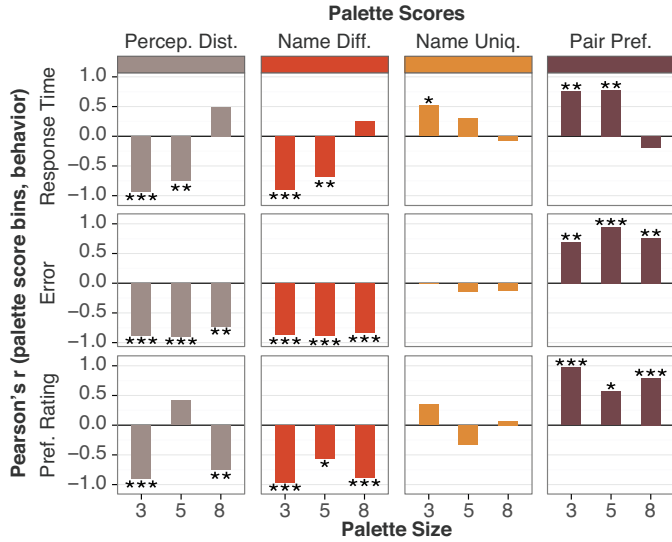


Fig. 5. Correlations between binned palette scores and responses on each measure for each palette size ( $p < 0.05$ ;  $p < 0.01$ ;  $p < 0.001$ ).

### 6.1.3 Procedure

**Discrimination Task.** Participants were first presented with an example display and were told that their task would be to indicate which half (left/right) of a map had more Neek counties using the left/right arrow keys. They were also told that the target Neek color would always be shown at the top of a legend and that answers would be marked incorrect if they did not respond within 3.5 seconds. Participants completed five practice trials using distinct displays from the 660 test maps they would see in the experiment, followed by 660 test trials. Maps were shown in random order in the center of the window. Trials were separated by a 500-ms inter-trial interval with a fixation cross displayed at the center of the screen. Optional breaks were given every 20 trials. This task took ~30 minutes to complete.

**Preference Task.** Participants were asked to rate their aesthetic preference for the color combination in each palette by clicking a point on a slider between the left (“not at all” preferable) and the right (“very much” preferable) ends (Fig. 4). To help them gauge what liking “not at all” and “very much” meant to them in the context of these color combinations, participants were shown an anchoring page containing 66 representative maps. They scrolled through the maps and considered how they would rate each map while using the full range of the scale. During the experiment, each map was presented one at a time in a random order (separated by a 250-ms blank pause screen). The preference slider appeared 1 second after the map appeared to encourage participants to consider their preference carefully before responding. This task took ~40 minutes to complete.

## 6.2 Results and Discussion

Before analysis, we pruned response time (RT) data by removing incorrect trials and then eliminating trials for each subject that were more than  $\pm 2.5$  standard deviations away from their mean RT [27]. On average, 129 errors (19.5%) and 24 outliers (3.6%) were removed.

Overall participant accuracy decreased as palette size increased (3-color average error: 79/660; 5-color: 119/660; 8-color: 190/660), indicating that displays with smaller color palettes were more discriminable (**P1, P3**). This result mirrored the increased participant discard rate for larger palette size conditions due to high error rates (Sec. 6.1.1) and is consistent with previous findings that showed visualizations with fewer color categories are more effective [5]. Between-subjects one-way ANOVAs testing for effects of palette size (3, 5, 8) within each measure indicated significant effects for number of errors ( $F(2,57) = 30.801, p < .001$ ) but not for RT or preference ( $F(2,57) = 1.574, 1.035; p = 0.216, 0.362$ , respectively).

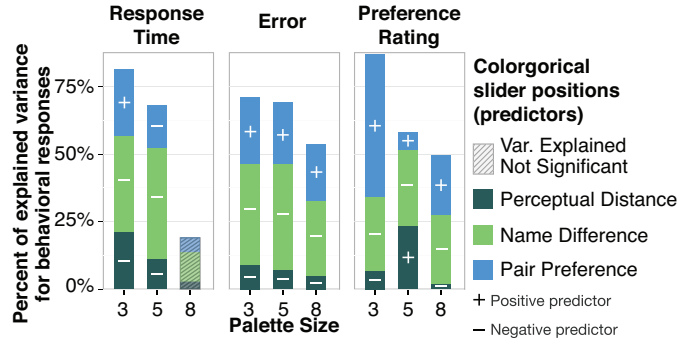


Fig. 6. Variance explained ( $R^2$ ) for the 9 models decomposed to look at the variance explained of behavioral data in terms of slider settings (i.e., palette score relative importance).

### 6.2.1 Palette score and behavioral measure correlations

Figure 5 shows the correlations between each type of palette score (Perceptual Distance, Name Difference, Name Uniqueness, and Pair Preference) and the three behavioral measures (RT, error rate, and preference ratings), averaged over participants. *Palette score* refers to the lowest *palette scoring function* value (Sec. 3) given all color pairs in a palette. To conduct these analyses, we first binned the behavioral data for each measure according to palette scores (15 equally-spaced bins) for each subject<sup>2</sup>. After, we averaged the data for all palettes that scored in the same bin and then averaged those values across participants. This binning was necessary prior to averaging across participants because each participant saw different palettes with slightly different scores (Supp. Mat.). For example, RT for palettes with a Pair Preference scores of 30.03 and 30.05 would be binned together.

We cross-checked the binned-score correlations by calculating the within-subject correlations for each behavioral measure and palette score and then used Fisher’s Z transform prior to calculating the between-subject average Pearson’s  $r$  for each measure and score combination. For the most part, these analyses showed the same pattern of results as the binned correlation statistics (Supp. Mat.).

The binned-score correlations are presented below (see Supp. Mat. for individual correlations on non-binned data). In summary, the Perceptual Distance, Name Difference, and Pair Preference scores had the predicted effects: RT and error rates decreased (i.e., better performance) as Perceptual Distance and Name Difference increased, but they increased (i.e., worse performance) as Pair Preference increased (**P2**). In contrast, preference decreased as Perceptual Distance and Name Difference increased and they increased as Pair Preference increased. Name Uniqueness had little to no effect.

**RT.** RT decreased as Perceptual Distance and Name Difference increased for 3- and 5-color palettes (Perceptual Distance:  $r(13) = -0.926, r(13) = -0.757; p \leq 0.001$  respectively; Name Difference:  $r(13) = -0.893, r(13) = -0.689; p \leq 0.005$  respectively). Similarly, Pair Preference followed **P2** for 3- and 5-color palettes with strong positive correlations with RT ( $r(13) = 0.755, 0.776; p = 0.001$ ). Name Uniqueness was significantly correlated with RT for 3-colors ( $r(13) = 0.521; p = 0.046$ ), but not for 5-colors ( $r(13) = 0.306; p = 0.268$ ). No scores were significantly correlated with RT for 8-color conditions ( $r(12) = 0.496, r(13) = 0.251, -0.071, -0.193; p \geq 0.071$  for Perceptual Distance, Name Difference, Name Uniqueness, and Pair Preference respectively). These findings largely support **P2** for 3 and 5 colors; however, 8-color palette correlations were not significant.

**Error.** Error rate correlations were significant for all sizes with Perceptual Distance ( $r(13) = -0.887, -0.898, r(12) = -0.731; p \leq 0.003$ , for 3-, 5-, and 8-colors respectively), Name Difference ( $r(13) = -0.874, -0.892, -0.838; p < 0.001$ ), and Pair Preference ( $r(13) = 0.697, 0.945, 0.761; p \leq 0.004$ ). Similar to RT correlations, Name

<sup>2</sup>The degrees of freedom for 8-color perceptual distance correlations is one less due to an empty bin, which is shown in the Supplementary Material.



Uniqueness was not significantly related to error measures ( $r(13) = -0.016, -0.141, -0.126; p \geq 0.616$ ).

**Preference Rating.** Preference rating trends were the opposite of error and RT, and consistent with **P2**. Increasing 3- and 8-color Perceptual Distance reduced preference ratings, ( $r(13) = -0.897, r(12) = -0.751; p \leq 0.002$ ) but not significantly so for 5-color palette ( $r(13) = 0.412; p = 0.127$ ). Preference ratings also decreased as Name Difference increased for 3-, 5-, and 8-colors ( $r(13) = -0.969, -0.57, -0.891; p \leq 0.026$ , respectively). Increasing Pair Preference increased preference ratings ( $r(13) = 0.971, 0.57, 0.796; p \leq 0.026$ ). Again, Name Uniqueness was not significantly related ( $r(13) = 0.346, -0.333, 0.073; p \geq 0.207$ ).

## 6.2.2 Predicting behavioral measures from slider settings

To test whether slider settings (i.e., relative importance of the *palette scoring functions*) significantly predict behavior (**P4**), we performed a series of multiple linear regressions that predicted behavioral measures as a function of changing sliders to 0%, 50%, or 100% (Fig. 6). Given that the correlational analyses above suggested that Name Uniqueness had little effect on behavior, we averaged slider configurations that would be equivalent if Name Uniqueness were ignored. For example, if Perceptual Distance and Name Uniqueness were both set to 50%, the new setting would be Perceptual Distance as 100% and would be averaged with other palettes where Perceptual Distance is 100%. This reduced the regression analysis to predict 20 unique slider settings rather than the previous 66. The data that were input to the correlations are graphed in the Supplementary Material.

Below we detail the results of the multiple linear regressions using slider relative importance to predict the three behavioral measures. More information about the relation between sliders, size, and behavioral measures is provided in the Supplementary Material. In summary, the slider settings were typically able to significantly predict the behavioral measures (**P4**).

**RT.** RT decreased (improved) as Perceptual Distance and Name Difference slider weights increased and RT increased (got worse) as Pair Preference slider weights increased (**P4**; Supp. Mat.). Name Difference was always the most predictive and Perceptual Distance and Pair Preference were similarly less predictive (Fig. 6). Although this pattern was present for all three palette sizes, the models were significant for the 3- and 5-color palettes ( $F(3, 16) = 23.442, 11.447; R^2 = 0.815, 0.682; p < 0.001$ , respectively), but not the 8-color palettes ( $F(3, 16) = 1.267, R^2 = 0.192, p = 0.319$ ). The lack of significance for 8-color palettes coincides with the oddity that response time was typically faster for 8-color palettes than 5-color ones; this is unexpected, given (1) past visual search research finding that more colors take longer to discriminate [5] and (2) the previously discussed palette size relation with accuracy. We suspect that this difference may be because participants tried less hard or the task became too difficult in the 8-color condition because they had higher overall error rates. Another possibility is that pair-based color discriminability scores (e.g., Perceptual Distance) may break down as the number of colors increases, which would create a need for higher-order combination discriminability scores. Each of these possibilities raise interesting future directions for studying the relation between palette effectiveness and number of colors.

**Error.** Slider relative importance analysis mirrored RT (**P4**; Supp. Mat.), except that Pair Preference was more important than Perceptual Distance (Fig. 6 and Supp. Mat.). The reason for this difference is unknown. The multiple linear regressions for all 3-, 5-, and 8-colors were all significant ( $F(3, 16) = 13.186, 11.964, 6.192; R^2 = 0.712, 0.692, 0.537; p \leq 0.005$ , respectively).

**Preference Rating.** Preference ratings increased with weights on the Pair Preference slider and decreased with weights on the Perceptual Distance and Name Difference sliders (**P4**; Supp. Mat. slider-behavior figure). Pair Preference was the most predictive slider for 3-colors, but not for 5- and 8-colors (Fig. 6 and Supp. Mat.); instead, Name Difference was most predictive. Perceptual Distance was more important than Pair Preference for 5-colors, but was otherwise the

least important slider. The multiple linear regressions for 3-, 5-, and 8-colors were all significant ( $F(3, 16) = 35.089, 7.396, 5.228; R^2 = 0.868, 0.581, 0.495; p \leq 0.01$ , respectively).

It is noteworthy that the model's ability to predict preference ratings decreased for 5-colors relative to 3-colors, suggesting that the mechanism behind human aesthetic preference ratings may deviate from pair-based preference predictions as the number of palette colors changes. Another difference for 5-color palettes, compared to 3- and 8-colors, was that all settings were rated either neutral or slightly negative. These results suggest that the assumption that pair-wise based preference models generalize to palettes of three colors might break down for larger palettes. The differences in preference ratings over palette sizes motivates the need for further research on the aesthetics of higher-order color combinations.

We also found that preference ratings decreased faster as Name Difference relative importance was increased compared to increases in Perceptual Distance relative importance (see Supp. Mat.). This asymmetry might be caused by differences in how Perceptual Distance and Name Difference measure distances in color space. It could be that Perceptual Distance is more supportive because it can generate color pairs that differ primarily in lightness (which is one of the terms in Pair Preference), whereas Name Difference might be more likely to favor differences in hue, which would be in opposition to Pair Preference's hue similarity term.

## 6.2.3 Lowest-Error and Highest-Preference settings

A main goal of Experiment 1 was to determine which Colorgorical settings to use to generate color palettes for comparison against current standards (Experiment 2). The combinatorial explosion of conditions prevented comparing all slider combinations to current standards. Therefore, we chose to select slider settings that either produced highly discriminable or highly preferable palettes (i.e., at either end of the previously-discussed discriminability-preference trade-off). Figure 8 shows the lowest discrimination error setting (subsequently called "Low-Error" palettes) and the highest preference rating setting ("Preferable" palettes) for each palette size. There were significantly fewer errors for Low-Error palettes than for Preferable palettes ( $t(19) = 3.322, 7.589, 3.15; p \leq 0.005$ , 3-, 5-, 8-colors). Preference ratings were significantly greater for Preferable palettes than for Low-Error palettes for 3- and 8-colors ( $t(19) = 4.610, 2.841, p \leq 0.01$ ), but not for 5-colors ( $t(19) = 0.499, p = 0.623$ ) (consistent issues about 5-color palettes discussed above).

## 6.2.4 Summary

Experiment 1's results largely support each of our four predictions and suggest that Colorgorical's sliders are effective at controlling the discriminability and preference of color palettes, although some 5- and 8-color conditions led to unexpected behavioral results. Discriminability performance typically improved (faster RT, fewer errors) as the Perceptual Distance and Name Difference palette scores increased (and with greater weight on their corresponding sliders) and Preference judgments typically increased as Pair Preference palette scores increased (with greater weight on its slider) (**P2**). There was also evidence for a tradeoff – discriminability decreased as both Pair Preference scores and scoring function weights increased, and preference judgments decreased as Perceptual Distance and Name Difference increased. This finding supports our earlier claim that care must be taken to design palettes that balance both discriminability and aesthetic preference. We also found that Name Difference, not Perceptual Distance, might better predict discriminability. This would also support Demiralp et al.'s previous findings that suggested Name Difference is a better measure of color distance than Perceptual Distance [3].

Additionally, our results suggest that smaller palettes are more discriminable (**P1**), that palette size modulates discriminability and preference ratings (**P3**), and that slider configurations significantly predict behavior (**P4**). We provide additional analysis and discussion for each prediction in the Supplementary Material.

Last, differences in discriminability and aesthetic preference trends

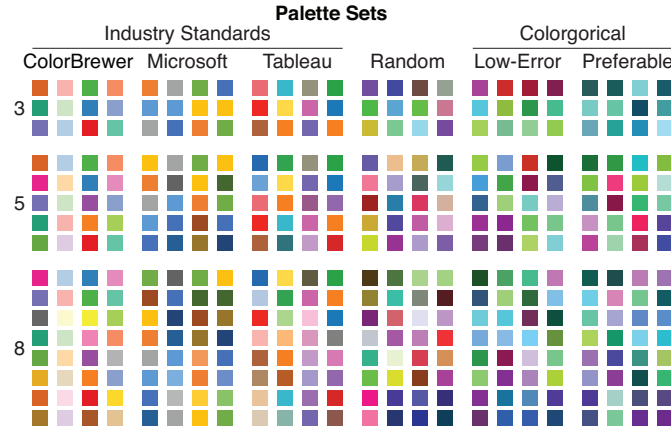


Fig. 7. Exp. 3 palettes: ColorBrewer (Dark2, Pastel1, Set1, Set2); Microsoft (all); Tableau (10, Blue-Red, Green-Orange, Purple-Gray); Colorgorical and Random palettes varied across participants.

across palette sizes motivate additional research beyond pairwise theoretical models of color discrimination and preference rating.

## 7 EXP. 2: COLORGORICAL-OTHERS BENCHMARK

Experiment 2 compares palettes generated by Colorgorical Low-Error and Preferable slider settings to commonly used “benchmark” palettes (ColorBrewer, Microsoft Excel, and Tableau; Fig. 7). We also included randomly sampled palettes with noticeably different colors to simulate palettes made by someone without design expertise who tried to choose colors that were not confusable. We predicted that:

- P Colorgorical Low-Error and Preferable settings would produce palettes that are at least as discriminable and typically more preferable compared to the majority of benchmarks

We based this prediction on expected outcomes of Colorgorical and benchmark palettes by applying regressions modeled on Experiment 1 palette scores and behavioral responses to the palette scores of Experiment 2 palette sets. As shown in Figure 8, Colorgorical palettes were expected to create more preferable palettes, with the exception of Microsoft 5- and 8-color palettes, which were predicted to outperform both Colorgorical settings. We also expected that Colorgorical would produce palettes with error rates similar to Tableau across all three sizes. We specified planned comparisons to test these predictions with the human-subject data from Experiment 1.

### 7.1 Methods

#### 7.1.1 Participants

75 participants (recruited through Amazon Mechanical Turk; paid \$1) completed the discrimination task and 60 completed the preference task. All gave informed consent, and the Brown University IRB approved the experiment protocol. All self-reported having normal color vision. 15 discrimination participants were less than 60% accurate and were discarded, per Experiment 1 procedure (3-colors:  $n = 0$ , 5-colors:  $n = 7$ , 8-colors:  $n = 8$ ). Participants were divided equally across size conditions ( $n = 20$  per size), and there was a significant effect between discard rate and size ( $\chi^2(2) = 6.878, p = 0.032$ ).

#### 7.1.2 Design, Displays, & Procedure

Palette size (3,5,8) varied between subjects and the rest of the factors varied within-subject. Participants in the discriminability task completed 96 trials (6 palette sets {Colorgorical Low-Error and Preferable, ColorBrewer, Microsoft, Tableau, Random}  $\times$  4 palettes taken from each set  $\times$  4 repetitions). Participants in the preference rating task were presented with 24 trials (6 palette sets  $\times$  4 palettes, no repetition).

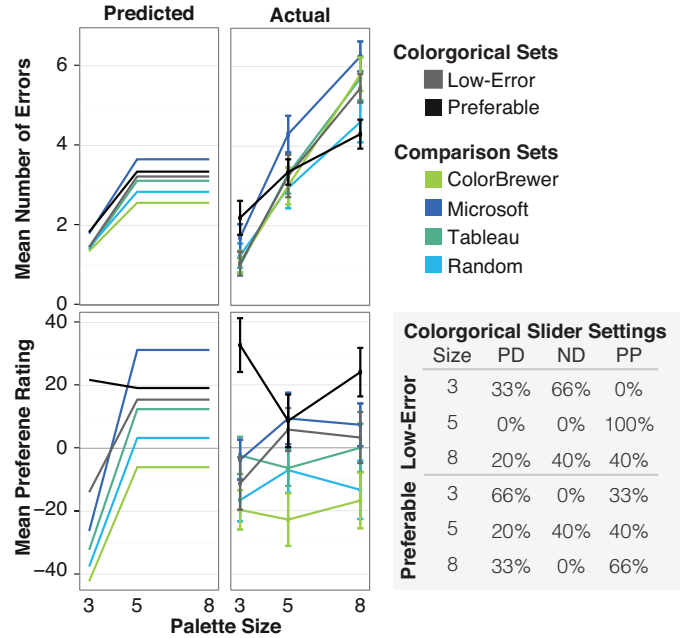


Fig. 8. Palette-behavior predictions and actual results for Experiment 2 (e.g., 4 errors = 25% error rate). Prediction models were trained on Experiment 1 palette scores and behavioral responses. Error bars show SEM. The table shows the Exp.2 Colorgorical slider settings (PD: Perceptual Distance; ND: Name Difference; PP: Pair Preference).

The benchmark palette sets included four palettes from each palette group’s larger collection (Fig. 7). Microsoft palettes included all four available palettes in Microsoft Excel for Mac (v.15.8). ColorBrewer palettes included four of the eight available palettes, including those with the greatest minimum Perceptual Distance and excluding palettes with niche purposes (e.g., “Paired”) [6]. Tableau palettes included the default Tableau 10 and the three palettes that were not designed for niche applications. We created random palettes by randomly sampling discriminable colors in RGB space for each participant (Sec. 4.3.1). All participants saw the same benchmark palettes aside from random. Each participant was given different random and Colorgorical palettes to test each palette type’s full potential variance. The Low-Error and Preferable palettes were made with settings described at the end of Section 6 (Fig. 8). Otherwise, the design, stimuli, and procedure were the same as Experiment 1. The discrimination task took ~5 minutes to complete and the preference rating task took ~10 minutes to complete.

### 7.2 Results and Discussion

We focused only on error and preference rating data (not RT) because error and RT results in Experiment 1 were similar and because we chose the Colorgorical palettes based on error rates and preference ratings. All reported  $t$ -tests were paired sample and two-tailed.

We first conducted two 6 palette set (within-subject)  $\times$  3 palette size (between-subject) mixed-design ANOVAs: one for error rates (averaged over replications) and a second for preference ratings. For error, there were main effects of palette set ( $F(5, 285) = 3.538, p = 0.004$ ), palette size ( $F(2, 57) = 59.34, p < 0.001$ ), and a 2-way interaction between them ( $F(10, 285) = 5.896, p = 0.01$ ). For preference ratings, there was a main effect of palette set ( $F(5, 285) = 13.235, p < 0.001$ ) with no effect of palette size ( $F(2, 57) = 0.258, p = 0.773$ ) and no interaction ( $F(10, 285) = 1.283, p = 0.239$ ). As shown in Figure 8, error increased with size, but preference ratings were more stable as size increased. Palette set differences are shown through the vertical separation of behavioral responses across palette sets. Our planned comparisons below delve into these effects, and they largely support the trends in our predictive models based on palette score with (although size does not show the predicted effect for preference ratings).



### 7.2.1 Colorgorical Low-Error vs. Preferable Palettes

We first tested whether the error and preference differences between Colorgorical-Low-Error and -Preferable palettes replicated the results of Experiment 1. As in Experiment 1, the Preferable palettes were preferred to the Low-Error palettes for 3- and 8-color palettes ( $t(19) = 3.573, -3.79; p = 0.002, 0.001$ ); but not for 5-color palettes ( $t(19) = -0.405, p = 0.690$ ). There were fewer errors for the 3-color Low-Error palettes than for the Preferable palettes ( $t(19) = 3.286, p = 0.004$ ), but there was no difference for the 5-color palettes ( $t(19) = 0.195, p = 0.847$ ). The only test that was inconsistent with our previous findings was that error rates for 8-colors were lower for Preferable palettes than for Low-Error palettes ( $t(19) = 2.113, p = 0.048$ ). The reason for this result is unknown.

### 7.2.2 Comparing Colorgorical to industry standard palettes

We next tested our prediction that Colorgorical palettes would be as discriminable and typically more preferable than the benchmark palettes. The tests were planned *a priori* based on predictions from Colorgorical and benchmark palette scores described below (Fig. 8). We conducted 48 paired two-sample *t*-tests comparing participants' discrimination error and preference ratings within the Colorgorical palettes and between the Colorgorical palettes and the four benchmark palette sets within each palette size (Fig. 8).

**Error rate.** Based on the model predictions (Fig. 8), we expected that error would not significantly differ between Colorgorical Low-Error palettes and all benchmarks except for Microsoft, where we predicted that Low Error palettes would elicit fewer errors. For 5- and 8-colors we predicted that Low-Error errors would be similar to Tableau, worse than ColorBrewer and Random, and better than Microsoft. We made the same predictions for Preferable palettes, except that 3-color error might only be as good as Microsoft, and 5- and 8-color error might be worse than Tableau.

Performance for Low-Error palettes was slightly better than expected. There were significantly fewer errors for 5-color Low Error than for 5-color Microsoft ( $t(19) = 2.396, p = 0.027$ ) and no significant difference from the other benchmarks ( $t(19) < 1.628, p \geq 0.12$ ).

Colorgorical-Preferable error also matched our predictions because there was always at least one benchmark that had non-significantly different error rates compared to the setting ( $t(19) \leq 1.898, p \geq 0.073$ ). Unexpectedly, Colorgorical-Preferable palettes led to significantly lower error than 8-color ColorBrewer and Microsoft palettes ( $t(19) \geq 2.910, p \leq 0.009$ ). However, consistent with our predictions, 3-color Colorgorical-Preferable led to significantly more errors than ColorBrewer, Tableau, and Random benchmarks ( $t(19) = 2.531, 3.644, 3.047; p = 0.020, 0.002, 0.007$ , respectively).

The fewer errors for random than for Colorgorical preferable may be surprising, but it is consistent with our earlier observations. There is a high likelihood that three randomly sampled colors will be far apart in our quantized CIELAB space, leading to very high discriminability but also low preference. As the number of randomly sampled colors increases, discriminability decreases, as shown in the non-significant comparisons to 5- and 8-color Colorgorical Preferable. Although the Colorgorical-Preferable settings produced less discriminable results in some conditions (e.g., 3-color error), there was always at least one benchmark that lacked significantly different error rates.

**Preference ratings.** We predicted that both Low-Error and Preferable palettes would be more preferable in all comparisons except to 5- and 8-color Microsoft palettes (Fig. 8).

Low-Error was significantly more preferred than 5-color ColorBrewer and 8-color Random ( $t(19) = 2.784, 2.279, p = 0.012, 0.034$ , respectively) and was never significantly less preferred than the other benchmarks ( $t(19) \leq 1.781, p \geq 0.091$ ). Colorgorical-Preferable palettes often led to significantly more preferable palettes (8 of 12, all but 5- and 8-color Microsoft, 5-color Random, and 8-color Tableau;  $t(19) \geq 2.105, p < 0.05$ ).

**Summary.** Colorgorical Low-Error and Preferable palettes are almost always as discriminable and often more preferable than the current standard visualization-specific categorical color palettes (P).

Low-Error palettes were sometimes more discriminable and more preferable or otherwise not significantly different than the benchmark palettes. Similarly, Preferable palettes often led to significantly higher preference ratings, and discriminability was not significantly different compared to at least one industry standard for all sizes. Thus, Colorgorical allows users without design expertise to create discriminable and preferable palettes that often do not have significantly different discriminability and that sometimes are more preferable than current pre-made standards.

## 8 OPEN RESEARCH AREAS

We found that Colorgorical palettes, based on models of aesthetics and discriminability, can be as effective as expert-made visualization palettes and even more aesthetically preferable. These findings lead to several future research directions. First, given that color combination discriminability and preference can be inversely related, how can discriminability and preference be automatically optimized? Second, what alternatives to the current pairwise theoretical models might better predict discriminability and aesthetic preference for higher-order combinations (e.g., 5- or 8-colors)? Third, how would color preference models that diverge from figure/ground preference alter palette construction? For instance, how might Lin et al.'s preferable palette generation technique that learns from artist-generated training palettes [12] compare to palettes made with Pair Preference? Fourth, would the same results hold if hue filters are applied when constructing Colorgorical palettes? Fifth, how might Colorgorical help designers foresee palettes that might be indiscriminable given color deficiencies [28]?

## 9 CONCLUSION

We presented Colorgorical, a model-driven approach to generating categorical color palettes for information visualizations by configuring palette discriminability and preference. Colorgorical uses an iterative, semi-random-sampling procedure to generate palettes of a specified size. User-defined configurations work by changing the relative importance of Perceptual Distance, Name Difference, and Pair Preference scoring functions. Users can further customize palette creation by modifying the number of colors, by defining which hues to sample from, and by providing an existing palette to build upon.

The novelty of our approach stems from our departure from previous palette creation strategies. Whereas previous palette creation tools focused primarily on discriminability or favored color relations in harmonic templates whose empirical validity is questionable (e.g., Adobe Color [21]), Colorgorical generates palettes with user-defined relative importances for discriminability and aesthetic preference (Sec. 3). Our color sampling approach also differs in strategy from pre-made palette sets such as ColorBrewer, in which categorical palettes are generated by first choosing colors representing different names and then varying each palette color's value [1].

Empirical tests show that each of Colorgorical's sliders, which are used to balance palette discrimination and preference, measure different aspects of color (Sec. 5) and modulate behavior as they were designed to do (e.g., weighting discriminability sliders increases discriminability performance) (Sec. 6).

Empirical tests that compare Colorgorical palettes and industry standards revealed that our model-derived palettes are as effective as, and sometimes better than, current categorical color palette standards. Our findings also indicate that the number of colors may alter the effectiveness of pair-based discriminability and preference scores. Colorgorical also improves upon industry standards by giving users the flexibility to create their own discriminable and preferable palettes while enforcing visualization design constraints. These results indicate that Colorgorical provides an effective way to create categorical visualization color palettes. Colorgorical is open-sourced at <http://vrl.cs.brown.edu/color>.

## ACKNOWLEDGMENTS

This material is based upon work supported by a National Science Foundation Graduate Research Fellowship under Grant No. DGE-1058262 and Brown University's Center for Vision Research.

## REFERENCES

- [1] C. A. Brewer, G. W. Hatchard, and M. A. Harrower. Colorbrewer in print: A catalog of color schemes for maps. *Cartography and Geographic Information Science*, 30(1):5–32, 2003.
- [2] M. E. Chevreul. *The principles of harmony and contrast of colours, and their applications to the arts*. Van Nostrand Reinhold, New York, NY, USA, 1987 (1839).
- [3] Çağatay. Demiralp, M. Bernstein, and J. Heer. Learning perceptual kernels for visualization design. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):1933–1942, Dec 2014.
- [4] C. C. Gramazio, K. B. Schloss, and D. H. Laidlaw. The relation between visualization size, grouping, and user performance. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):1953–1962, Dec 2014.
- [5] S. Haroz and D. Whitney. How capacity limits of attention influence information visualization effectiveness. *IEEE Transactions on Visualization and Computer Graphics*, 18(12):2402–2410, Dec 2012.
- [6] M. Harrower and C. A. Brewer. Colorbrewer.org: an online tool for selecting colour schemes for maps. *The Cartographic Journal*, 40(1):27–37, 2003.
- [7] C. G. Healey. Choosing effective colours for data visualization. In *Visualization*, pages 263–270, Oct 1996.
- [8] J. Heer and M. Stone. Color naming models for color selection, image editing and palette design. In *ACM Human Factors in Computing Systems (CHI)*, pages 1007–1016, New York, NY, USA, 2012. ACM.
- [9] G. Hu, Z. Pan, M. Zhang, D. Chen, W. Yang, and J. Chen. An interactive method for generating harmonious color schemes. *Color Research and Application*, 39(1):70–78, 2014.
- [10] J. Itten. *The art of color*. Van Nostrand Reinhold, New York, NY, USA, 1961.
- [11] S. Lin, J. Fortuna, C. Kulkarni, M. Stone, and J. Heer. Selecting semantically-resonant colors for data visualization. *Computer Graphics Forum*, 32(3pt4):401–410, 2013.
- [12] S. Lin, D. Ritchie, M. Fisher, and P. Hanrahan. Probabilistic color-by-numbers: Suggesting pattern colorizations using factor graphs. *ACM Trans. Graph.*, 32(4):37:1–37:12, July 2013.
- [13] M. R. Luo, G. Cui, and C. Li. Uniform colour spaces based on ciecam02 colour appearance model. *Color Research & Application*, 31(4):320–330, 2006.
- [14] M. R. Luo, G. Cui, and B. Rigg. The development of the cie 2000 colour-difference formula: Ciede2000. *Color Research & Application*, 26(5):340–350, 2001.
- [15] Y. Matsuda. *Color Design*. Asakura Shoten, 1995.
- [16] B. A. Maxwell. Visualizing geographic classifications using color. *The Cartographic Journal*, 37(2):93–99, 2000.
- [17] B. J. Meier. Ace: A color expert system for user interface design. In *Proceedings of the ACM SIGGRAPH Symposium on User Interface Software*, UIST ’88, pages 117–128, New York, NY, USA, 1988.
- [18] B. J. Meier, A. M. Spalter, and D. B. Karelitz. Interactive color palette tools. *IEEE Computer Graphics and Applications*, 24(3):64–72, May 2004.
- [19] R. Munroe. Color survey results, May 2010.
- [20] A. H. Munsell. *A grammar of color*. Van Nostrand Reinhold, New York, NY, USA, 1969 (1921).
- [21] P. O’Donovan, A. Agarwala, and A. Hertzmann. Color compatibility from large datasets. *ACM Trans. Graph.*, 30(4):63:1–63:12, 2011.
- [22] W. Ostwald. *Colour Science (Vol. II)*. Winsor and Newton, Ltd., London, UK, 1933.
- [23] L.-C. Ou, P. Chong, M. R. Luo, and C. Minchew. Additivity of colour harmony. *Color Research & Application*, 36(5):355–372, 2011.
- [24] L.-C. Ou and M. R. Luo. A colour harmony model for two-colour combinations. *Color Research & Application*, 31(3):191–204, 2006.
- [25] S. E. Palmer and W. S. Griscom. Accounting for taste: Individual differences in preference for harmony. *Psychonomic Bulletin & Review*, 20(3):453–461, 2012.
- [26] S. E. Palmer and K. B. Schloss. An ecological valence theory of human color preference. *Proceedings of the National Academy of Sciences*, 107(19):8877–8882, 2010.
- [27] R. Ratcliff. Methods for dealing with reaction time outliers. *Psychological bulletin*, 114(3):510, 1993.
- [28] K. Reinecke, D. R. Flatla, and C. Brooks. Enabling designers to foresee which colors users cannot see. In *ACM Human Factors in Computing Systems (CHI)*, pages 2693–2704, New York, NY, USA, 2016. ACM.
- [29] K. B. Schloss, D. Hawthorne-Madell, and S. E. Palmer. Ecological influences on individual differences in color preference. *Attention, Perception, & Psychophysics*, 77(8):2803–2816, 2015.
- [30] K. B. Schloss and S. E. Palmer. Aesthetic response to color combinations: preference, harmony, and similarity. *Attention, Perception, & Psychophysics*, 73(2):551–571, 2011.
- [31] V. Setlur and M. Stone. A linguistic approach to categorical color assignment for data visualization. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):698–707, Jan 2016.
- [32] G. Sharma, W. Wu, and E. N. Dalal. The ciede2000 color-difference formula: Implementation notes, supplementary test data, and mathematical observations. *Color Research and Applications*, 30(1), 2005.
- [33] M. Stone, D. A. Szafr, and V. Setlur. An engineering model for color difference as a function of size. In *Proceedings of the Color and Imaging Conference*, pages 228–233, 2014.
- [34] F. Szabó, P. Bodrogi, and J. Schanda. Experimental modeling of colour harmony. *Color Research & Application*, 35(1):34–49, 2010.
- [35] C. Taylor and A. Franklin. The relationship between color–object associations and color preference: Further investigation of ecological valence theory. *Psychonomic Bulletin & Review*, 19(2):190–197, 2012.
- [36] K. Yokosawa, K. B. Schloss, M. Asano, and S. E. Palmer. Cross-cultural studies of color preferences: Us and japan. *Cognitive Science*, 2015.
- [37] A. Zeileis, K. Hornik, and P. Murrell. Escaping rgbland: Selecting colors for statistical graphics. *Computational Statistics and Data Analysis*, 53(9):3259 – 3270, 2009.