
M

Citation: Jameson, K. A., Benjamin, N. A., Chang, S.M., Deshpande, P. S., Gago, S., Harris, I. G., Jiao, Y., and Tauber, S. (2015). Mesoamerican Color Survey Digital Archive. In Encyclopedia of Color Science and Technology, (Ronnier Luo, Ed.). Springer Berlin Heidelberg. ISBN: 978-3-642-27851-8 (Online). DOI 10.1007/978-3-642-27851-8.

Mesoamerican Color Survey Digital Archive

Kimberly A. Jameson¹, Nathan A. Benjamin², Stephanie M. Chang², Prutha S. Deshpande³, Sergio Gago⁵, Ian G. Harris⁴, Yang Jiao⁴ and Sean Tauber¹

¹Institute for Mathematical Behavioral Sciences, University of California, Irvine, Irvine, CA, USA

²Calit2, Computer Science, University of California, Irvine, Irvine, CA, USA

³Cognitive Sciences, University of California, Irvine, CA, USA

⁴Computer Science, University of California, Irvine, CA, USA

⁵Calit2, School of Engineering, University of California, Irvine, Irvine, CA, USA

Definition

The Mesoamerican Color Survey (MCS) collected color-naming and categorization data from approximately 900 speakers from each of 116 indigenous languages from regions in Mesoamerica or Central America. Analyses of these data were originally reported by Dr. Robert E. MacLaury, principal investigator of the survey. The MCS data exist as a public-access color categorization and naming digital archive, in conjunction with other color categorization data collected by MacLaury, and are available at <http://colcat.calit2.uci.edu/>.

Introduction

Human categorization behavior is widely studied across the behavioral sciences. It underlies many cognitive functions, including concept formation, decision making, learning, and communication. Color appearance, similar to other natural categorization domains, has distinctive features or properties that vary along continuous dimensions. Semantic color categories, their formation, their best exemplars and boundaries, and the influence of these on human behavior have been the topics of much empirical study, receiving considerable attention from anthropologist, linguists, cognitive scientists, and psychologists. The general aim of such research is to understand human color categorization and how it is cognitively and culturally represented across languages. For the case of color categories, the literature suggests that there are “universal trends” in human color representations with the category best exemplars being predictable across languages [1–7], and there are culturally specific color categorization influences [8–15], as well as human evidence [16–22] and simulated color category evolution evidence [23–27], suggesting combinations of universal, cultural, and pragmatic influences on color categorization and naming behaviors.

The Mesoamerican Color Survey (MCS) is one of the two existing databases (the other being the World Color Survey or WCS; see “► [World Color Survey](#)”) which directly

investigated, on a large scale, color naming and categorization across many linguistic societies. The MCS and the WCS employ nearly identical standardized procedures for evaluating large numbers of color stimuli, languages, and informants. The Mesoamerican Color Survey was conducted by Robert E. MacLaury during the years 1978–1981. MCS data were collected by MacLaury himself, or by research associates and colleagues in Mesoamerica, whom MacLaury trained and directed. The MCS data represent interviews with 900 speakers of some 116 Mesoamerican languages. It is estimated that more than 100 indigenous languages are spoken in Mexico and Central America. Like most languages each MCS language has a color lexicon that partitions environmental color appearance stimuli according to a pattern that is specifically relevant to a given language’s speakers. But every local MCS system of color categorization also shares characteristics with systems observed for other Mesoamerican languages and with those of languages elsewhere in the world. MacLaury published analyses of the MCS data in the context of his *Vantage Theory* modeling approach ([28–30], see ► [Vantage Theory of Color](#)).

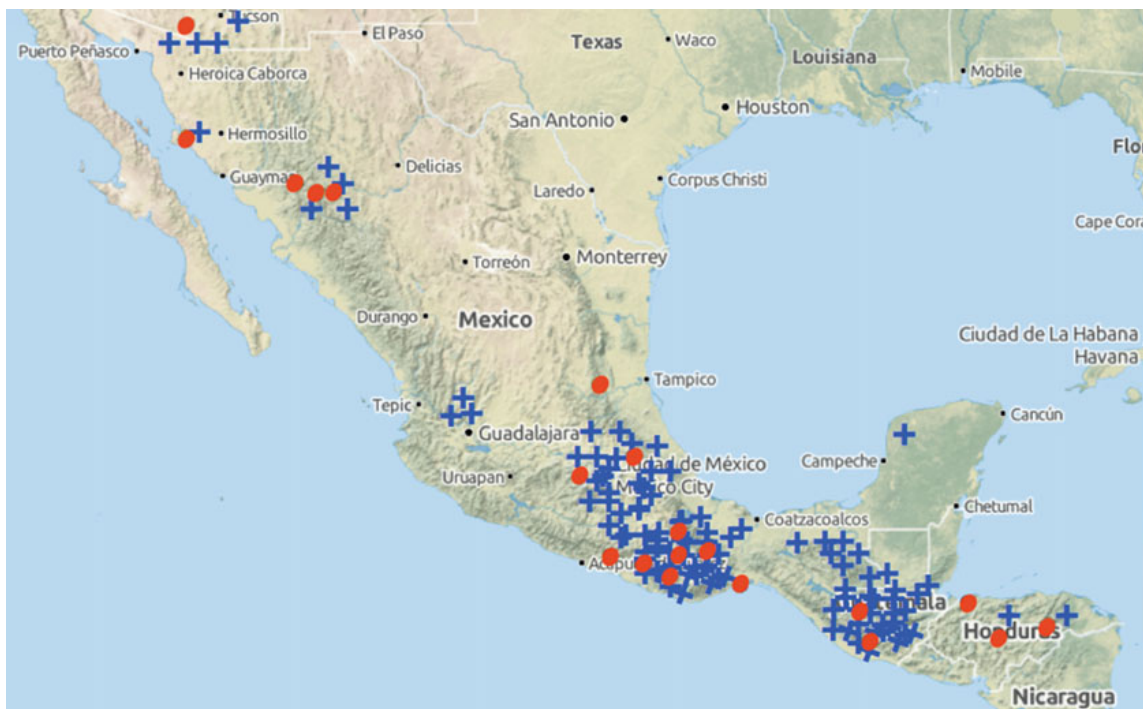
MCS and Theories of Color Categorization

Humans can discriminate on the order of 10^6 different colors, many more colors than any individual can name reliably. Most of the 7,000 living languages have some form of a color lexicon. Various theories in the literature have sought to account for empirical findings on how each language represents color concepts and how such representations evolve, changing over time, and successively form hierarchical linguistic representations of color appearance similarity. In 1969 Brent Berlin and Paul Kay were among the first to systematically empirically examine how different linguistic societies represented color and theorize how color lexicons might similarly evolve ([1], see “► [Berlin & Kay Theory](#)”). In the data they collected from 20 unrelated language families, they found (a) that there exist

universal crosslinguistic constraints on color naming and (b) that basic color terminology systems tend to develop in a partially fixed order.

Berlin and Kay made two conjectures about the evolution of color lexicons: (1) There is a limited set of basic color terms (abbreviated as BCTs) in most languages, which are distinct from other color terms that an individual might use to name colors. (2) Color lexicons evolve from simple to complex, along highly constrained paths, starting from two BCTs corresponding to warm-or-light and dark-or-cool categories in the most reduced lexicons and extending to 11 BCTs as seen in English and other languages spoken in industrialized societies. The Berlin and Kay theory is the basis for a mainstream approach to color categorization theory, along with its subsequent formulations using the World Color Survey (WCS) data ([2–4, 7], see “► [World Color Survey](#)”).

While much empirical support exists for the Berlin and Kay/WCS theoretical approach, alternative explanations have also been suggested. Researchers directly involved in the WCS project, who collaborated with Berlin and Kay, have suggested that some linguistic societies might follow alternative patterns of color lexicalization, suggesting possible deviations from a hue-based model of color category universality that was central to the Berlin and Kay and WCS research program. The suggestions raised, for example, questioned whether a universal BCTs conjecture positing “a total universal inventory of exactly 11 basic color categories exists from which the 11 or fewer basic color terms of any language are always drawn” ([1], p. 2) was the best model to describe the extensive amount of data accumulated by the Berlin and Kay and WCS efforts. Robert E. MacLaury was among those who sought to explore alternative explanations, which was the impetus for his direction of the MCS project [11, 28–30]. In the 1970s MacLaury was an anthropologist working with Brent Berlin and Paul Kay, conducting fieldwork in Oaxaca, Mexico, and later working on the World Color Survey, and the Mesoamerican Color Survey was the basis for his PhD dissertation obtained from the University of California, Berkeley in 1986. In



Mesoamerican Color Survey Digital Archive, Fig. 1 Map depicting a portion of the geographic region surveyed by the MCS from the border of the United States to Nicaragua. Approximate locations of 116 MCS languages sampled are indicated as *blue crosses*, and 20 *red*

dots show WCS languages that duplicate some of the MCS languages surveyed. Of note are the dense samples of linguistically related MCS languages found in the areas of Oaxaca (37 languages), Guatemala (30 languages), and Mexico City (33 languages)

1997 MacLaury successfully used the MCS data as the basis for an alternative color-naming and categorization theory, called Vantage Theory [28]. That approach presents a model of color categorization at the lexeme level where color categories are constructed as vantages (see “► [Vantage Theory of Color](#)”). In the context of this approach, MacLaury’s book provides an extensive inventory of the organization and semantics of color categorization in Mesoamerica.

The MCS: History, Methodology, and Data

Historical Background

As noted earlier, the Mesoamerican Color Survey (MCS) includes 116 surveyed societies, each with a different Mesoamerican language or dialect, assessing approximately 900 native-language speakers from indigenous populations

in Mesoamerica. The MCS employed color stimuli from the WCS stimulus palette shown elsewhere in this volume (see “► [World Color Survey](#)”). Figure 1 above visually illustrates how the MCS substantially extends the geography of areas surveyed and thus extends beyond the wealth of information provided by the widely cited and valued WCS. It does so by replicating an estimated 20 languages found in the WCS, providing observations for 96 additional languages using the WCS’ standardized data collection procedures. Such an extensive survey of indigenous languages is likely no longer possible in 2015, due to dramatic increases in exposure over the last 25 years of native-language monolingual speakers to English-language broadcast media and entertainment.

The original MCS data was handwritten on paper, and while MacLaury published analysis and theory of MCS data (notably, [28]), thousands of pages of raw MCS data have never been publicly available and have not been

subjected to the kinds of highly informative methodological approaches that have more recently been used to reveal the insights learned from the WCS data (e.g., [3–6, 20, 21, 31, 32]). MacLaury focused on MCS language groups across Northwest Mesoamerica for a particular reason: to investigate what he conceived as an alternative theory of color category systems that emphasized “brightness”-based categories in naming – as an alternative to historically prominent theories emphasizing hue-based categories in naming – and to test hypotheses and give interpretations of the data grounded in a relativistic “framing” idea, which MacLaury called Vantage Theory.

MacLaury employed the WCS methodology from 1978 up until approximately 2004. Subsequent to 2004 the archive containing the MCS was inactive, unlike that of the WCS, with its raw data inaccessible to researchers while it was preserved in private storage by the MacLaury estate. Beginning in 2010 funding to initiate preservation and digitization of the archive was obtained. The initial preservation of the MCS archive, including its estimated 70 additional worldwide languages, came in the form of funding to purchase the archive from MacLaury’s estate, supported by the University of California Pacific Rim Research Program (award to K.A. Jameson, PI), and was augmented by a 2014 National Science Foundation funding (#SMA-1416907, K.A. Jameson, PI) that in part aimed to establish a public-access digital archive of transcribed MacLaury data and construct a user-friendly Wiki to generally facilitate collaborative color categorization research.

MCS Methodologies

Using individual samples and a color stimulus array shown in the WCS stimulus palette (see Fig. 1, ► [Kay and Cook’s ECST](#) entry), the MCS included three independent tasks done by every native-language speaker surveyed. These consist of (1) a “naming task” that involved the free naming of 330 color samples presented in a fixed random order; (2) a “focus task,” identification of named category best exemplars or category “foci”; and (3) a “mapping task” that

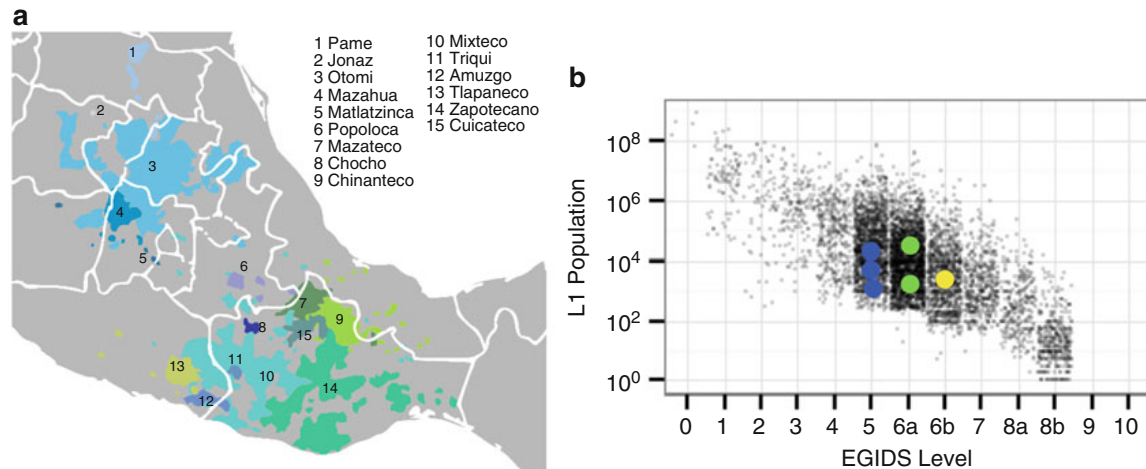
involved the demarcation of named category boundaries, producing a “boundary map.” The latter “mapping task” was used in Berlin and Kay’s original work, but for simplicity was omitted from the WCS. In the MCS, naming and focus data were collected for 900 individuals, and the mapping task method was reintroduced and assessed on an estimated 365 individuals. Data collection procedures used are detailed in works of MacLaury [28, 29] and are also similarly reviewed in WCS materials [2].

These MCS methods, along with the selection of languages it assessed, permit an unprecedented study of concept formation and linguistic representation, because of the following: (1) Many MCS languages are phylogenetically related to varying degrees. (2) Many MCS languages emphasize different perceptual dimensions (e.g., either primarily brightness or hue) and permit new analyses of dimensional salience. (3) Most MCS languages that satisfy (1) and (2) are related spatially, or geographically, and are thus known or expected candidates for linguistic borrowing, drift, and bilingual influences.

MCS Data

The research scope of the MCS data is substantial because it has features that can be exploited for addressing a range of heretofore unexplored color categorization research questions concerning how color categorization systems evolve, how semantic meaning drifts, and how meaning may be impacted by the influences of neighboring systems.

For example, Fig. 2’s representation of one language family (Otomanguean, Glottolog: otom 1299) shows one example that is well represented in the MCS, and, like that shown in panel (a), the MCS generally contains clusters of phylogenetically related languages, usually from contiguous geographic regions. Thus, for this one language family, as Fig. 2a illustrates, the MCS provides an estimated 369 surveyed participants, for which most subdivisions shown in panel (a) were assessed. Moreover, panel (b) illustrates that for some subdivisions, using Chinantec as an example, one may find several variations: That is, the MCS has six variations



Mesoamerican Color Survey Digital Archive, Fig. 2 Example of the typical level of representation for one language family (Otomanguean) found in the MCS. Panel (a) shows a map of Oaxaca and Mexico City areas shown above in Fig. 1. Numbered color regions represent Otomanguean family language branches listed in inset (Retrieved on 12/02/2014 from “https://commons.wikimedia.org/w/index.php?title=File:Otomanguean_Languages.png&oldid=89424714”. Image: wikipedia.org. https://commons.wikimedia.org/wiki/File:Otomanguean_Languages.png) exemplifying spatial density of related languages that is typical of the MCS. Panel (b) shows one of the family branches found in (a), namely, #9 Chinantec. Panel (b) shows a level of language diversity

and complexity typical of the MCS. In this example, six forms of Chinantec data are illustrated in the MCS dataset. Panel (b)’s six color dots locate MCS Chinantec languages within the cloud of all living languages (i.e., *small dots*) in relation to a language’s population (*vertical axis*) and its level of development or endangerment (*horizontal axis*). *Blue dots* are Developing (EGIDS 5) showing examples of Chinantec in vigorous use, *green dots* are Vigorous (EGIDS 6a) – unstandardized and in vigorous use among all generations. *Yellow dots* are In trouble (EGIDS 6b-7) (Chinantec living language data adapted and cited with permission from www.ethnologue.com [33]. Used by permission, © SIL International)

that cover three different levels of Chinantec language development/endangerment.

In addition, the MCS also has many languages that emphasize different dimensions of color appearance space compared to the hue-based systems emphasized by languages typically investigated in the literature. For each language surveyed, the actual MCS datasets typically include detailed investigator notes or ethnographies, as well as the datasheets on the three tasks. Variations on information contained and common formats of MCS datasheets are now described and illustrated:

1. **Naming task:** Informants were asked to name 330 loose color chips in a fixed random order. While investigators conducting the survey tended to employ two datasheet formats for collecting naming data, individual variations do exist in the raw datasheets. Figure 3 shows two common raw datasheet formats used to

collect naming task data. Image scans of raw MCS datasheets, and for a currently limited number of cases their transcribed digital data files, are provided on the archive website (detailed below).

Figure 3a, b show two common examples of MCS naming task datasheet formats as raw pdf image scans of the data. Panel (a) illustrates only one page (samples 1–81) surveyed from informant #1 (assessed in Guarijio language) showing a “list format” in which the investigator sequentially listed a participant’s responses to 330 samples evaluated in the fixed random order in which they were assessed. Panel (b) shows a second common format in which the investigator reports the participant’s responses to 330 samples, evaluated in the same fixed random order, in the appropriate row/column location on the MCS stimulus palette (see Fig. 1 in “► [World Color Survey](#)”). The palette format

a

oktoma : Guarizjo de Senora Baja
 fangwaga

20. 21. 22. 23. 24. 25. 26. 27. 28. 29. 30. 31. 32. 33. 34. 35. 36. 37. 38. 39. 40.

| | | | |
|-----------------------|-----------------------|------------------------|-------------------------|
| 1. sioxomúriame | 29. tisa sioxoname | 49. seltoname | 69. tisa sioxomúriame |
| 2. teesa sioxéname | 30. tisa seltoname | 50. warosahérame | 70. [tisa seltoname] |
| 3. seltoname | 31. tisa tesaname | 51. ohéoname | 71. tisa sioxoname |
| 4. ohéoname | 32. kawé sioxoname | 52. tesamúrame | 72. setapórame |
| 5. -sioxoname | 33. tisa sioxomúriame | 53. sioxomúriame | 73. tosa pórame |
| 6. warósa | 34. tisa oéomúriame | 54. wetapórame | 74. sioxomúriame |
| 7. kawé sioxoname | 35. tisa setamúrame | 55. ohéomúriame | 75. setamúriame |
| 8. setomúriame | 36. waróherame | 56. warósa hérame | 76. tisa ... wetapórame |
| 9. warosahérame | 37. setamúrame | 57. sioxomúriame | 77. [tosa pórame] |
| 10. kawé sioxoname | 38. okorá hérame | 58. kawé ohéoname | 78. - tosa seltoname? |
| 11. -oéomúriame | 39. sioxomúriame | 59. tisa wetapórame | 79. sioxoname |
| 12. sawaine | 40. setapórame | 60. tosa múrame | 80. setapórame |
| 13. -tisa ohéoname | 41. -sioxomúriame | 61. setapórame | 81. tisa sioxoname |
| 14. kawé ohéoname | 42. kawé sioxomúriame | 62. -sioxomúriame | 82. ohéoname |
| 15. seta pórame | 43. sawamúriame | 63. tisa sioxoname | 83. tisa warósa hérame |
| 16. tosa pórame | 44. oéomúriame | 64. tisa warósa hérame | 84. sioxoname |
| 17. tisa seltoname | 45. tisa sioxomúriame | 65. oéomúriame | 85. seta pórame |
| 18. -ohéoname (i) | 46. seltoname | 66. wetapórame | 86. [ohéoname (kawé)] |
| 19. tisa sioxomúriame | 47. sioxoname | 67. tosa pórame | 87. [setamúriame] |
| 20. kawé sioxoname | 48. sioxomúriame | 68. -ohéoname | 88. tisa sioxoname |

b

Guarizjo de Senora Baja - St. Jorda

| | | | | | |
|---|-----------------|-----------------------------|--------------------|--------------------|----------------|
| | | 1 | 2 | 3 | 4 |
| A | tokórame | | | | |
| B | tokórame | tisa warósa | tosa pórame | tesapórame | |
| C | tosa-pórame | setapórame tisa setamúriame | tisa warósa-hérame | tisa warósa-hérame | |
| D | tisa wetapórame | tisa - / setapórame | sawamúriame | tisa sawamúriame | |
| E | tosa wetapórame | setapórame | tisa seltoname | tisa seltoname | setamúriame |
| F | wetapórame | setapórame | setamúriame | setamúriame | seta name |
| G | tisa wetapórame | seltoname | kawé seltoname | setapórame | setapórame |
| H | kawé ohéoname | seltoname (kawé) | seltoname | setamúriame | setamúriame |
| I | ohéoname | [?] | kawé seltoname | kawé seltoname | [?] setapórame |
| J | ohéoname | [?] | | | |

Mesoamerican Color Survey Digital Archive, Fig. 3 MCS naming task datasheet formats shown as image scans of the raw data

of the naming task data often has participant's "focus" choices embedded in the datasheet as can be seen in panel (b)'s column 0 which has a "circled-X" symbol shown in cells **A0** and **J0**, recording the "foci" indicated by this participant for color appearances which in English language would gloss as the focus of "white" and the focus of "black," respectively.

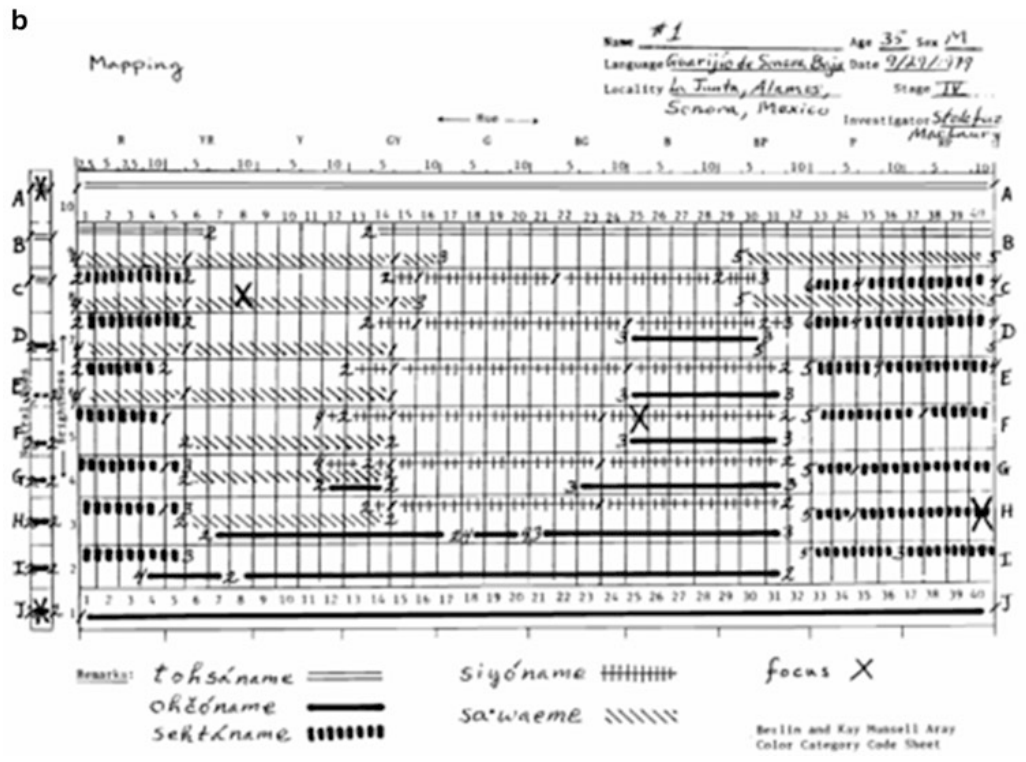
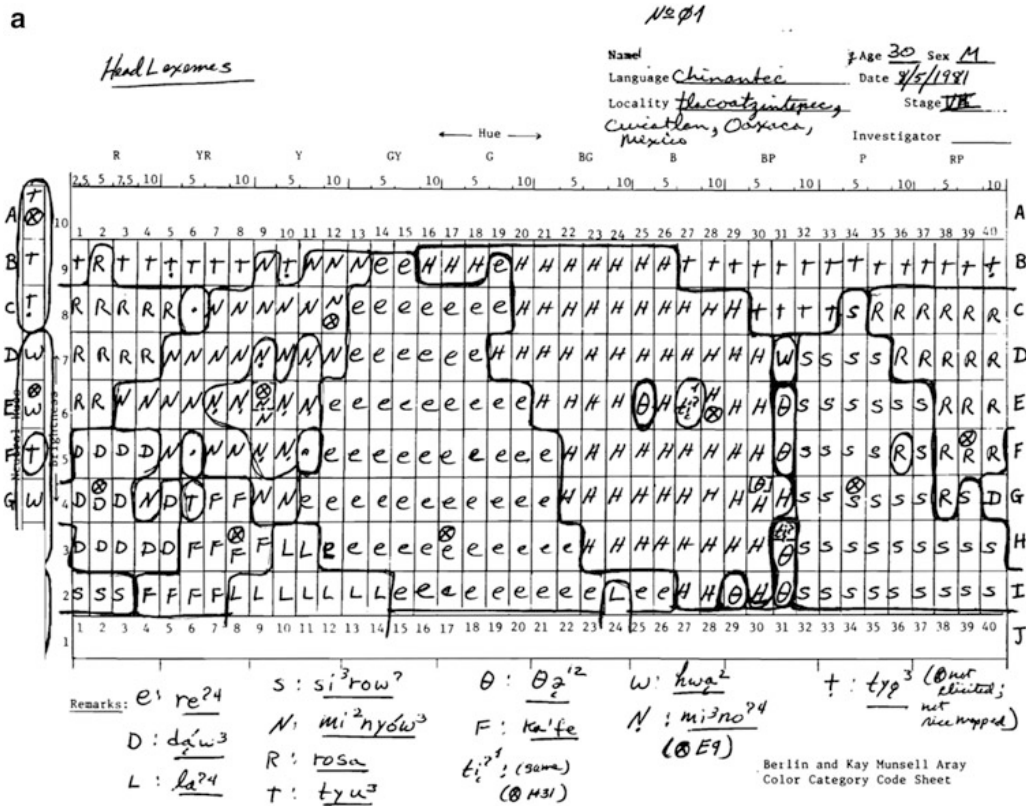
2. **Focus selection task:** Informants were shown the entire 330 color stimulus palette and asked to select the best example, or focus, of each different color name elicited in the naming task (see [28, 29] for method details). Again, individual variations do exist in the raw datasheets. Figure 4 shows two common raw datasheet formats used to collect focus selection task data. In panel (a) focus task data is reported in conjunction with naming task data, whereas in panel (b) focus selection task data is combined with mapping task data (mapping task is detailed below). In both formats shown in Fig. 4, focus selections are indicated by either a "circled-X" symbol (panel (a)) or by a plain "X" symbol alone (as in panel (b)). In focus task datasheets, it is common to find a dictionary of color terms elicited at the bottom of the sheet (as seen in both Fig. 4 panels). Mapping task datasheets most frequently delineate mapped category areas using hatched markings (e.g., Fig. 4b); however these data were also often recorded using color pencils (as shown in Fig. 5 and described in the next section). As with the naming task data, image scans of raw MCS datasheets, and transcribed digital data files for focus data, are provided on the archive website.
3. **Category mapping task:** For the same color category terms as in the focus selection task, informants were shown the entire 330 color stimulus palette and asked to indicate the regions of the stimulus palette in which the appearance of a given color term was represented. The informants in this way would map the color category stimulus regions glossed by each color term for which an investigator solicited a mapping (see [28, 29] for method details). This task is also referred to as the rice mapping task in the

literature, because grains of rice were often provided to informants to mark the many stimuli that were denoted by a given color term. Figure 5 shows two raw datasheets, for Guarijio and Chinantec languages, illustrating, above and beyond data variations, slight variations in the ways investigators recorded the surveyed information. Again, in some cases Focus locations are recorded, whereas in other Mapping Task datasheets Foci are not recorded. In mapping task datasheets, it is common to find a dictionary of color terms elicited at the bottom of the sheet and a pseudo-color code legend is given to interpret the regions defined using a particular color pencil as belonging to a specific color term (as seen in both Fig. 5 panels). (In such pseudo-color codes, the color of the ink used is not to be construed as resembling the color of the category data that it encodes.) As with the other MCS data, image scans of raw MCS datasheets are provided on the archive website.

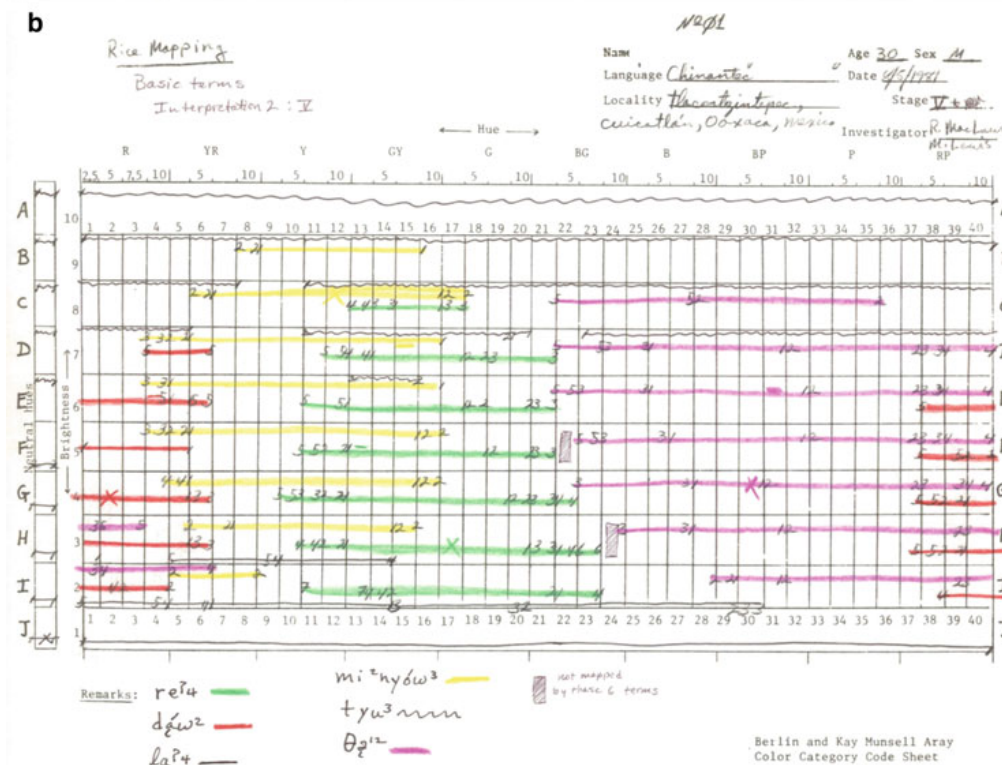
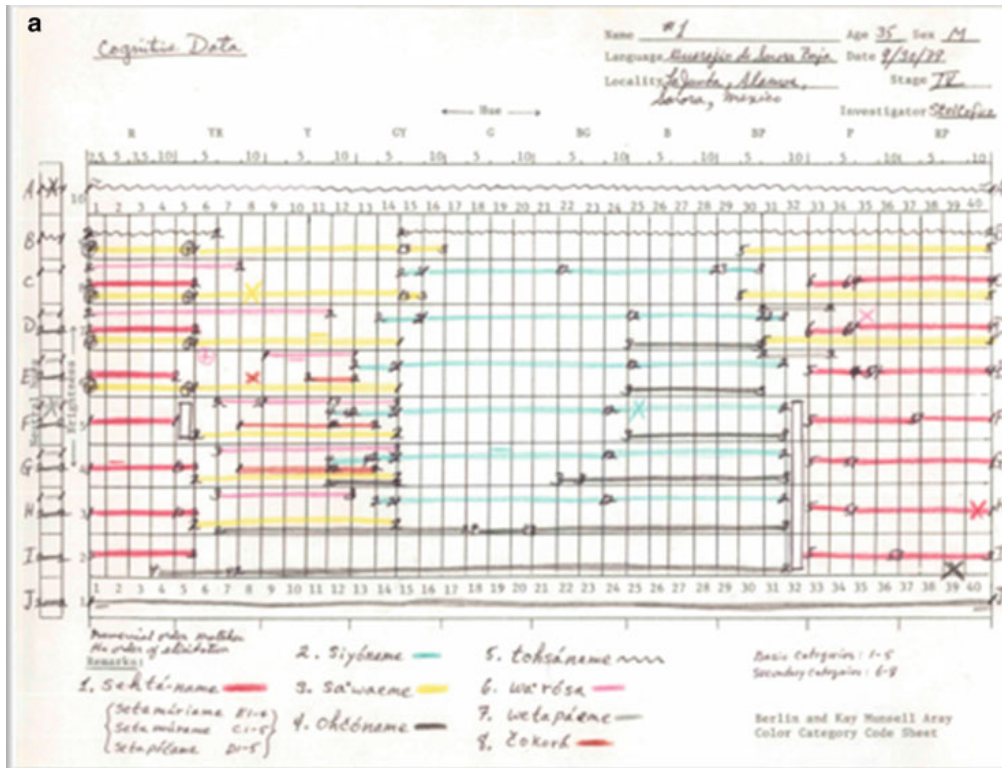
The above details of MCS data, including Figs. 3, 4, and 5, convey the kind of coding variation found throughout the datasheets in the MCS archive. The variation is not wholly unexpected as a large number of the 116 languages of the MCS were collected by different investigators, at different locations, over the approximately 5 years during which the MCS was conducted. Notwithstanding this amount of idiosyncratic coding style, the MCS raw data, as well as the rest of the MacLaury archive, appear to be of rather consistent and thorough formatting concerning the data that is collected, especially considering the mass of information that exists in the estimated 20,000 pages of the archive.

The MCS and the Robert E. MacLaury Color Categorization Digital Archive (ColCat)

Recognizing the value in the MCS data, the interdisciplinary *ColCat* research group at UC Irvine sought to convert the paper copy of MacLaury's



Mesoamerican Color Survey Digital Archive, Fig. 4 MCS focus and mapping task datasheet formats shown as image scans of the raw data



Mesoamerican Color Survey Digital Archive, Fig. 5 Two mapping task datasheet formats (panel (a) Guarijio and (b) Chinantec languages) presented as image scans of the raw data

research archive into a public-access digital database, similar to that developed for the WCS by Paul Kay and colleagues [2]. The entirety of MacLaury's MCS data is included in this *ColCat* digital archive [34]. A major aim of the ColCat project is to develop a platform to make the archive, including the MCS data, readily available to the scientific research and teaching community. The project presently makes available organized sets of scanned image files of datasheet pages in the archive acquired from MacLaury's estate and continues to implement an intuitive Internet-based user interface to serve as the front end of a collaborative research space with which researchers can readily explore and investigate their own basic research questions concerning the MCS data, as well as the color categorization data from other languages that MacLaury assessed. As mentioned earlier, neither the raw MCS data nor the MacLaury archive's additional non-Mesoamerican languages have been systematically organized for public use or previously published in an unanalyzed form. The additional 70 color categorization surveys (many with only a single informant, although others with frequently many more, and the largest with 40 informants) are valuable for their diversity in that they include native speakers from a wide range of languages including several Slavic languages, Hungarian, several Salishan languages of the Pacific Northwest United States, Zulu and several other South Africa/Zimbabwe languages, several native American languages, Germanic languages, European languages, Asian languages, and more.

The *Robert E. MacLaury Color Categorization (ColCat) Digital Archive* uses a content management system (CMS) for organizing, editing, and publishing the contents of the archive to a web-based graphical user interface (GUI) hosted at <http://colcat.calit2.uci.edu>. A partial list of features provided for the MCS through ColCat is summarized below:

1. User access level permissions and controls for the archive areas of the Wiki. User logins that define various levels of access for researchers, including those who may only want unregistered access to view database contents. Member level access for registered users who want to interact with the database. Editor level access as a trusted contributor to the database. Administrator access for systems operators who maintain database changes and review and approve contributions made by editor level users. Proper controls are implemented to maintain long-term integrity of the archive and to forestall loss and corruption of datafiles as well as accidental losses of data.
2. Facilitated search and query tools for exploring aspects of the datasets such as language type, language family, specific color terms, participant group subsets (e.g., by gender, age, etc.), survey geographic locations, and other filters that have been specified by content in the archived data.
3. Transcription quality indicators for evaluating conversion accuracy of handwritten data content into digitally addressable datasets. Because there are several ways the raw handwritten data in the archive can be converted into data addressable digital files (e.g., manual transcription, crowdsourced transcription, optical character recognition or OCR-based transcription, linguistic expert transcription, etc.), a useful feature of the ColCat Wiki includes utilities that permit assessment of each dataset's transcription quality and comparison across multiple forms of transcription of the same dataset. Transcription quality indicators appear on every surveyed language's archive page, and, where appropriate, indicate the degree of transcription completion, and computed degree accuracy, for each language's transcribed datasets.
4. To encourage users to explore the archive's original raw datasheets, the *ColCat Wiki* provides each survey's scanned image files of the archive's raw data organized in a browsable catalog/library of pdfs (like those seen in Figs. 3, 4, and 5 above). Onboard pdf files of raw datasheets permit visual inspection of investigator's ethnographic notes and embellishments; verification of digitized data content; and querying pdfs by language family,

language type, location and so forth, as well as querying by combinations of such tags.

5. An onboard suite of supporting resources for color categorization research. For example, (a) on every surveyed language's archive page (where appropriate and provided for by the archive's contents), maps of the geographic regions are provided where surveys were conducted, permitting visualization of geographical relatedness among the archive's datasets, as well as allowing for search and identification of database contents by exploring geographical locations (via a Google Map plugin) where the actual survey data originated. (b.) Bibliographic resources highlighting specific languages or language families, including lists of peer-reviewed journal publications, books, relevant media, and reference materials, provided for each of the archive's languages surveyed. (c.) Active links out to relevant Internet resources, such as *World Atlas of Language Structures Database* (<http://wals.info>), and the public-access site *Ethnologue: Languages of the World Database* (<http://www.ethnologue.com>).
6. As mentioned above a major aim of the archive project was to develop a platform to make the archive's data readily available to the scientific research and teaching community. Toward this end, the archive is designed to function as a collaborative research space, with user contributions enabled for trusted contributors, an informal blog feature to encourage a community of research exchange, and a platform for sharing and communicating results and for collective problem solving.

Additionally important was the aim of making the digital archive's data accessible to a broad audience comprised of users with varying levels of expertise in handling and analyzing large amounts of data. For this reason, a set of *ColCat data tools* was developed which allow users to search, filter, and analyze data in ways that best meet their research goals and technical abilities.

The planned data-handling tools and utilities aim to facilitate, for all user levels, the use of the

archive's contents for original color categorization research studies.

The ColCat digital archive, including the MCS, contains data at three levels:

1. Digital image scans, in pdf format, of approximately 12,500 handwritten pages of color-naming data, ethnographic notes, and support documents from 116 Mesoamerican languages, plus pdf image scans of an additional 10,000 pages surveyed globally from other languages and nations
2. Raw transcription data of these surveys converted and collected through OCR, crowd sourcing, and expert transcribers, and
3. "Official" transcription conversions for the archive's color-naming and categorization data

The Wiki's onboard data tools are designed for visualizing and analyzing the computer-addressable archive data at levels 2 and 3. Visualization tools for transcription data (level 2) include quantitative indicators of completeness and quality of transcriptions for each language and visualizations of the variance between or within transcription methods for entire languages or subsets of images within a language. Visualization tools for color-naming data (level 3) will include visualizations of the aggregated color-name mappings for each language including information about the proportion of agreement among all respondents or subsets of respondents through the use of filters – i.e., based on gender, age, etc.

Data export tools will allow users to download both transcription and naming data (levels 2 and 3) in formats conducive to off-line analysis (i.e., CSV, Excel, etc.). Although information about aggregated color-name mappings (level 3) will be based on data from all subjects in a language, users may be interested in computing and/or comparing consensus color-name mappings for subpopulations based on their own criteria. To that end, downloadable tools will allow users to take the same analyses and visualizations that are used on the Wiki for complete populations and apply them to subpopulations that were exported by the

user. The expectation is that a variety of user types will employ the Wiki's visualization tools to help identify subsets of data in the archive that meet specific research goals, while other researchers will be more inclined to export data in a familiar format (i.e., in an established WCS format) and apply their own off-line analyses. Regardless of the case, the ColCat data tools make the entire archive's data accessible to a wide range of researchers, students, and other scientists who would like to explore and learn from the archive but are not well equipped to manage and process large quantities of data.

Uses of the MCS Archive

As a complementary database to the World Color Survey database, the addition of the MCS archive allows for an even greater, more extensive, assessment of behaviors of cross-cultural color cognition and color naming. For example, reconsider the Fig. 2(a) discussion above, concerning the 369 surveyed Otomanguean language family informants for which MCS data exists. These MCS informants can now be independently compared to, or even analyzed in conjunction with, the 175 previously surveyed WCS informants from the Otomanguean language family. In addition to the novel surveys it adds, it is this kind of duplication and augmentation that the MCS provides as a complementary database to the WCS that will increase the value of both archives and present new opportunities to greatly extend understanding of how the data in the MCS and the WCS inform on the questions originally asked by the surveys' directors, namely, what specifically are the universal and culturally specific patterns of color naming and categorization that can be empirically demonstrated for large samples of the world's diverse languages and how do they inform us on the larger question of general human tendencies for categorizing and naming objects perceived in the world.

Cross-References

- ▶ [Berlin and Kay Theory](#)
- ▶ [Centroid and Boundary Colors](#)
- ▶ [Color Categorical Perception](#)
- ▶ [Color Category Learning in Naming-game Simulations](#)
- ▶ [Color Dictionaries and Corpora](#)
- ▶ [Dynamics of Color Category Formation and Boundaries](#)
- ▶ [Effect of Color Terms on Color Perception](#)
- ▶ [Multilingual/Bilingual Color Naming/ Categories](#)
- ▶ [Vantage Theory of Color](#)

References

1. Berlin, B., Kay, P.: *Basic Color Terms: Their Universality and Evolution*. University of California Press, Berkeley (1969)
2. Kay, P., Berlin, B., Maffi, L., Merrifield, W.: *The World Color Survey*. Center for the Study of Language and Information, Stanford (2003)
3. Kay, P., Regier, T.: Resolving the question of color naming universals. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 9085–9089 (2003)
4. Regier, T., Kay, P., Cook, R.: Focal colors are universal after all. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 8386–8391 (2005)
5. Lindsey, D.T., Brown, A.M.: Universality of color names. *Proc. Natl. Acad. Sci. U. S. A.* **103**, 16609–16613 (2006)
6. Lindsey, D.T., Brown, A.M.: World color survey color naming reveals universal motifs and their within-language diversity. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 19785–19790 (2009)
7. Kay, P., Regier, T.: Color naming universals: the case of Berinmo. *Cognition* **102**, 289–298 (2007)
8. Davidoff, J., Davies, I.R.L., Roberson, D.: Color categories in a stone-age tribe. *Nature* **398**, 203–204 (1999)
9. Roberson, D., Davies, I.R.L., Davidoff, J.: Color categories are not universal: replications and new evidence from a stone age culture. *J. Exp. Psychol. Gen.* **129**, 369–398 (2000)
10. Roberson, D., Hanley, J.R.: Color categories vary with language after all. *Curr. Biol.* **17**, 605–606 (2007)
11. MacLaury, R.E.: Color-category evolution and shuswap yellow-with-green. *Am. Anthropol.* **89**(1), 107–124 (1987)
12. Paramei, G.V.: Singing the Russian blues: an argument for culturally basic color terms. *Cross-Cult. Res.* **39**(1), 10–38 (2005)

13. Dedrick, D.: Color language universality and evolution: on the explanation for basic color terms. *Philos. Psychol.* **9**(4), 497–524 (1996)
14. Jameson, K.A.: Culture and cognition: what is universal about the representation of color experience? *J. Cogn. Cult.* **5**(3–4), 293–347 (2005)
15. Alvarado, N., Jameson, K.A.: Confidence judgments and color category best exemplar salience. *Cross-Cult. Res.* **39**(2), 134–158 (2005)
16. Jameson, K.A.: Why GRUE? An interpoint-distance model analysis of composite color categories. *Cross-Cult. Res.* **39**(2), 159–194 (2005)
17. Jameson, K.A.: Where in the world color survey is the support for the hering primaries as the basis for color categorization? In: Cohen, J., Matthen, M. (eds.) *Color Ontology and Color Science*, pp. 179–202. The MIT Press, Cambridge (2010)
18. Davies, I.R.L., Corbett, G.G.: A cross-cultural study of color grouping: evidence for weak linguistic relativity. *Br. J. Psychol.* **88**(3), 493–517 (1997)
19. Komarova, N.L., Jameson, K.A.: A quantitative theory of human color choices. *PLoS One* **8**(2), e55986 (2013). doi:10.1371/journal.pone.0055986
20. Bimler, D.: Are color categories innate or internalized? Hypotheses and implications. *J. Cogn. Cult.* **5**(3), 265–292 (2005)
21. Bimler, D.: From color naming to a language space: an analysis of data from the world color survey. *J. Cogn. Cult.* **7**(3), 173–199 (2007)
22. Bimler, D., Uusküla, M.: Clothed in triple blues: sorting out the Italian blues. *J. Opt. Soc. Am. A* **31**, A332–A340 (2014)
23. Narens, L., Jameson, K.A., Komarova, N.L., Tauber, S.: Language, categorization, and convention. *Adv. Complex Syst.* **15**(03n04), 1150022 (2012)
24. Jameson, K.A., Komarova, N.L.: Evolutionary models of color categorization. I. Population categorization systems based on normal and dichromat observers. *J. Opt. Soc. Am. A*, **26**(6), 1414–1423. Featured Reprint in *The Virtual J. Biomed. Opt.* **4**(8), (2009)
25. Jameson, K.A., Komarova, N.L.: Evolutionary models of color categorization. II. Realistic observer models and population heterogeneity. *J. Opt. Soc. Am. A*, **26**(6), 1424–1436. Featured Reprint in *The Virtual J. Biomed. Opt.* **4**(8), (2009)
26. Komarova, N.L., Jameson, K.A.: Population heterogeneity and color stimulus heterogeneity in agent-based color categorization. *J. Theor. Biol.* **253**, 680–700 (2008)
27. Komarova, N.L., Jameson, K.A., Narens, L.: Evolutionary models of color categorization based on discrimination. *J. Math. Psychol.* **51**, 359–382 (2007)
28. MacLaury, R.E.: *Color and Cognition in Mesoamerica: Constructing Categories as Vantages*. University of Texas Press, Austin (1997)
29. MacLaury, R.E.: *Color in mesoamerica. Vol. 1: a theory of composite categorization*. Doctoral dissertation. University of California, Berkeley. UMI University Microfilms, No. 8718073, Ann Arbor (1986)
30. MacLaury, R.E.: From brightness to hue: an explanatory model of color category evolution. *Curr. Anthropol.* **33**(2), 137–186 (1992)
31. Regier, T., Kay, P., Khetarpal, N.: Color naming and the shape of color space. *Language* **85**, 884–892 (2009)
32. Webster, M., Kay, P.: Individual and population differences in focal colors. In: MacLaury, R., Paramei, G., Dedrick, D. (eds.) *Anthropology of Color*, pp. 29–53. John Benjamins, Amsterdam (2007)
33. Paul, L.M., Simons, G.F., Fennig, C.D. (eds.): *Ethnologue: Languages of the World*, Eighteenth edition. Dallas, Texas: SIL International. Online version: <http://www.ethnologue.com> (2015)
34. Jameson, K.A., Gago, S., Deshpande, P.S., Benjamin, N.A., Chang, S.M., Tauber, S., Jiao, Y., Harris, I.G., Xiang, Z., Bhakta, H.R., MacLaury, R.E.: *The Robert E. MacLaury Color Categorization (ColCat) Digital Archive*. <http://colcat.calit2.uci.edu/>. The California Institute for Telecommunications and Information Technology (*Calit2*), UC Irvine (2015)