

# Functional Semantic Categories for Art History Text - Human Labeling and Preliminary Machine Learning <sup>\*</sup>

Rebecca J. Passonneau<sup>1</sup>, Tae Yano<sup>2</sup>, Tom Lippincott<sup>3</sup> and Judith Klavans<sup>4</sup>

<sup>1</sup> Center for Computational Learning Systems, Columbia University, USA

<sup>2</sup> Department of Computer Science, Carnegie Mellon University, USA

<sup>3</sup> Department of Computer Science, Columbia University, USA

<sup>4</sup> College of Information Studies, University of Maryland, USA

**Abstract.** The CLiMB project investigates semi-automatic methods to extract descriptive metadata from texts for indexing digital image collections. We developed a set of functional semantic categories to classify text extracts that describe images. Each semantic category names a functional relation between an image depicting a work of art historical significance, and expository text associated with the image. This includes description of the image, discussion of the historical context in which the work was created, and so on. We present interannotator agreement results on human classification of text extracts, and accuracy results from initial machine learning experiments. In our pilot studies, human agreement varied widely, depending the labeler's expertise, the image-text pair under consideration, the number of labels that could be assigned to one text, and the type of training, if any, we gave labelers. Initial machine learning results indicate the three most relevant categories are machine learnable. Based on our pilot work, we implemented a labeling interface that we are currently using to collect a large dataset of text that will be used in training and testing machine classifiers.

## 1 Introduction

The work presented here was developed in the context of the Computational Linguistics for Metadata Building (CLiMB) research project, which has been investigating methods for automated support to image catalogers and other image professionals for locating subject matter metadata in electronic versions of scholarly texts [5]. The CLiMB project is developing a Toolkit for image catalogers that would allow them to access electronic versions of the types of texts they consult manually, and that could thus lead to improved access through richer indexing. Toolkit design goals include proposing terms for the cataloger to review. Because the Toolkit is co-evolving with changes in practice concurrent with the growth of digital image collections, we have followed an

---

<sup>\*</sup> We thank the current and past members of the project members; numerous advisors who reviewed our labeling categories and tested the interface; and especially, volunteer labelers from University of Maryland, Columbia University, Drexel University, and Indiana University.

iterative design process [11], where we elicit feedback from image professionals and catalogers along the way. Here we continued this approach in addressing the question of how to classify text associated with images into functional semantic categories.

Figure 1 shows an image taken from the ARTstor Art Images for College Teaching collection (AICT): <http://www.arthist.umn.edu/aict/html/ancient/EN/EN006.html>. It depicts a sunk relief portrait of Akhenaten and his family. Also shown is an extract from a few paragraphs taken from an art history survey text describing an image of the same work. Note that if the terms **Akhenaten** and **shrine** were used to index the image, it would not be clear whether the image depicts a shrine or Akhenaten or both. The word Akhenaten occurs in a sentence about Akhenaten's role in fostering the Amarna style, and in another sentence indicating that he is depicted in the work. The word shrine occurs in a sentence indicating how the depicted work was used. Our goal is to automatically tag sentences like these prior to semi-automatic or automatic extraction of the bold face terms, and for the extracted terms to be associated with tags corresponding to our semantic categories; see the bottom of Figure 1. This could permit terms to be filtered or prioritized during the term selection process, depending on the semantic tag they are occur with. It could also facilitate search; for a user who wants an image of a shrine, it would be possible to exclude cases where *shrine* does not come from text that describes the content of the image. In consultation with



Historical Context: Akhenaten  
Image Content: Akhenaten  
Historical Context: shrine

Of the great projects built by **Akhenaten** hardly anything remains . . . . Through his choice of masters, he fostered a new style. Known as the Amarna style, it can be seen at its best in a sunk relief portrait of **Akhenaten** and his family. The intimate domestic scene suggests that the relief was meant to serve as a **shrine** in a private household.

**Fig. 1.** Illustration of an image and associated descriptive text.

experts, we developed a set of seven categories to apply paragraphs or sentences extracted from art history survey texts, where the text extracts are about a specific image. A larger number of categories would lead to much sparser data; a smaller number would lead to categories that are less distinct. Two of the categories, for example, are **Image Content**, defined as text that describes the objective content of the image, and **Implementation**, text that describes the manner in which the work was created (technique, style, technical challenges, etc.).

In a set of four pilot studies, interannotator agreement among humans varied widely, depending on the labeler's expertise, the image-text pair under consideration, the number of labels that could be assigned to one text, and the type of training, if any, we gave

labelers. Human agreement improved if annotators could select multiple labels, which is consistent with our previous results on a lexical semantic annotation task [10]. We also found that agreement was higher among domain experts, and that the consistency of the labeling depended heavily on the image/text pair under consideration.

Our seven semantic categories vary in relevance, frequency and distinguishability. Thus we do not anticipate attempting to apply machine learning to every category. Using texts labeled during our pilot studies, we have initial results on three of the classes. For example, using a separate naive Bayes classifier for each category, we have been able to achieve 80% accuracy. This indicates that a larger scale learning study is feasible.

In section 2, we summarize related work on inter-annotator agreement among human indexers and annotators. Section 3 describes the two art history survey texts we draw from and our datasets. Section 4 describes our pilot studies on human labeling, and our current large scale effort. We present preliminary learning results in section 5. We conclude in section 6 with general observations about the prospect for adding subject descriptors to large image collections.

## 2 Related Work

There are relatively few discussions of inter-annotator or inter-indexer consistency for image indexing and classification tasks. Two works that address the topic deeply and broadly are [8] and [4]. In the twenty plus years since Markey's analysis of forty years of inter-indexer consistency tests, no comparable review has appeared, and her observations still hold. Although her goal was to use the conclusions from previous work to sort through the issues involved in indexing visual material, all the tests referenced in her paper were on indexing of printed material. She notes significant variability, with accuracy (percent agreement) results ranging from 82% to a low of 4%.

The Giral and Taylor study [4] concerned indexing overlap and consistency on catalog records for the same items in architectural collections, including an analysis of subject descriptors. On large ( $\geq 1400$ ) samples of records from the Avery Index to Architectural Periodicals and the Architectural Periodicals Index, they compare proportions of items in their samples that match according to a variety of criteria, and compute 90% confidence intervals. Only 7% of items match entirely, and they find some element of overlap in descriptors in only about 40% of the remaining cases ( $\pm 3\%$ ).

Markey noted two features of documents that affect inter-indexer consistency: document length, and the complexity of the document, which is difficult to quantify. Our image/text pairs, which correspond to Markey's documents, are quite short. We did not attempt to measure complexity, but we did find wide variation in labeling consistency depending on the image/text pair being labeled. This indicates that the image/text pairs have inherent properties that make them more or less difficult for humans to agree on.

Markey found no significant difference between the indexers with or without experience in using such schemes. She found no higher levels of inter-indexer consistency among subject specialists, as compared with non-specialists. This is in contrast to our results. In our pilot studies, the two developers of the categories (the first two co-authors) were the most familiar with them, and had the highest interannotator agreement. In our current large-scale labeling effort, the highest agreement is found among the

most expert pairs of labelers. We also found that across studies, agreement increased when we provided more training. Markey found that using a standardized scheme led to higher inter-indexer consistency, ranging from 80% to 34% (in contrast to 4%; see above). This is roughly the range we find, using a different metric but a similar scale.

### 3 Texts

The domain of digital images and texts we focus on parallels the ARTstor *Art History Survey Collection (AHSC)*. ARTstor is a Mellon funded non-profit organization developing digital image collections and resources. The AHSC is a collection of 4,000 images that is the product of a collaboration with the Digital Library Federation's Academic Image Cooperative project and the College Art Association. One of our motivations for focusing on the AHSC is that it is based on thirteen standard art history survey texts, thus there is a strong correlation between the images and texts that describe them. The AHSC images all have metadata providing the name of the work, the artist, date, and so on, but very few have subject matter metadata.

We are currently using two of the texts from the AHSC concordance of thirteen art history survey volumes. Both books have a similar lineup of chapter topics, though there are some differences in text layout. They cover a broad time range, from Neolithic art to late 20th century. Each text contains roughly thirty chapters (approximately five megabytes in digital format), with twenty to forty color images each.

For research purposes, we created electronic versions of the two texts, encoded in TEI compliant xml. TEI is a widely used interdisciplinary standard of text representation. The rules are defined in the TEI Lite customized schema. (See [http://www.tei-c.org/Lite/teiu5\\_split\\_en.html](http://www.tei-c.org/Lite/teiu5_split_en.html) for more detail of this schema.) Chapters, subdivisions, and paragraphs (but not sentences) have distinctive xml tags.

To facilitate the construction of image/text pairs for our text labeling experiments, we employed software that had been created as a module for importing text into an image indexer's Toolkit we have implemented. The software module relies primarily on the relative position of xml tags for image plates, major text divisions, and paragraph boundaries. It takes a chapter as input, and produces a list of all the plates in the chapter, with each plate number associated with a sequential list of associated paragraph numbers. We manually correct the output before importing the data into our labeling interface. Using Google image search, we locate non-copyrighted images of the works depicted in the book plates.

## 4 Text Labeling Experiments

### 4.1 Semantic Category Labels

Our current guidelines give four pieces of information per semantic category: the category name, one or two questions the labeled text should answer, one or two paragraphs describing the category, and four image/text pairs that exemplify each category. For the Image Content category (or label), the questions are *Does the text describe what the art work looks like? What conventional use of symbols does the artist rely on?*

Over a period of four months, we developed a set of functional semantic categories for classifying paragraphs and sentences in our art history survey texts. Three criteria motivated the classification. Most important, we did not attempt to develop an independent set of categories based on existing image indexing work. We took the information in the texts as our starting point. Second, the set of classes were designed to apply to all chapters regardless of time period, and to allow most paragraphs or sentences to fall into a specific category, rather than to a default *Other* class. Finally, we worked with an image librarian at Columbia University and a metadata expert to arrive at a relevant set.

Table 1 summarizes our seven semantic categories. The column on the left indicates the name of the label, and the column on the right gives a highly abbreviated description of the type of textual content that should be assigned a given label. The labels appear here in the same order that they appear in the interface, which puts the most central category first (Image Content), and which lists categories that have a similar focus together. Thus the first three categories are all about the depicted art work (form, meaning, manner); Biographic and Historical Context are both about the historical context.

**Table 1.** Seven Functional Semantic Categories for Labeling Text Extracts.

| Category Label        | Description   |
|-----------------------|---|
| Image Content         | Text that mentions the depicted object, discusses the subject matter, and describes what the artwork looks like, or contains.   |
| Interpretation        | Text in which the author provides his or her interpretation of the work.  |
| <i>Implementation</i> | Text that explains artistic methods used to create the work, including the style, any technical problems, new techniques or approaches, etc.  |
| Comparison            | Text that discusses the art object in reference to one or more other works to compare or contrast the imagery, technique, subject matter, materials, etc.   |
| Biographic            | Text that provides information about the artist, the patron, or other people involved in creating the work, or that have a direct and meaningful link to the work after it was created.             |
| Historical Context    | Text describing the social or historical context in which the depicted work was created, including who commissioned it, or the impact of the image on the social or historical context of the time. |
| Significance          | Text pointing to the specific art historical significance of the image. This usually applies to a single sentence, rather than to an entire paragraph.  |

During the first month, we arrived at a provisional set of six categories consisting of everything in Figure 1 apart from the italicized category, which now has the name *Implementation*, and developed our first set of guidelines. We added the seventh category after a month or so of pilot work. During the remaining three months we created versions of our labeling guidelines, each revising the category names and definitions.

#### 4.2 Materials: Datasets, Annotation Constraints, Annotators, and other Task Parameters

We created three sets of image/text pairs, and we used them in the experiments listed in Table 2. The second column of the table shows for each experiment which of the three image/text sets was used. Set 1 consisted of thirteen images and 52 associated paragraphs. Set 2 consisted of nine images and 24 associated paragraphs. Set 3 consisted of ten images taken from two new chapters, and was used in for sentence labeling (159 sentences) as well as paragraph labeling (24 paragraphs).



**Table 2.** Annotation Task Parameters.

| Exp | Set | Images | Units | Label Set | Labels/Par | Annotators |
|-----|-----|--------|-------|-----------|------------|------------|
| 1   | 1   | 13     | 52    | 6         | any        | 2          |
| 2   | 2   | 9      | 24    | 7         | any        | 2          |
| 3   | 2   | 9      | 24    | 7         | two        | 5          |
| 4a  | 3   | 10     | 24    | 7         | one        | 7          |
| 4b  | 3   | 10     | 159   | 7         | one        | 7          |

Labelers were recruited from the team of project researchers, their acquaintances, and colleagues at other institutions involved in image indexing.

The two parameters of most interest for comparing the experiments appear in columns five (Labels/Par) and six (Annotators). For the first two experiments, the first two co-authors were the annotators, and the number of labels that could be assigned to a single paragraph was unrestricted. In experiment 1, the maximum number of labels for a single paragraph was three; each annotator used three labels twice; 99% of the labelings consisted of one or two labels. In experiment 2, 71% of all labels from both annotators were one or two labels; the maximum of four labels occurred once per annotator.

Due to the relative infrequency of more than two labels in experiments 1 and 2, we added a restriction in experiment three that only two labels could be used. In experiment four, we restricted the paragraph level labeling further to a single label, but expanded the labeling task to include sentences.

For experiments 1 through 3, the labeling was done with pen and paper. For experiment 4, we implemented a custom browser-based labeling interface that included the guidelines, training materials, and labeling task. Based on our experience with this browser interface, we developed a much more flexible web-based labeling interface using the Django python environment.

In our pilot studies and current data collection effort, labelers worked independently at remote sites, and could suspend and resume work at will. After experiment 3, labelers were required to go through a training sequence (approx. one hour). Up to four paragraphs were associated with each image, but in most cases there were one or two paragraphs. Paragraphs were presented one at a time along with the corresponding image. When we began using sentences as well as paragraphs, labelers would first select a paragraph label; then the labeler would be presented with the same paragraph in a sentence-by-sentence format, in order to label the sentences. Labelers were given the opportunity to review and revise their choices.

### 4.3 Evaluation Metrics

We report interannotator agreement using Krippendorff’s  $\alpha$  [6], which factors out chance agreement. It ranges from 1 for perfect agreement to values close to -1 for maximally non-random disagreement, with 0 representing no difference from chance distribution. An advantageous feature of  $\alpha$  is that instead of treating agreement as a binary distinction, it permits the use of a distance metric to weight the degree of agreement from 0 to 1. Because annotators could make multiple selections, we used a distance metric we refer to as MASI [9]. It is intended for set-based annotations, and gives partial

agreement credit when the annotators’ sets overlap. Our experiments typically allowed annotators to assign multiple labels to the same text. If one annotator assigns the single label **{Image Content}** to the same text that another annotator labels **{Image Content, Implementation}**, a non-weighted agreement measure would assign a score of 0 for non-agreement. In contrast, MASI would assign a weighting of  $\frac{1}{3}$  (see [9] for details).

#### 4.4 Human Labeling Pilot Studies

**Table 3.** Interannotator consistency of paragraph labeling under multiple conditions.

| Exper. | Dataset | Label Set | #Choices | #Labelers | Alpha <sub>MASI</sub> |
|--------|---------|-----------|----------|-----------|-----------------------|
| 1      | Set 1   | 6         | any      | 2         | 0.76                  |
| 2      | Set 2   | 7         | any      | 2         | 0.93                  |
| 3      | Set 2   | 7         | two      | 5         | 0.46                  |
| 4a     | Set 3   | 7         | one      | 7         | 0.24                  |
| 4a'    | Set 3   | 7         | merge 4b | 7         | 0.36                  |
| 4b     | Set 3   | 7         | one      | 7         | 0.30                  |

Results for the four pilot experiments appear in Table 3. Experiment 2, with the final labeling set of seven labels, the first two co-authors as the sole annotators, and any number of label choices, had the best results. It improved on experiment 1, which used an earlier, less well-defined set of labels, but also a larger set of units (52 rather than 24 paragraphs) from two texts, rather than from a single text. Otherwise, the labeling criteria were the same: multiple labels could be assigned to each paragraph.

Experiment 3 was the first attempt to use a larger set of annotators. We hypothesized that with each new annotator, the number of distinct combinations of labels would increase, with the result that a large number of annotators would result in a large set of distinct classes, and correspondingly smaller classes. In order to guard against this possibility, we restricted the number of labels that annotators could apply to two. The resulting  $\kappa$  score of 0.46 reflects the relative unfamiliarity of a majority of the five annotators with the labeling categories and domain. When we computed interannotator consistency for all combinations of annotators from the set of five, we found that the two experienced annotators had values on the three measures (0.88, 0.88, 0.90) that were consistent with the results of experiment 2.

We collected sentence labelings for the first time in experiment 4: 4a pertains to the paragraph labels, and 4b to the sentence labels. For experiment 4a', we computed agreement on paragraphs based on merging the sentence labels. We created a relatively short label consisting of each distinct type of label applied to any sentence in the paragraph; if three sentences of a five-sentence paragraph were labeled Image Content and two were labeled Historical Context, the paragraph level label we compute is the multi-label consisting of these two labels.

Experiments 4a and 4b yielded the poorest results, which we attribute to the constraint that annotators could only apply one label. The seven labelers included the first two co-authors, plus five new annotators. As in experiment 3, we computed interanno-

tator agreement metrics for all combinations of annotators in experiment 4a. For all 21 pairs of annotators, agreement ranged from a low of 0.15 to a high of 0.32.

In addition to a great deal of variation in reliability across annotators, we find wide variation depending on the individual units consisting of images/text sets. For the ten units, agreement ranged from 0.12 to 0.40.

#### 4.5 Initial Results of Large Scale Human Labeling

A key feature of our new labeling interface is that labelers can work concurrently on distinct labeling tasks. We plan to collect data on between six and ten datasets. We have currently collected labelings from six annotators on the first dataset consisting of 25 images (45 paragraphs, 313 sentences).

In the new interface, annotators can choose any number of labels. We recruited four new labelers, and used one previous labeler. Results for the first dataset, which consists of 25 images and 48 associated paragraphs (313 sentences), are better than experiments 3 and 4 where we also used multiple annotators. We believe the improvement is due to the training provided in the interface, and the lack of constraint on the number of labels annotators could pick. However, the expert annotators disagreed with the training labels, which were selected from the best annotators on our pilot data, who were not experts. As we continue to collect data, we will revise the training labels using the best consensus from the best annotators, and investigate the impact.

As in experiment 4, sentence labeling had a higher agreement than for paragraphs. For sentences the overall  $\alpha$  measure was 0.45, compared with 0.40 for paragraphs. For all combinations of 2 to 5 coders, paragraph labeling agreement ranged from 0.56 to 0.27, and sentence labeling agreement ranged from 0.55 to 0.33 for sentences. The two coders who are experts in the area of image indexing had the highest interannotator agreement. As in the pilot studies, there was a significant variation in agreement, depending on the unit, ranging from a high of 0.70 to a low of 0.16.

The most frequent label combination for both paragraphs and sentences was the single label Image Content. There were 47 distinct combinations of labels for sentences, of which 34 were label pairs and five were triples; the remaining 8 unigram labels were the seven labels plus the default "Other". There were 38 combinations for paragraphs: 7 singletons, 20 pairs, 10 triples, and 1 combination of four labels.

## 5 Preliminary Machine Learning Results

Using data from our pilot studies of human labeling, augmented by an additional set of images labeled by one of the co-authors, we investigated the learnability of three categories: Image Content, Historical Context and Implementation. There were insufficient examples from the other categories. All learning was done using WEKA [13], a Java-based toolkit that implements a wide range of machine-learning algorithms using a standard input format.

We created three types of feature sets. Set A consisted of word features selected on the basis of significant chi-square tests. Set B consisted of hand-picked features in approximately half a dozen groups, such as words and phrases characteristic of the



art history domain (e.g., *masterpiece*), and words and phrases referring to parts of the human body. Set C consisted of the union of Sets A and B.

We tested several types of learners, including naive bayes, SVM and tree-based classifiers. Naive bayes performed best overall. On ten-fold cross-validation, the highest classification accuracy on Image Content relied on feature set C, and achieved 83% accuracy, compared with 63% for Historical Context and 53% for Implementation. The highest accuracy for Historical Context used feature set A: 70%. Using a random forest classifier for the Implementation class, we achieved an accuracy of 80%.

## 6 Conclusions and Future Work

We have presented a detailed analysis of the development of a functional semantic labeling for art history texts, and have identified some of the problems that arise in achieving consistently high agreement scores across multiple annotators. One issue, the variance across texts, is more difficult to address unless we alter the texts. The other key issue is that annotators with expertise are much more consistent with each other than non experts are. As we continue collecting data, and updating our training with the expert consensus on previously labeled examples, we hope to learn something about training and experts. However, we have found that we can still achieve high accuracy with machine learning. As pointed out in [12], the relationship between interannotator agreement and learnability is not a predictable one.

We believe the initial learning results are quite promising. One difficulty for learning functional semantic categories is that many of the content words are not relevant features, since they will be different for descriptions of different images. In contrast, for topical text classification, content words are often sufficient for automatic classification, which is the intuition behind approaches such as latent semantic indexing. By using features such as verb tense, which distinguishes the **Image Content** class from others, we have achieved high results on relatively small datasets. On the other hand, since our categories are functional, they may transfer more easily to texts that are substantially different from our training and test materials.

As illustrated in the introduction, we anticipate that classifying text into functional semantic categories can provide more control over selection of metadata. Our categories have a rough correspondence with categories discussed in the image indexing literature [7, 3, 2]. As a result, it should be possible to map between our categories and the types of controlled vocabularies used in university visual resource centers. The external knowledge sources our project has examined include the three Getty resources (Art and Architecture Thesaurus, Thesaurus of Geographic Names, Union List of Artist Names), the Library of Congress Authorities and Library of Congress Thesauri for Graphic Materials, and ICONCLASS, a library classification for art and iconography.

## References

1. R. Artstein and M. Poesio.  $\text{Kappa}^3 = \text{alpha}$  (or beta). Technical Report NLE Technote 2005-01, University of Essex, Essex, 2005.

2. M. Baca. *Practical Issues in Applying Metadata Schemas and Controlled Vocabularies to Cultural Heritage Information*. The Haworth Press, Inc., 2003. Available through Library Literature, last accessed July 25, 2006.
3. H. Chen. An analysis of image queries in the field of art history. *Journal of the American Society for Information Science and Technology*, pages 260–273, 2001.
4. A. Giral and A. Taylor. Indexing overlap and consistency between the Avery Index to Architectural Periodicals and the Architectural Periodicals Index. *Library Resources and Technical Services* 37(1):19-44, 1993.
5. J. Klavans. Using computational linguistic techniques and thesauri for enhancing metadata records in image search: The CLiMB project. Article in preparation.
6. K. Krippendorff. *Content Analysis: An Introduction to Its Methodology*. Sage Publications, Beverly Hills, CA, 1980.
7. S. S. Layne. Some issues in the indexing of images. *Journal of the American Society for Information Science*, pages 583–8, 1994.
8. K. Markey. Interindexer consistency tests: a literature review and report of a test of consistency in indexing visual materials. *Library and Information Science Research*, pages 155–177, 1984.
9. R. Passonneau. Measuring agreement on set-valued items (MASI) for semantic and pragmatic annotation. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC)*, 2006.
10. R. Passonneau, N. Habash and O. Rambow. Inter-annotator agreement on a multilingual semantic annotation task. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC)*, 2006.
11. R. J. Passonneau, D. Elson, R. Blitz, and J. Klavans. CLiMB Toolkit: A case study of iterative evaluation in a multidisciplinary project. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC)*, 2006.
12. D. Riedsma and J. Carletta. Reliability measurement: there's no safe limit To appear in *Computational Linguistics*.
13. I. Witten and E. Frank. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann: San Francisco, 2000.