

## Optical Music Recognition in Manuscripts from the Ricordi Archive

FEDERICO SIMONETTA, GSSI - Gran Sasso Science Institute, Italy

RISHAV MONDAL, LUCA ANDREA LUDOVICO, and STAVROS NTALAMPIRAS, University of Milan, Italy



Fig. 1. Example of handwritten score from the Archivio Storico Ricordi from the manuscript score of “La Bohème”, by Giacomo Puccini.

The Ricordi archive, a prestigious collection of significant musical manuscripts from renowned opera composers such as Donizetti, Verdi and Puccini, has been digitized. This process has allowed us to automatically extract samples that represent various musical elements depicted on the manuscripts, including notes, staves, clefs, erasures, and composer’s annotations, among others. To distinguish between digitization noise and actual music elements, a subset of these images was meticulously grouped and labeled by multiple individuals into several classes. After assessing the consistency of the annotations, we trained multiple neural network-based classifiers to differentiate between the identified music elements. The primary objective of this study was to evaluate the reliability of these classifiers, with the ultimate goal of using them for the automatic categorization of the remaining unannotated data set. The dataset,

Authors’ Contact Information: Federico Simonetta, GSSI - Gran Sasso Science Institute, L’Aquila, Italy, federico.simonetta@gssi.it; Rishav Mondal, rishav.mondal@studenti.unimi.it; Luca Andrea Ludovico, luca.ludovico@unimi.it; Stavros Ntalampiras, stavros.ntalampiras@unimi.it, University of Milan, Milan, Italy.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2024 Copyright held by the owner/author(s).

Manuscript submitted to ACM

Manuscript submitted to ACM

complemented by manual annotations, models, and source code used in these experiments are publicly accessible for replication purposes.<sup>1</sup>

CCS Concepts: • **Applied computing** → *Sound and music computing*; **Digital libraries and archives**; **Arts and humanities**; • **Computing methodologies** → *Supervised learning by classification*; **Neural networks**.

Additional Key Words and Phrases: Optical Music Recognition, Neural Networks, Computer Vision, Music

#### ACM Reference Format:

Federico Simonetta, Rishav Mondal, Luca Andrea Ludovico, and Stavros Ntalampiras. 2024. Optical Music Recognition in Manuscripts from the Ricordi Archive. In *Audio Mostly 2024 - Explorations in Sonic Cultures (AM '24)*, September 18–20, 2024, Milan, Italy. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3678299.3678324>

## 1 INTRODUCTION

The Ricordi Archive, or *Archivio Storico Ricordi* in Italian, is a vast repository of historical documents amassed by the Italian publisher, Ricordi. The archive is renowned for its digitized manuscripts of distinguished opera composers such as Donizetti, Verdi, and Puccini. These manuscripts, digitized and systematically cataloged in a database maintained by the author's institution,<sup>2</sup> are a significant asset for musicological and historical research.

This study aims to annotate the entire database with pertinent musical symbols, thereby improving the accessibility and discoverability of these priceless manuscripts. To achieve this, we developed a suitable Optical Music Recognition (OMR) methodology.

OMR is a subfield of computer science dedicated to converting music notation from a visual format, such as scanned images or printed music sheets, into a digital form that can be manipulated by software. It serves as a bridge between physical music representations and their digital equivalents, with the aim of automating the interpretation of music symbols for various applications, including content retrieval, digital music libraries, and musicological research [3, 14, 15].

The evolution of OMR, particularly for printed music, is largely attributed to advancements in image processing and machine learning. Although early efforts relied heavily on rule-based systems [6], contemporary strategies use modern neural architectures. These architectures excel in feature extraction and help to accurately identify and classify musical symbols across various datasets and notation styles [2, 8, 11].

Handwritten Music Recognition (HMR) focuses on the recognition of handwritten scores, adding an extra layer of complexity due to the unique styles and nuances inherent in individual handwriting. Recent methodologies suggest various strategies, including data augmentation and transfer learning, to enhance system performance and address the challenges specific to HMR [3, 18].

Automating the analysis of handwritten music presents significant obstacles due to the intricacy and variability of human handwriting; here, progress has been driven by the implementation of advanced machine learning models, specifically Convolutional Recurrent Neural Networks. Such models effectively capture both the spatial characteristics of the image and the sequential nature of music notation, which are crucial for the development of successful OMR systems for handwritten scores [1].

The development of large-scale datasets has been fundamental to the advancement of OMR technology. A key resource in OMR research is the MUSCIMA++ dataset [9]. This dataset, comprising 140 pages of handwritten music scores, is meticulously annotated with over 91000 symbols across 107 classes. It facilitates a wide range of tasks,

<sup>1</sup>Code repository: <https://github.com/LIMUNIMI/RicordiArchiveOMR>; dataset repository: <https://zenodo.org/doi/10.5281/zenodo.11186095>

<sup>2</sup><https://ricordi.lim.di.unimi.it>

including symbol classification and notation graph assembly. MUSCIMA++ is an extension of the CVC-MUSCIMA dataset, which includes 1,000 music sheets from 50 musicians. The detailed annotations of musical symbols and their interrelations in MUSCIMA++ are crucial for training and evaluating OMR systems. Its extensive coverage of musical notations, from notes to articulation marks, and structural annotations detailing symbol relationships, make it an invaluable tool for machine learning applications in music symbol detection and recognition. Moreover, MUSCIMA++ serves as a benchmark in the OMR field, contributing significantly to the development of technologies that convert sheet music into digital formats, thereby enhancing the accuracy and reliability of OMR systems.

The DeepScores [17] dataset is another valuable resource, consisting of high-quality images of printed music divided into approximately 300,000 sheets of musical scores with nearly a hundred million small objects. It provides ground truth for object classification, detection, and semantic segmentation, focusing on the recognition of small objects and intricate details in musical scores.

Carefully curated and annotated datasets, such as MUSCIMA++, are instrumental in advancing the field of HMR. They not only improve the accuracy of recognition models, such as Convolutional Recurrent Neural Networks, but also establish rigorous benchmarks for evaluating these models against the complexities of real-world musical notations. By offering a comprehensive collection of samples that cover the full range of human handwriting variability, these datasets are essential to refine the capabilities of end-to-end OMR systems. Consequently, these systems can accurately interpret the intricate nuances of handwritten music scores, addressing a significant academic and practical challenge.

In this study, we address a series of the challenges in OMR by providing a new dataset of musical symbols from real-world annotated manuscripts, and by training and evaluating several neural classifiers to distinguish between these symbols. Moreover, we expect that the current work will be an important step towards the automatic annotation of the entire Ricordi Archive. The data set is published online.<sup>3</sup>

The rest of the paper is structured as follows: in Sec. 2 we provide an insight into the history and the activities of the Ricordi Archive. Sec. 3 presents the dataset of images focusing on preprocessing and annotation processes, while in Sec. 4 we conduct a series of experiments assessing the ability of statistical models in music symbol identification. Sec. 5 discusses the obtained results, and, finally, in Sec. 6 we draw the conclusions.

## 2 BACKGROUND

The Ricordi Archive originated alongside the publishing house Casa Ricordi, established in 1808. Regarded as a paramount private musical repository, it safeguards the original handwritten scores of 23 out of Verdi's 28 operas, all operas by Giacomo Puccini (except *La Rondine*), and numerous works by composers such as Bellini, Rossini, Donizetti, as well as contemporary composers like Nono, Donatoni, Sciarrino, and Bussotti.

The archive's exceptional significance lies in the diversity of its materials, offering an articulated view of Italian culture, industry, and society. This archive preserves an extensive collection of visual materials associated with numerous premieres worldwide and locally, encompassing set and costume designs, photography compilations, correspondence, and business records. These resources empower researchers to reconstruct the inception of significant operas and the evolution of the musical publishing industry during the 19<sup>th</sup> and early 20<sup>th</sup> Centuries. Furthermore, the visual collection covers various artistic domains such as painting, stage design, and decorative arts, offering insights into costume history, jewelry design, stage properties, and the broader publishing landscape. It also sheds light on the relationship between publishers and artists across different fields and provides glimpses into the theatrical realm. Scholars can trace the

<sup>3</sup><https://zenodo.org/doi/10.5281/zenodo.11186095>

personal and professional trajectories of numerous composers from their earliest works, such as Verdi’s *Oberto Conte di San Bonifacio* and Puccini’s *Le Villi*, to their most important operas like Verdi’s *Falstaff* and Puccini’s unfinished *Turandot*.

The Ricordi Archive preserves approximately 8,000 scores, over 16,000 letters exchanged among musicians, librettists, singers, and other stakeholders, around 10,000 set and costume designs, more than 9,000 librettos, 6,000 historical photographs, and a substantial collection of Art Nouveau and Art Deco posters crafted by prominent artists of the era.

The digitization initiative of the historical archive stemmed from a collaborative effort involving the Italian Ministry of Culture, the National Department of Archives and Libraries, the Italian Supervisory Council for Libraries and Cultural Institutions, Casa Ricordi, Biblioteca Nazionale Braidense, and the Laboratory of Music Informatics (*Laboratorio di Informatica Musicale*, LIM) of the University of Milan. This project, initiated in 2006, adheres to the standards established by the Italian National Library Service (*Servizio Bibliotecario Nazionale*, SBN), which is overseen by the Central Unified Catalogue Institute (*Istituto Centrale per il Catalogo Unico*, ICCU). Given the archive’s artistic and historical significance, its preservation is subject to the regulations and oversight of the Ministry of Culture.

### 3 DATASET

The original core of the digitization campaign of Ricordi Archive consisted of about 3000 digitized images, mainly handwritten scores by Donizetti, Puccini, Verdi, and Respighi.

#### 3.1 Preprocessing

The creation of the dataset necessitated preliminary processing to identify pertinent objects and reduce the annotation effort in its initial phase. This process entailed the following steps:

- Staff Line Removal – A neural autoencoder-based algorithm [7] was employed to identify staff lines. Given the distinct clarity of the staff lines in the 19th-century documents from the archive, this method proved highly effective. The staff lines, being printed rather than handwritten, were easily distinguishable from the musical symbols, thereby enhancing the reliability of this step;
- Blob Detection – The Difference of Gaussians (DoG) method, as implemented by the `scikit-image` Python module [12], was utilized to identify the musical symbols in the images. We used  $\sigma \in [10, 50]$  and a threshold of 0.1. Although this step does not guarantee the detection of all relevant objects in the images, it was tuned to be particularly sensitive to the ink regions. As a result, a large number of false positives were included to minimize the occurrence of false negatives, i.e., relevant objects not included in the dataset;
- Rescale and Save – The grayscale images of the detected blobs were stored after rescaling their intensity values to  $[0, 255]$ .

The construction of our inter-referencing database, which utilizes JSON files, began with the collection of blob images. In total, 473,238 blobs were extracted. These images were then systematically stored to facilitate easy access and reference.

The initial step in this process involved generating a grayscale image for each image in the original Ricordi Archive by removing the staff lines. Each of these grayscale images was then associated with a JSON file, which contained a reference to the original image, the path to the grayscale image without staff lines, and a list of JSON files associated with the blobs detected in the image.



Fig. 2. Screenshot of the interface used for annotating the dataset. Texts are in Italian.

Subsequently, a set of blobs was detected for each image from which staff lines had been removed. Each detected blob was stored as a grayscale image, and a JSON file was created for each blob. This file contained a reference to the parent image (the image without staves) and the bounding box of the blob. This systematic approach ensured that each blob and its associated data could be easily traced back to its parent image.

### 3.2 Annotation

The annotation phase was facilitated by 15 local high-school students with music reading skills, who were divided into seven groups of two or three. We developed a custom interface that enabled the students to assign labels to each blob image.

For reference, each detected blob was highlighted with a bounding box within a larger excerpt of the original image. Additionally, an HTML link to the original image was provided for further examination if necessary. A screenshot of the annotation interface is depicted in Figure 2.

We identified 16 classes of objects that could be recognized as blobs. These included page border, erasure, smudge, printed and handwritten text, rest, single or multiple notes, single or multiple chords, alterations, clefs, ornaments, multiple categories (with and without music signs), and an “other” category (with and without music signs) for objects that did not fit into any of the other categories.

To assess the accuracy of the annotations, a sample of 500 blobs was randomly selected for cyclic annotation by all annotators. This process involved the selection of a control blob with a 20% probability at each annotation cycle. The control blobs were initially used during the annotation process to gamify the labeling work by providing the annotators with simple scores that reflect the quality of their work. The score was calculated based on the average of two factors: the Spearman correlation coefficient of the annotator’s labels, which represents intra-agreement, and the Spearman correlation coefficient between a) the average of the annotator’s labels for each control blob and b) the average of the annotations already stored in the database, provided by other annotators, representing inter-agreement.

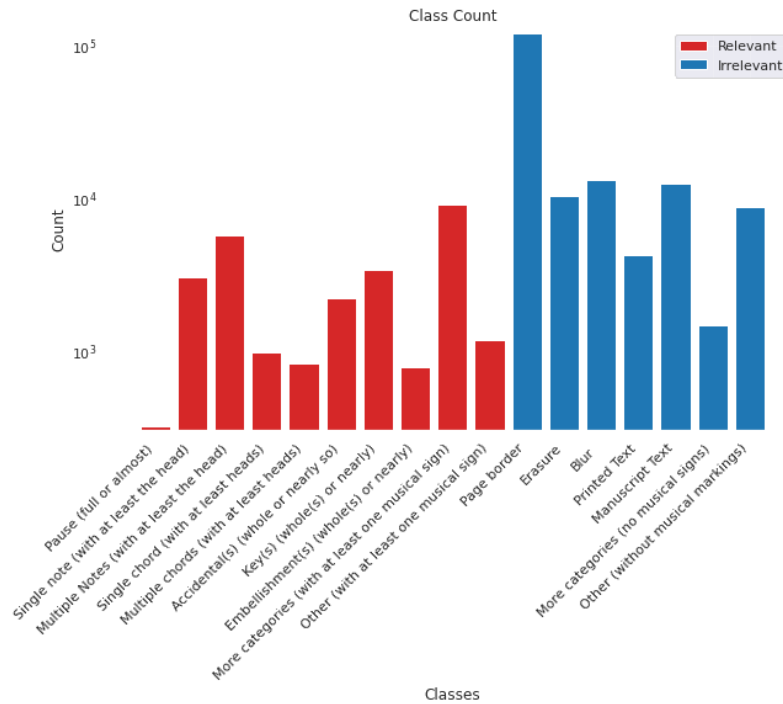


Fig. 3. Original distribution of the blob images across the classes before merging the less frequent ones. Note that the Y axis is in log scale.

The same control blobs were then used to compute the inter and intra-annotator agreement in order to assess the annotation quality. We first calculated a reference label for each control blob  $i$  and annotator  $j$  as the mode of the ratings given by annotator  $j$  to the control blob  $i$ . On these sets of ratings, we computed Krippendorff’s alpha (0.72), indicating that the raters generally agreed on the representation of the symbols. The intra-rater agreement was computed using Krippendorff’s alpha over each annotator’s labels separately, and was found to be between 0.52 and 0.71 depending on the annotator, indicating that the annotators were generally coherent on the representation of the symbols.

To identify the reasons for the partial disagreement, we first identified a reference annotation for each control blob  $i$  across all annotators using the mode of the labels given by all annotators. We then analyzed the normalized confusion matrix resulting from the annotated labels and the reference labels. If a label was annotated correctly on average, the maximum value of each row was along the diagonal. The remaining values were then in reference to such value, so that a value near to 1 would mean that there was confusion among the annotators about the meaning of that label. We merged the classes that had a normalized confusion value greater than 0.5. We performed this step – i.e. computation of the confusion matrix, normalization, and class merging – iteratively until no classes were merged. This procedure led us to merge the “multiple notes” and “multiple chords” labels, as well as the labels “blurs” and “multiple categories without musical signs”, thus resulting in 14 classes. The class merging increased Krippendorff’s alpha to 0.84 for the inter-rater agreement, while the the intra-rater agreement raised to [0.63, 0.79]. The distribution is shown in Figure 3.

To enhance the training of machine learning models, we addressed the issue of class imbalance by consolidating certain classes into a single category, termed as “Remaining”. This amalgamation involved classes with sample sizes less

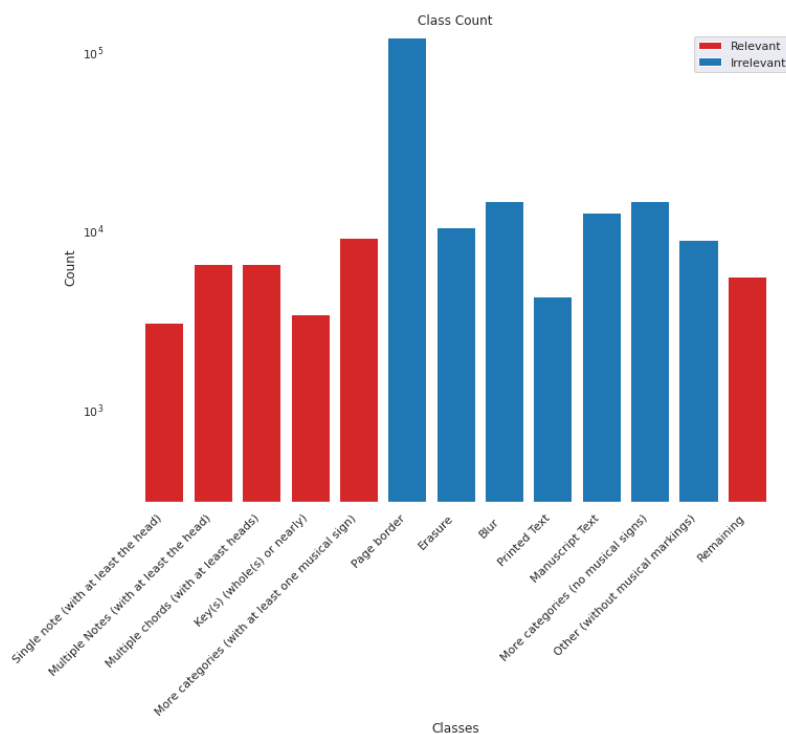


Fig. 4. Distribution of the images across the classes, after the merging of the less frequent classes. Note that the Y axis is in log scale.

than  $0.75 \times m$ , where  $m$  represents the median class cardinality. Consequently, pauses, embellishments, single chords, accidentals, and “Other (with musical signs)” were merged, resulting in a total of 11 classes. The revised distribution of samples across these classes is depicted in Figure 4.

For the purpose of simplifying the classification task, we categorized the labels into two clusters: “musically relevant” and “musically irrelevant”. This distinction signifies the presence or absence of musical signs in the blob. The inter-rater agreement for this binary annotation was measured using Krippendorff’s alpha, yielding a value of 0.89. The intra-rater agreement ranged between 0.74 and 0.79.

The final class distribution in the proposed dataset, comprising a total of 198,159 annotated blobs, is illustrated in Figure 4. Examples of blobs are shown in Figure 5. We further provide predefined splits for training and testing sets, applicable to both binary and multiclass classification tasks.

## 4 EXPERIMENTS

To evaluate the efficacy of statistical models in recognizing musical symbols, we fine-tuned three renowned deep learning classifiers: ResNet, DenseNet, and GoogleNet. We also employed an advanced AutoML method [5] to compare neural networks with conventional machine learning techniques.

Considering the significant imbalance depicted in Figure 3, we initially subsampled the classes with the highest cardinalities to achieve perfectly balanced training and validation sets. This subsampling involved randomly selecting  $n$  samples from the largest categories, where  $n$  corresponds to the number of samples in the smallest category. While this

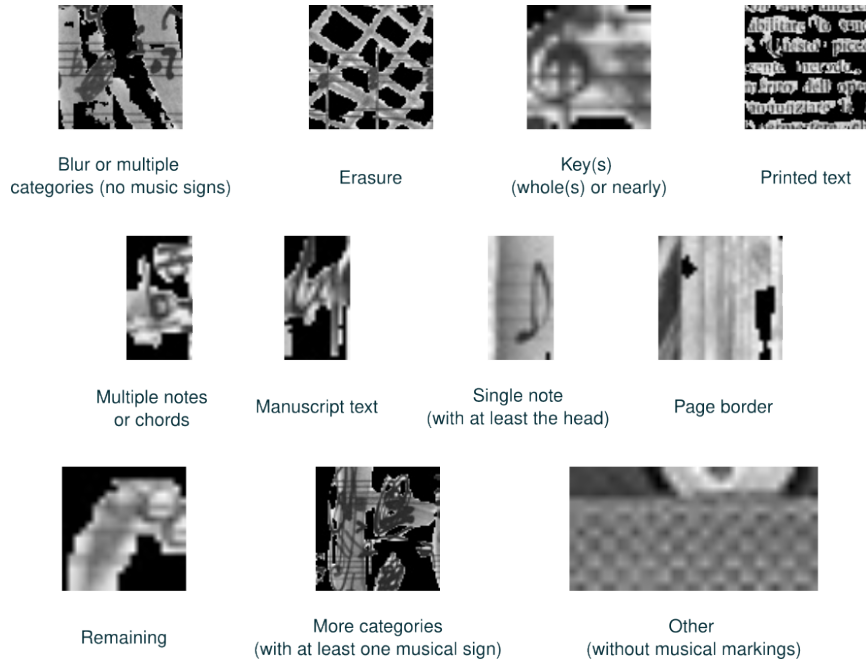


Fig. 5. Examples of blobs for each class in the dataset.

method yields balanced training and validation sets, the test sets remain highly imbalanced. Therefore, it is crucial to implement appropriate validation measures that account for class imbalance during testing, as shown in Tables 1 and 2.

In all instances, we enhanced the training set by applying random rotations of up to  $\pm 10$  degrees, random flips with probability of 0.5, brightness, contrast, and saturation jitters with factor 0.25. To improve the quality of the images, we also implemented Gaussian blur denoising with kernel of 3 and sigma 1.5 and contrast correction to 1.5 of the original contrast. To ensure compatibility with the ImageNet pre-trained weights, we renormalized the image channels and resized all images to  $256 \times 256$  pixels.

For the deep-learning classifiers, we utilized the pre-trained weights available in the `torchvision` library, obtained from the ImageNet 1K dataset [4, 13]. We re-trained all models using standard cross-entropy loss and the 1 cycle learning rate policy [16] with Stochastic Gradient Descent, setting the maximum learning rate at 0.01. The models were trained for 500 epochs with early stopping based on validation loss, exhibiting a patience of 20 epochs. This resulted in approximately 40 epochs of actual training, with a maximum of 70 epochs for GoogleNet in binary classification. We used a batch size of 64 and allocated 68%, 17%, and 15% of the dataset for training, validation, and testing, respectively.

We placed particular emphasis on uncertainty analysis by examining the activations in the networks' final layer. Ideally, a confidence close to 1 suggests that the input sample is located in a region of the feature space familiar to the network, while a confidence near 0 indicates unfamiliarity. Consequently, we can disregard low-confidence predictions to minimize the risk of incorrect classifications. In Bayesian statistics, uncertainties are categorized as epistemic and aleatoric [10]. Epistemic uncertainty pertains to the model parameters, indicating that the model parameters have learned the under-represented region of the data. Aleatoric uncertainty, on the other hand, relates to the data itself, suggesting that the data is inherently noisy.



Table 1. Neural network performances for the **binary** classification task. Note that the balanced accuracy is equivalent to the average recall. The best average values for each measure are highlighted in bold. The symbol “-” means that no data was retained for that class at that level of confidence.

DenseNet															
Confidence level	Precision					Recall					F1-score				
	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%
Irrelevant	0.98	0.99	0.99	1.00	1.00	0.82	0.89	0.93	0.97	0.99	0.89	0.94	0.96	0.98	1.00
Relevant	0.44	0.56	0.66	0.77	0.87	0.87	0.94	0.95	0.98	0.98	0.58	0.70	0.78	0.86	0.92
Average	0.71	0.78	0.83	0.88	0.93	0.85	0.91	0.94	0.97	<b>0.98</b>	0.74	0.82	0.87	0.92	<b>0.96</b>
ResNet															
Confidence level	Precision					Recall					F1-score				
	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%
Irrelevant	0.98	0.99	0.99	1.00	1.00	0.82	0.90	0.95	0.98	1.00	0.89	0.94	0.97	0.99	1.00
Relevant	0.45	0.58	0.69	0.82	0.90	0.87	0.94	0.94	0.97	0.95	0.59	0.72	0.79	0.89	0.92
Average	0.71	0.79	0.84	0.91	<b>0.95</b>	0.85	0.92	0.95	<b>0.98</b>	0.97	0.74	0.83	0.88	0.94	<b>0.96</b>
GoogleNet															
Confidence level	Precision					Recall					F1-score				
	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%
Irrelevant	0.98	0.99	0.99	1.00	1.00	0.82	0.90	0.94	0.97	1.00	0.89	0.94	0.97	0.99	1.00
Relevant	0.45	0.58	0.68	0.78	0.90	0.88	0.94	0.95	0.98	0.97	0.59	0.71	0.79	0.87	0.93
Average	0.71	0.78	0.84	0.89	<b>0.95</b>	0.85	0.92	0.95	<b>0.98</b>	<b>0.98</b>	0.74	0.83	0.88	0.93	<b>0.96</b>

In this study, we employed the entropy of the neural output as a foundation for the confidence score, serving as a comprehensive measure of both epistemic and aleatoric uncertainty. Mathematically, given the network outputs  $y_i, i \in [1, N]$  for classifying  $N$  classes, the entropy is calculated as:

$$H = \sum_{i=1}^N \text{SoftMax}(y_i) \times \log_N(\min(1, \text{SoftMax}(y_i) + \epsilon)),$$

Here,  $\min(\cdot)$  and  $\epsilon$  are incorporated to circumvent numerical instability, and  $\text{SoftMax}$  is conventionally defined as:

$$\text{SoftMax}(y_i) = \frac{e^{y_i}}{\sum_{j=1}^N e^{y_j}}$$

For  $N = 2$ ,  $H$  aligns with the classical Shannon entropy computed using bits as information units. Generally, it always lies within  $[0, 1]$ , allowing the computation of a confidence score as  $1 - H$ .

We conducted all experiments for both the binary classification task, which involves distinguishing between “musically relevant” and “musically irrelevant” blobs, and the multi-class classification task, which involves differentiating among the 14 classes of objects.

## 5 RESULTS

For the binary classification task, all three deep learning models achieved a balanced accuracy of 85% and an f1-score of 74%, as detailed in Table 1. In contrast, the constant predictor yielded a balanced accuracy of 50% and an f1-score of 46%, underperforming the random guessing.

By considering varying confidence levels, we observed a consistent monotonic trend for both accuracy and the proportion of retained test data. This indicates that higher accuracies can be attained by predicting fewer samples, as illustrated in Fig. 6. For example, with GoogleNet, a balanced accuracy of 95% and an f1-score of 88% can be achieved by retaining only samples with a confidence exceeding 50%, which constitutes 64% of the test set.

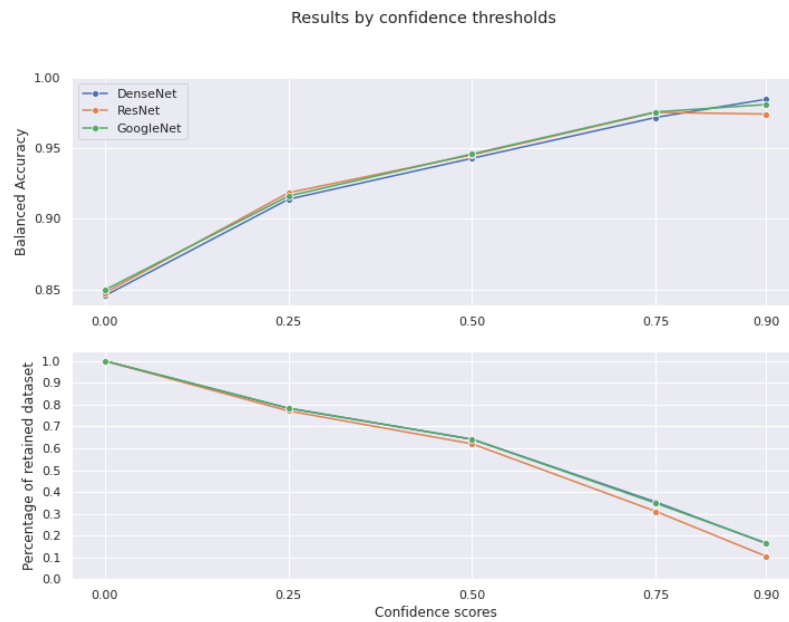


Fig. 6. Trend of the balanced accuracy and percentage of retained test data for various level of confidences for the **binary** task.

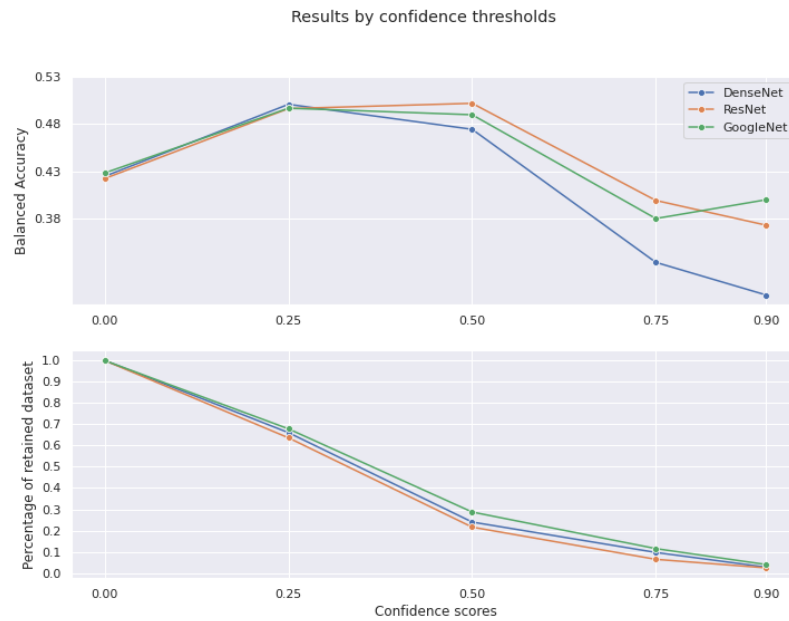


Fig. 7. Trend of the balanced accuracy and percentage of retained data for various level of confidences for the **multiclass** classification task.

Table 2. Neural network performances for the **multiclass** task. Note that the balanced accuracy is equivalent to the average recall. The best average values for each measure are highlighted in bold. The symbol “-” means that no data was retained for that class at that level of confidence.

DenseNet															
Confidence level	Precision					Recall					F1-score				
	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%
Single note (with at least the head)	0.11	0.16	0.00	0.00	-	0.31	0.27	0.00	0.00	-	0.16	0.20	0.00	0.00	-
Manuscript Text	0.42	0.57	0.86	1.00	0.00	0.37	0.46	0.46	0.13	0.00	0.40	0.51	0.60	0.24	0.00
Remaining	0.15	0.21	0.00	0.00	0.00	0.21	0.18	0.00	0.00	0.00	0.17	0.19	0.00	0.00	0.00
Printed Text	0.56	0.69	0.85	0.92	0.94	0.67	0.86	0.97	0.99	1.00	0.61	0.77	0.90	0.96	0.97
Key(s) (whole(s) or nearly)	0.50	0.72	0.91	1.00	1.00	0.77	0.94	1.00	1.00	1.00	0.61	0.82	0.95	1.00	1.00
Blur or multiple categories (no music signs)	0.26	0.30	0.59	0.33	0.00	0.49	0.58	0.50	0.05	0.00	0.34	0.39	0.54	0.09	0.00
Multiple notes or chords	0.29	0.38	0.79	0.00	0.00	0.38	0.59	0.56	0.00	0.00	0.33	0.46	0.65	0.00	0.00
Page border	0.95	0.95	0.98	0.99	0.99	0.77	0.89	0.98	1.00	0.99	0.85	0.92	0.98	0.99	0.99
Other (without musical markings)	0.23	0.31	0.43	0.75	0.00	0.18	0.27	0.47	0.50	0.00	0.20	0.29	0.45	0.60	0.00
Erasure	0.26	0.34	0.67	0.00	0.00	0.41	0.39	0.27	0.00	0.00	0.32	0.37	0.39	0.00	0.00
More categories (with at least one musical sign)	0.28	0.36	0.00	0.00	0.00	0.11	0.09	0.00	0.00	0.00	0.16	0.14	0.00	0.00	0.00
Average	0.36	0.45	0.55	0.45	0.29	0.42	0.50	0.47	0.33	0.30	0.38	0.46	0.50	0.35	0.30
ResNet															
Confidence level	Precision					Recall					F1-score				
	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%
Single note (with at least the head)	0.12	0.16	0.00	0.00	-	0.33	0.30	0.00	0.00	-	0.18	0.21	0.00	0.00	-
Manuscript Text	0.44	0.56	0.76	1.00	0.00	0.38	0.50	0.50	0.13	0.00	0.41	0.53	0.60	0.24	0.00
Remaining	0.14	0.19	0.00	0.00	0.00	0.15	0.08	0.00	0.00	0.00	0.15	0.11	0.00	0.00	0.00
Printed Text	0.61	0.70	0.87	0.93	0.95	0.65	0.83	0.97	0.99	1.00	0.63	0.76	0.91	0.96	0.98
Key(s) (whole(s) or nearly)	0.38	0.58	0.85	0.95	0.99	0.78	0.94	1.00	1.00	1.00	0.51	0.72	0.92	0.97	1.00
Blur or multiple categories (no music signs)	0.26	0.32	0.65	1.00	0.00	0.42	0.49	0.39	0.10	0.00	0.32	0.38	0.48	0.18	0.00
Multiple notes or chords	0.28	0.34	0.87	0.00	0.00	0.37	0.53	0.81	0.00	0.00	0.32	0.42	0.84	0.00	0.00
Page border	0.95	0.95	0.98	0.99	1.00	0.74	0.88	0.97	0.99	0.98	0.83	0.91	0.98	0.99	0.99
Other (without musical markings)	0.16	0.20	0.36	0.71	0.00	0.22	0.29	0.53	0.80	0.00	0.18	0.24	0.43	0.75	0.00
Erasure	0.24	0.34	0.56	1.00	-	0.42	0.43	0.35	0.38	-	0.31	0.38	0.43	0.55	-
More categories (with at least one musical sign)	0.32	0.44	0.00	0.00	-	0.18	0.19	0.00	0.00	-	0.23	0.27	0.00	0.00	-
Average	0.35	0.43	0.54	<b>0.60</b>	0.37	0.42	<b>0.50</b>	<b>0.50</b>	0.40	0.37	0.37	0.45	<b>0.51</b>	0.42	0.37
GoogleNet															
Confidence level	Precision					Recall					F1-score				
	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%	0%	25%	50%	75%	90%
Single note (with at least the head)	0.11	0.15	0.00	-	-	0.34	0.31	0.00	-	-	0.17	<b>0.21</b>	0.00	-	-
Manuscript Text	0.43	0.57	1.00	0.00	0.00	0.34	0.41	0.42	0.00	0.00	0.38	0.48	<b>0.59</b>	0.00	0.00
Remaining	0.14	0.18	0.00	0.00	0.00	0.15	0.08	0.00	0.00	0.00	0.15	0.11	0.00	0.00	0.00
Printed Text	0.58	0.69	0.83	0.92	0.95	0.67	0.84	0.96	0.99	1.00	0.62	0.76	0.89	0.95	0.97
Key(s) (whole(s) or nearly)	0.46	0.70	0.92	1.00	1.00	0.79	0.93	0.99	1.00	1.00	0.58	0.80	0.96	1.00	1.00
Blur or multiple categories (no music signs)	0.28	0.32	0.67	1.00	0.00	0.42	0.49	0.36	0.04	0.00	0.33	0.39	0.47	0.07	0.00
Multiple notes or chords	0.27	0.34	0.71	0.00	0.00	0.43	0.64	0.76	0.00	0.00	0.33	0.44	0.73	0.00	0.00
Page border	0.94	0.95	0.98	0.99	0.99	0.78	0.89	0.99	1.00	1.00	0.85	0.92	0.98	0.99	0.99
Other (without musical markings)	0.20	0.24	0.50	0.80	1.00	0.18	0.26	0.56	0.57	1.00	0.19	0.25	0.53	0.67	1.00
Erasure	0.26	0.36	0.61	1.00	0.00	0.44	0.48	0.34	0.20	0.00	0.33	0.42	0.44	0.33	0.00
More categories (with at least one musical sign)	0.28	0.37	0.00	0.00	0.00	0.15	0.11	0.00	0.00	0.00	0.20	0.17	0.00	0.00	0.00
Average	0.36	0.44	0.56	0.57	0.39	0.43	<b>0.50</b>	0.49	0.38	0.40	0.38	0.45	<b>0.51</b>	0.40	0.40

In the multi-class classification task, the deep learning models achieved a balanced accuracy of 43% and an f1-score of 38%. For comparison, random guessing and constant predictors were used as baselines, yielding a balanced accuracy of 9% and an f1-score of 6%.

Upon considering confidence levels, we observed a convex accuracy curve with a peak between 25% and 50% for the three models. Specifically, GoogleNet and ResNet can achieve a balanced accuracy and f1-score of 50% and 51% respectively, by retaining 68% of the data, as depicted in Fig. 7. For confidence levels larger than 0.5, only printed text, keys, and page border can be identified satisfactorily, with f1-scores near to 1. However, a little number of samples are misclassified, leading to a fall in balanced accuracy, which is computed as the arithmetic mean of the per-class recall.

The per-class f1-score values, presented in Table 2, reveal that the models are generally more proficient at predicting the most common classes. Specifically, the classes “Page Border”, “Printed Text”, and “Keys” achieved an f1-score of 85%, 63%, and 61% respectively. These results are particularly beneficial for reducing the number of samples requiring manual annotation for dataset expansion.

The AutoML classifier, while not reaching the accuracy of the neural networks, achieved a balanced accuracy of 37% and an f1-score of 31% in the multi-class case. In the binary classification, it achieved a maximum of 84% and 70% respectively. The resulting models, composed of large ensembles of random forests, gradient boosting, support vector machines, linear and quadratic discriminant analysis, coupled with various pre-processing steps, are described in detail in the notebooks available in the source code repository. These architectures generate large models of several gigabytes. However, existing tools cannot leverage GPU processing like neural network frameworks, making the training and inference of such AutoML models more memory and time-intensive than the neural transfer learning approach.

For all the aforementioned metrics, detailed values and per-class statistics can be found in the notebook in the source code repository.

## 6 CONCLUSIONS

This study presents a comprehensive methodology for OMR applied to historical and handwritten music scores, with a particular focus on the Ricordi Archive. This prestigious archive, housing significant musical manuscripts from eminent opera composers, has been digitized and meticulously annotated to generate a novel dataset of musical symbols. This dataset, along with the models and source code employed in our experiments, is publicly available, thereby contributing to the wider research community.<sup>4</sup>

We have addressed several OMR challenges by training and evaluating multiple neural classifiers to differentiate between these symbols. Three renowned deep learning classifiers, namely ResNet, DenseNet, and GoogleNet, were fine-tuned, and a robust AutoML approach was utilized as a baseline. The deep learning models demonstrated promising results, achieving a balanced accuracy of 85% in the binary classification task. By leveraging the confidence of the models, even higher accuracies were attained.

The primary contribution of this work lies in the creation of a unique dataset of musical symbols derived from real-world annotated manuscripts, which can be utilized to train and evaluate OMR models. Additionally, our work outlines a comprehensive methodology for preprocessing, annotating, and classifying musical symbols, which can be replicated and expanded upon in future research.

Future work will involve using the trained models to annotate additional data, discarding irrelevant sub-images and focusing on images where the model exhibits low confidence. This strategy will enable automatic pixel-wise classification of all pages, followed by a focus on image regions with lower confidence. The ability to identify musical objects will facilitate the provision of more specific labels for such objects. This approach will significantly simplify the annotation of the full corpus, providing the research community with an updated version of the dataset.

## ACKNOWLEDGMENTS

We would like to express our sincere gratitude to the Ricordi Archive for granting permission to utilize the graphical materials included in this paper. In addition, we extend our appreciation for the fruitful mutual collaboration that has

---

<sup>4</sup><https://zenodo.org/doi/10.5281/zenodo.11186095>

taken place over the last 20 years. We gratefully acknowledge the support of NVIDIA Corp. with the donation of two Titan V GPUs.

## REFERENCES

- [1] Arnau Baró, Pau Riba, Jorge Calvo-Zaragoza, and Alicia Fornés. 2019. From Optical Music Recognition to Handwritten Music Recognition: A Baseline. *Pattern Recognit. Lett.* 123 (2019), 1–8. <https://doi.org/10.1016/j.patrec.2019.02.029>
- [2] Jorge Calvo-Zaragoza, Antonio Javier Gallego, and A. Pertusa. 2017. Recognition of Handwritten Music Symbols with Convolutional Neural Codes. *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)* 01 (2017), 691–696. <https://doi.org/10.1109/ICDAR.2017.118>
- [3] Jorge Calvo-Zaragoza, Jan Hajič Jr., and Alexander Pacha. 2020. Understanding Optical Music Recognition. *Comput. Surveys* 53, 4 (July 2020). <https://doi.org/10.1145/3397499>
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>
- [5] Matthias Feurer, Katharina Eggensperger, Stefan Falkner, Marius Lindauer, and Frank Hutter. 2022. Auto-Sklearn 2.0: Hands-Free AutoML via Meta-Learning. *The Journal of Machine Learning Research* 23, 1 (Jan. 2022), 261:11936–261:11996.
- [6] Ichiro Fujinaga. 1997. *Adaptive Optical Music Recognition*. Ph. D. Dissertation.
- [7] Antonio-Javier Gallego and Jorge Calvo-Zaragoza. 2017. Staff-line removal with selectional auto-encoders. *Expert Systems with Applications* 89 (2017), 138 – 148. <https://doi.org/10.1016/j.eswa.2017.07.002>
- [8] Carlos Garrido-Munoz, Antonio Rios-Vila, and Jorge Calvo-Zaragoza. 2022. A Holistic Approach for Image-to-Graph: Application to Optical Music Recognition. *International Journal on Document Analysis and Recognition (IJDAR)* 25 (2022), 293–303. <https://doi.org/10.1007/s10032-022-00417-4>
- [9] Jan Hajič and Pavel Pecina. 2017. The MUSCIMA++ Dataset for Handwritten Optical Music Recognition. *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)* 01 (2017), 39–46. <https://doi.org/10.1109/ICDAR.2017.16>
- [10] Alex Kendall and Yarin Gal. 2017. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf)
- [11] Yixuan Li, Huaping Liu, Qiang Jin, Miaomiao Cai, and Peng Li. 2023. TrOMR:Transformer-Based Polyphonic Optical Music Recognition. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 1–5. <https://doi.org/10.1109/ICASSP49357.2023.10096055>
- [12] David G. Lowe. 2004. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision* 60, 2 (Nov. 2004), 91–110. <https://doi.org/10.1023/b:visi.0000029664.99615.94>
- [13] Sébastien Marcel and Yann Rodriguez. 2010. Torchvision the machine-vision package of torch. In *Proceedings of the 18th ACM International Conference on Multimedia (Firenze, Italy) (MM '10)*. Association for Computing Machinery, New York, NY, USA, 1485–1488. <https://doi.org/10.1145/1873951.1874254>
- [14] Federico Simonetta, Carlos Eduardo Cancino-Chacón, Stavros Ntalampiras, and Gerhard Widmer. 2019. A Convolutional Approach to Melody Line Identification in Symbolic Scores. In *2019 Proceedings of the 20th International Society for Music Information Retrieval Conference (ISMIR)*. 924–931.
- [15] Federico Simonetta, Stavros Ntalampiras, and Federico Avanzini. 2021. Audio-to-Score Alignment Using Deep Automatic Music Transcription. In *2021 IEEE 23rd International Workshop on Multimedia Signal Processing (MMSp)*. 1–6. <https://doi.org/10.1109/MMSp53017.2021.9733531>
- [16] Leslie N Smith and Nicholay Topin. 2019. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, Vol. 11006. SPIE, 369–386.
- [17] Lukas Tuggener, Yvan Putra Satyawan, Alexander Pacha, Jürgen Schmidhuber, and Thilo Stadelmann. 2020. The DeepScoresV2 Dataset and Benchmark for Music Object Detection. *2020 25th International Conference on Pattern Recognition (ICPR) (2020)*, 9188–9195.
- [18] Yusen Zhang, Zhiqing Huang, Yanxin Zhang, and Keyan Ren. 2023. A Detector for Page-Level Handwritten Music Object Recognition Based on Deep Learning. *Neural Computing and Applications* 35 (2023), 9773–9787. <https://doi.org/10.1007/s00521-023-08216-6>