# Classification for Multiformat Object of Cultural Heritage using Deep Learning

Ridwan Andi Kambau
*Faculty Of Computer Science*
*University of Indonesia*
Depok, Indonesia
ridwan.andi@ui.ac.id

Zainal Arifin Hasibuan
*Faculty Of Computer Science*
*University of Indonesia*
Depok, Indonesia
zhasibua@cs.ui.ac.id

M.Octaviano Pratama
*Faculty Of Computer Science*
*University of Indonesia*
Depok, Indonesia
m.octaviano61@ui.ac.id

*Abstract*— The growth of information in the last two decades is dominated by multimedia data such as text, image, audio, and video. Multimedia data with low-level features should be represented in a high-level concept that is easily understood by a human. Classification of the multi-format object is a technique that is used to represent a multi-source object like text, image, audio, and video at once. The object t features are extracted then categorized in several specified classes or concepts. This paper adopts Deep Learning Techniques: (1) Convolutional Neural Networks (CNN) techniques for classifying an image, audio, and video, (2) Recurrent Neural Networks (RNN) technique for classifying text. The experiment uses small data of Indonesian cultural heritage domain. As supervised learning form, the output model is grouped into five classes based on Indonesian ethnic groups (Toraja, Bali, Batak, Dayak, Betawis. The result, this classification model can be implemented in the Multimedia Information Retrieval System and Recommender System for Indonesia Cultural Heritage.

*Keywords*— *classification multi-format Object; convolutional neural network; recurrent neural network; deep learning, Cultural Heritage.*

## I. INTRODUCTION

Classification of the multi-format object appears in many application and a system like a computer vision, object recognition, information retrieval system or recommender system [1]. The idea is to utilize a various source of modalities to improve overall performance that is compared with a single modality mainly using machine learning or deep learning to classify multiple sources. Deep Learning is a state of the art of classification and a subset of machine learning for learning data representation with multiple levels of representation [2]. For classification task, representation on the higher layer amplify important aspect of the input for discrimination and suppress irrelevant variations. Supervised learning is the most common form of deep learning when the model provides desired categories or labels for the result of classification [3]. In the early several single classifications like the image, video, audio, and text classification using deep learning have been tested and has achieved a good performance, but when combining them, the performance slightly decreases. For Image Classification in a single source, Convolutional Neural Network (CNN), one of deep learning model, has been proven can produce good performance [4]. The experiment of 1.2 million images with 1000 class achieved top-1 and top-5 error rates of 37.5% and 17.0% [4]. Another image classification strategy named cross-level that establish a convolutional layer between two nonadjacent levels, aiming to extract more sufficient features with multiple scales at each feature representation level is potential in strategy for the future image classification [5]. Other modal like Audio Classification also uses Convolutional Neural Network technique but in different features. The features used in Audio Classification using deep learning can be done by the high-level feature of Spectrogram [6] [7], then CNN model is performed to the Spectrogram features to categorize into several class [8]. Video Classification studies multiple approaches for extending the connectivity of a CNN in the time domain to take advantage of the local spatiotemporal information. The performance improves the baseline features accuracy. Using 1 million Youtube videos that label 487 classes dataset in empirical evaluation [9]. Recurrent Convolutional Neural Network (RCNN) can perform text classification without human-designed features. The model applies recurrent structure to capture contextual information as far as possible when learning word representations. The experiment is conducted using dataset commonly and proving the method is to outperform the state of the art methods [10].

This research proposes a classification of the multi-format object that is consists of the single modal with deep learning method of image, audio, video, and audio classification to generate unified representation multi-format data with CNN and RNN techniques. Using dataset with object Indonesia Cultural Heritage for each media format in the training phase, the model generates five classes (concepts) contain group ethnic of Batak, Betawi, Bali, Toraja, dan Dayak. This paper is organized as follow; section 2 provides related works, section 3 contains Research Methodology, section 4 is Experiment Settings, section 5 is the Result and Analysis, and section 6 is the conclusion. These research contributions are: (1) we collect Multi-Format Object of Cultural Heritage of 5 tribes in Indonesia Dataset, (2) Multi-Format Classification of Indonesia Cultural Heritage using deep learning model.

## II. RELATED WORKS

### A. Multi-Format Classification

Multi-Format classification utilizes multiple sources (modalities) like texts, images or videos measuring the same instances. Various applications such as multimedia information retrieval system, social media analysis, or recommender system are using Support Vector Machine [11], Deep Boltzmann Machine[3] and Hidden Markov Model [12] techniques for multi-format classification in early. The technique is called shallow learning. Deep learning extends new technique in the multimodal classification with Convolutional Neural Network(CNN) and Recurrent Neural Network (RNN) technique [13]. The new techniques increase the performance and accuracy of the system outperform the shallow learning techniques.

*B. Deep Learning*

Convolutional Neural Network (CNN) [14] and Recurrent Neural Network (RNN) [15] techniques are part of deep learning and become popular in the last five years as a good classification technique. Beside CNN and RNN, there are Deep Belief Network, Deep Boltzmann Machine, and Deep auto-encoder techniques as a part of deep learning. every technique has unique characteristic and function depend on the system requirement. The CNN with various technique has proven as very useful in image classification, and varied techniques are developed sustainably [8]. Not only for image classification but modified CNN also is used for audio and video classification.

Deep Learning is a subset of machine learning that process multiple layers to learn data representation with multiple levels of abstraction [16]. Deep learning makes the breakthrough to solve the old problem that exists in the artificial intelligence community. Deep learning has good performance to find complex structure in high dimensional of data, therefore applicable to many domains of science, government, and business. For the future deep learning will accelerate progress in artificial neural network field [17].

*C. Image and Video Classification*

Before deep learning emergence, some techniques in shallow learning like Support Vector Machine (SVM), K-Nearest Neighbor (KNN) [18], or Naïve Bayes [19] have been widely used. They are a good classifier before the existing of deep learning. Convolutional Neural Network (CNN) as one of the techniques of deep learning is the most popular technique used as an image classifier. This is due to CNN trying to imitate the image recognition system in the human visual cortex[20]. Therefore, it has the ability to cultivate information of the image. A CNN consists of several layers, based on the illustration [14] there are five main layers; Input Layer, Convolutional Layer, Pooling Layer, Fully Connected Layer and Output Layer.
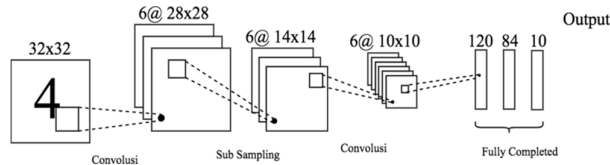

Fig.1. Illustration of Image Classification using CNN

Input layer receives an image dataset that is used as a source or data. The data is processed in the Convolutional Layer. Convolution means applying a kernel on the data image in all offset that allowed. The kernel moves from upper left corner to bottom right. Pooling or subsampling is a reduction process an image size of data. Max pooling is usually used in CNN that works to divide output from convolutional layer become several small grids and take maximum value from each grid to arrange image matrix which has been reduced. The purpose of convolution in image data is to extract features from the input image. Convolution will produce a linear transformation of the input data according to spatial information on the data. The      weights on that layer specify the kernel convolution used, so the kernel convolution can be trained based on inputs on CNN. Furthermore, every neuron in the convolutional layer is transformed to a single dimension of data before being forwarded to the Fully Connected Layer. It caused the data will lose spatial information and not reversible. Hence fully connected layer is only implemented at the end of the layer. Besides the layers, there is activation function or non-linear function that allow CNN is transforming input data to a higher dimension, so it is possible to cut simple hyperplane that allows performing classification. Sigmoid, Tangent Hyperbolic, and Rectified Linear Unit are activation function that used by CNN.

Performing Video Classification is quite tricky because there are three components that have to be concerned; image component on the frame, audio component and temporal or motion component. In this case, we only use image component on the frame by capturing several frames that keep have a significant influence on the video. Further, the image will be performed like Image classification with Convolutional Neural Network. Actually, there are several techniques to classify video like extending connectivity of CNN in the time domain and using shallow learning classifier like SVM, GMM or HMM. For reducing complexity and increase training acceleration we only capture still image from the frames and processing like image classification.

*D. Audio Classification*

In early 2000, Support Vector Machine (SVM), Gaussian Mixture Model (GMM) and Restricted Boltzmann Machine (RBM) become a reference to audio classification. Some audio features like MFCC (Mel Frequency Cepstral Coefficient), Spectrogram, Spectral Flux, and Zero-Crossing Rate is used to classify audio, including sound, speech, and music [21]. The existing of deep learning change the technique reference. Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) as a part of deep learning is a better classifier for audio if we compare with old techniques. Besides Image Classification, CNN also performs audio classification with a spectrogram feature using VGG architecture on CNN [22]. Spectogram features are a visual or 2D representation of spectrum frequency of sound or another signal as they vary with time. Spectrograms are sometimes called sonography, voiceprints, or voicegrams. CNN captures the feature as a 2D visual feature that can be read as an image.
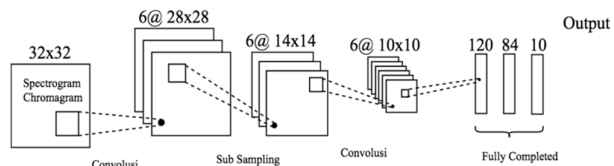

Fig.2. Illustration of Audio Classification using CNN

*E. Text Classification*

Text classification usually relies on many human design features, such as knowledge-based with corpus, dictionary or thesaurus. The problems in text classification are a feature representation like Bag of Word (BoW) model, Latent Semantic Indexing (LSI)[23] or Latent Dirichlet Allocation (LDA) [24]. In deep learning techniques, Recurrent Neural Networks has been proven to be efficient in term of constructing sentence representation. There are several researchers that used RNN as a model for classification like Lai et al. [10].

III. RESEARCH METHODOLOGY

This research starting from creating a multi-format dataset of Indonesia Cultural Heritage, then preprocessing is performed to make dataset can be used in deep learning. Dataset is also separated into training, validation and testing data. Creating a model is done by using Deep Learning (CNN

or RNN) to training data that already split. Validation data is used to measure accuracy and loss when training process perform. After producing an appropriate model, testing data is used to test the model.

### A. Dataset Collection and Preprocessing

Indonesia cultural heritage object with text, image, video, and audio format is used as a dataset. We collected objects of photos, sounds, texts and videos respectively from five group ethnics (Batak, Betawi, Bali, Toraja, and Dayak) manually. Objects of images, sounds, texts, and videos are preprocessed in order to make data feasible to be used in deep learning model. For images data, the data is enriched using data generator with rotating, mirroring and reconstructing current images, then alpha channel images are removed, and the size of images resize into the same size. The sound is also preprocessed by changing raw music file into high-level features Spectrogram. We used techniques from Ramachandran to extract raw sounds into Spectrogram [25]. For text, the dataset is also preprocessed by implementing several techniques like Stemming, Stop-words and Tokenizing words [26]. All of the datasets is separated into 70% as training data, 10% as validation data and 20% as testing data. The dataset also changed into vector format in order to be read by deep learning model.

### B. Model

We used the Convolutional Neural Networks (CNN) architecture for Image classification, audio classification, and video classification, on the other side text classification, is done by using Recurrent Neural Networks (RNN). For image classification, it contains a dataset of five tribes in Indonesia. For each tribe, it contains 100 images respectively. The images then feed-forwarded into CNN which received by the input layer $x = (x_1, x_2, x_3 ..., x_T)^T$. The input layer receives 3 layers of 70x70 image size that will be forwarded into convolutional layer $f(C(x, W))$ with kernel size 5x5. Output from convolutional layer will be forwarded into max-pooling $f(P(x, W))$ layer with kernel size 2x2. For each function f(x), we use whether Rectified Linear Unit (ReLU), Sigmoid or Tangent Hyperbolic (TanH) depends on output from previous layers. In order to make feature-vector in 1D size, fully-connected layer will be the the the performed to produce 1D vector $h_j = f(P(x, W))$, the last step Softmax classification will be performed to produce class prediction $y_j = f(h_j)$. Wea a a also compute loss function using Mean Square Error $MSE$ then $MSE$ will be used in backpropagation algorithm to compute new weight. A new weight $W_{new}$ will be used in the networks and feedforward networks will be repeated as a time until convergence.

$$C(x, W) = \sum_{i=1}^{n} \sum_{j=1}^{m} x_{ij} W_{(i-m, j-n)} \quad (1)$$

$$P(x, W) = \sum_{i=1}^{n} \sum_{j=1}^{m} C(x, W) W_{(i-m, j-n)} \quad (2)$$

$$Sigmoid = \frac{1}{1 + e^{-x}} \quad (3)$$

$$ReLU = \begin{cases} 0, x < 0 \\ x, x \geq 0 \end{cases} \quad (4)$$

$$tanH = \frac{2}{1 + e^{-2x}} - 1 \quad (5)$$

$$Softmax(y_i) = \frac{\exp(h_i)}{\sum_{i=1}^{n} \exp(h_j)} \quad (6)$$

$$MSE = \frac{1}{n} \sum_{j=1}^{n} (y_j - y_o)^2 \quad (7)$$

$$W_{new} = W_{old} \eta - \nabla W \quad (8)$$

This CNN model can be used for image classification and sound classification with different architecture and experimental design. For video classification, the architecture will be the same with image classification but in time framthe e. It means when the video size is 3 second, the number of image wilimages60 frame/secframesiply by 3 second that resulting 180 images. 180 images then will be feed-forwarded into input layer same with image classification. For sound classification also the same with video and image classification but the input is different. Spectrogram as a feature for sound modal, it will be forwarded into CNN architecture then classification is performed to obtain class prediction.

For text classification, RNN model is used to receive text input. Text input is changed into one-hot vector. Then vector $x = (x_1, x_2, x_3 ..., x_T)^T$ is forwarded into architecture wthe hich put in $h_i$ to predict class $y_i$. Loss function is also use MSE like in eq.7 and classification is done by using Softmax like in eq. 6. Activation function is also used for each output from $h_i$ whether ReLU (eq, 4), sigmoid (eq 6) or TanH (eq.5).

$$h_i = RNN(x_i, h_{i-1}) \quad (9)$$

### C. Evaluation

We used various evaluation methods. It started from dataset creation evaluation using Kappa Score (eq. 13) which measure the agreement between two or more annotators. We also use training and validation accuracy to determine training g process is not overfitting. For testing a new or unseen data, we used testing accuracy. The last step of evaluation, we use Precision (Eq.10), Recall (Eq. 11) and F-score (eq. 12) to determine our model can perform better.

$$Precision\ (P) = \frac{|Relevan \cap Retrieved|}{|Relevant|} \quad (10)$$

$$Recall\ (R) = \frac{|Relevan \cap Retrieved|}{|Retreived|} \quad (11)$$

$$F\text{-}Score\ (F) = 2\frac{P \cdot R}{P + R} \quad (12)$$

$$Kappa = \frac{P(A) - P(E)}{P(E) - 1} \quad (13)$$

## IV. RESEARCH EXPERIMENT

We used Cultural Heritage dataset in this research. We used 100 extracted feature of the image, a feature of video, Spectrogram, and Text respectively.

Based on the table below, we use four kinds of features: images data, audios data, videos data and texts data. For each feature, Support Vector Machine model is used as a benchmarking model meanwhile deep learning model CNN and RNN is performed into the dataset. For detail architecture, CNN audio Spectrogram architecture detail can be seen in table 1. Meanwhile, CNN image architecture detail can be seen in table 2.

TABEL 1. CNN AUDIO SPECTROGRAM ARCHITECTURE DETAIL

| Layer | Type | Input | Filter | Kernel Size | Stride | Output |
|---|---|---|---|---|---|---|
| 1 | Input Layer | 3x70x70 | - | | | 3x70x70 |
| 2 | Convolutional | 3x70x70 | 32 | 5x5 | 1 | 32x66x66 |
| 3 | Max-Pooling | 32x66x66 | 32 | 2x2 | 2 | 32x33x33 |
| 4 | Convolutional | 32x33x33 | 32 | 5x5 | 1 | 32x29x29 |
| 5 | Max-Pooling | 32x29x29 | 32 | 2x2 | 2 | 32x14x14 |
| | Convolutional | 32x14x14 | 32 | 5x5 | 1 | 32x10x10 |
| | Max-Pooling | 32x10x10 | 64 | 2x2 | 2 | 64x5x5 |
| 6 | Fully Connected | 64x5x5 | 64 | | | 64 |
| 7 | Softmax | 64 | 1 | | | 1 |

TABEL 2. CNN IMAGE ARCHITECTURE DETAIL

| Layer | Type | Input | Filter | Kernel Size | Stride | Output |
|---|---|---|---|---|---|---|
| 1 | Input Layer | 3x32x32 | - | | | 3x32x32 |
| 2 | Convolutional | 3x32x32 | 32 | 5x5 | 1 | 32x28x28 |
| 3 | Max-Pooling | 32x28x28 | 32 | 2x2 | 2 | 32x14x14 |
| 4 | Convolutional | 32x14x14 | 32 | 5x5 | 1 | 32x10x10 |
| 5 | Max-Pooling | 32x10x10 | 64 | 2x2 | 2 | 64x5x5 |
| 6 | Fully Connected | 64x5x5 | 64 | | | 64 |
| 7 | Softmax | 64 | 1 | | | 1 |

The experiment is tried in GPU computer using Python language and various Python library like Tensor flow, Scikit-learn, Librosa and others. Here detail of research experiments:

TABEL 3. DESIGN EXPERIMENT

| Feature | Model | Dataset |
|---|---|---|
| Images | *Support Vector Machine (Baseline) for Benchmarking Model* | Cultural Heritage Dataset |
| Spectrogram | | |
| Texts | | |
| Images | *Deep learning (CNN and RNN)* | |
| Spectrogram | | |
| Texts | | |

## V. RESULT AND ANALYSIS

### A. Result

We created a dataset of cultural heritage with 100 data for each tribe consists of images, audios, texts, and videos. Most of the image of five tribes consists of traditional houses, most of the audio consists of traditional songs, most of the texts consists of traditional lyrics of songs and most of the video consists of cultural heritage videos. The dataset has been pre-processed and fits with Deep Learning architecture. Dataset is also annotated with 2 annotators. We obtained 0.76 Kappa score which means both annotators are correctly annotated dataset. The sample of Image and Audio dataset can be seen in figure 3.

We have divided the dataset into Training data, Validation data and Testing data. The training process spends 10 epochs for each modal.
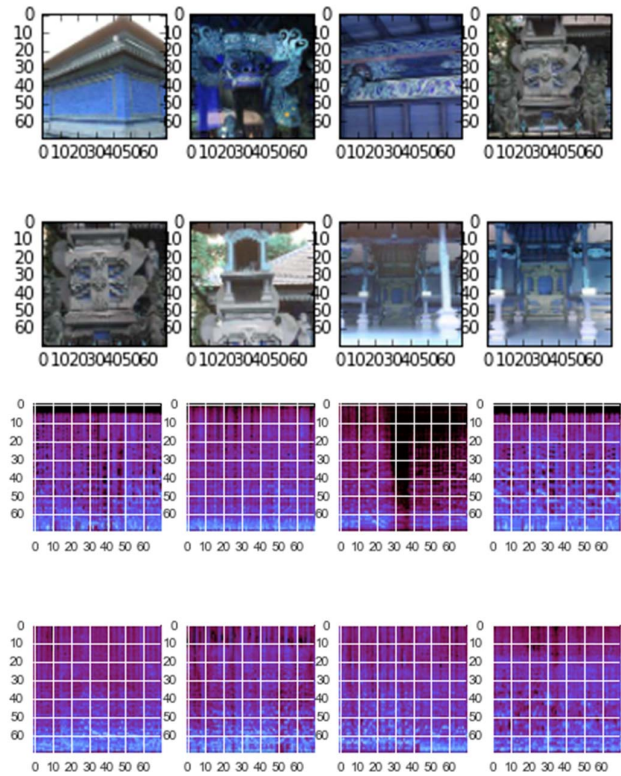


Fig. 3. Image and Audio Dataset Representation

TABEL 4. ACCURACY (A), PRECISION (P), RECALL ®, F-SCORE (F1) MEASUREMENT EVALUATION

| Tribes | Image Classification | | | | | Audio Classification | | | | | Video Classification | | | | | Text Classification | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | P | R | F1 | Size | A | P | R | F1 | Size | A | P | R | F1 | Size | A | P | R | F1 | Size |
| Bali | 0.76 | 0.77 | 0.83 | 0.80 | 20 | 0.57 | 0.62 | 0.69 | 0.67 | 20 | 0.76 | 0.77 | 0.83 | 0.80 | 20 | 0.92 | 0.98 | 0.99 | 0.98 | 20 |
| Batak | 0.76 | 0.95 | 0.85 | 0.78 | 20 | 0.63 | 0.64 | 0.72 | 0.67 | 20 | 0.78 | 0.95 | 0.85 | 0.78 | 20 | 0.95 | 0.97 | 0.99 | 0.98 | 20 |
| Toraja | 0.72 | 0.64 | 0.88 | 0.74 | 20 | 0.67 | 0.60 | 0.64 | 0.58 | 20 | 0.72 | 0.64 | 0.88 | 0.74 | 20 | 0.92 | 0.91 | 0.95 | 0.96 | 20 |
| Dayak | 0.74 | 0.78 | 0.82 | 0.78 | 20 | 0.68 | 0.67 | 0.67 | 0.66 | 20 | 0.78 | 0.77 | 0.87 | 0.78 | 20 | 0.94 | 0.92 | 0.95 | 0.96 | 20 |
| Betawi | 0.70 | 0.79 | 0.85 | 0.76 | 20 | 0.69 | 0.66 | 0.69 | 0.59 | 20 | 0.77 | 0.88 | 0.86 | 0.77 | 20 | 0.95 | 0.95 | 0.92 | 0.93 | 20 |

After the training process performs, we have ensured the process was not overfitting indicating by figure 4 that both training and validation accuracy have almost similar results.
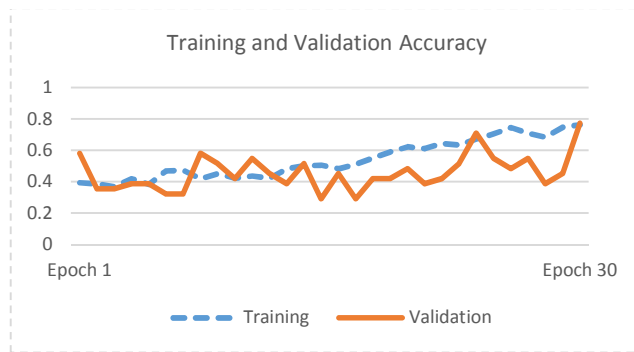


Fig.4. Image Classification Training and validation Accuracy

We can see from figure 4 training and validation accuracy that accuracy increase as time, on the other side loss decrease loss as a time that indicating our model is not overfitting. After creating a model, we use testing data to test our model. The result for image classification, audio classification, video classification and text classification are 77%, 83%, 55%, and 93% respectively. We also use measurement evaluation Precision, Recall, F-score to each model as can be seen in table 4. Comparing with the benchmarking model, Support Vector Machine. The accuracy of the model reaches 55%, 59%, 30% and 77% for image classification, audio classification, video classification and text classification. Deep Learning model performs better than SVM model in Cultural Heritage dataset. The result of classification higher 10 – 30% than SVM.

*B. Analysis*

Cultural Heritage dataset consists of 100 data for each of five Indonesian Tribes (Bali, Batak, Toraja, Dayak, and Betawi). The dataset involves image, audio, video and text data which has already pre-processed. The dataset can be used in deep learning by augmented more data or use data generator. We have shown that our dataset can be used as a benchmark in Cultural Heritage research. The dataset has already correctly annotated that indicated by Kappa score higher than 0.70.

Based on the experiment, Deep Learning model performs better performance comparing with SVM model that indicating by testing accuracy 77%, 83%, 55% and 93% for image classification, audio classification, video classification and text classification respectively. The result is better 10 – 30% than benchmarking model SVM. Deep Learning model performs better because the data could not be separately linear. By performing deep learning model, we have to use a large amount of data by augmenting data generator or other techniques. We have shown that classification multi-format object images, audios, videos and texts of cultural heritage can produce a better performance when using Deep Learning indicating by measurement evaluation.

In table IV shown that when predicting Bali Tribes to Bali Tribes Ground Truth in image classification, 76 images were correctly predicted and merely 24 images was incorrectly predicted. For Audio classification, when predicting Betawi Tribes to Betawi Tribe Ground Truth, 69 audio files were correctly predicted, and merely 31 audio files were incorrectly predicted. The result of table IV indicated our model could perform better in determining tribes of given file.

## VI. CONCLUSION

We created the dataset of multi-format object of Images, Audios, Videos and Texts of Indonesia Cultural Heritage. The Classification of multi-format object by using Deep Learning (CNN and RNN) can produce better performance comparing to SVM with overall accuracy better 10 – 30% than benchmarking model. We have shown our model can predict appropriate tribes given among images, audios, videos or texts of cultural heritage.

Our future work is to integrate the model in Unified Multimedia Information Retrieval or Recommender System for searching and giving a recommendation for Indonesia's cultural heritage object.

## REFERENCES

[1] H. Zhang, V. M. Patel, and R. Chellappa, "Hierarchical Multimodal Metric Learning for Multimodal Classification."

[2] Y. Lecun and Y. Bengio, "Deep learning - review."

[3] N. Srivastava and R. Salakhutdinov, "Multimodal Learning with Deep Boltzmann Machines," *Adv. neural Inf. Process. Syst.*, pp. 2222–2230, 2012.

[4] A. Krizhevsky, I. Sutskever, and H. Geoffrey E., "ImageNet Classification with Deep Convolutional Neural Networks," *Adv. Neural Inf. Process. Syst. 25*, pp. 1–9, 2012.

[5] Y. Liu, B. Yin, J. Yu, and Z. Wang, "Image classification based on convolutional neural networks with cross-level strategy," *Multimed. Tools Appl.*, vol. 76, no. 8, pp. 11065–11079, 2017.

[6] G. Keunwoo Choi and K. C. orgy Fazekas, Mark Sandler, "Convolutional Recurrent Neural Networks for Music Classification," *arXiv Prepr. arXiv1609.04243v3*, pp. 1–5, 2016.

[7] T. Lidy and A. Schindler, "Parallel Convolutional Neural Networks for Music Genre and Mood Classification," *MIREX 2016*, pp. 1–4, 2016.

[8] S. Hershey *et al.*, "CNN ARCHITECTURES FOR LARGE-SCALE AUDIO CLASSIFICATION," *IEEE Int. Conf. Acoust. Speech Signal Process.*, pp. 131–135, 2017.

[9] A. Karpathy, T. Leung, and at all, "Large-scale Video Classification with Convolutional Neural Networks," *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1725–1732, 2014.

[10] S. Lai, L. Xu, K. Liu, and J. Zhao, "Recurrent Convolutional Neural Networks for Text Classification," *Twenty-Ninth AAAI Conf. Artif. Intell.*, pp. 2267–2273, 2015.

[11] A. Fleury, M. Vacher, N. Noury, and S. Member, "SVM-Based Multi-Modal Classification of Activities of Daily Living in Health Smart Homes : Sensors , Algorithms and First Experimental Results," vol. 2, pp. 274–283, 2010.

[12] M. Glodek, S. Scherer, and F. Schwenker, "Conditioned Hidden Markov Model Fusion for Multimodal Classification," no. August, pp. 2269–2272, 2011.

[13] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu, "CNN-RNN : A Unified Framework for Multi-label Image Classification," pp. 2285–2294.

[14] Y. Le Cunn, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based Learning Applied to Document Recognition," *Proceeding IEEE*, 1998.

[15] S. Fern, A. Graves, and T. U. Munich, "Sequence Labelling in Structured Domains with Hierarchical Recurrent Neural Networks," *Proceeding 20th Int. Join Conf. Artif. Intell. IJCAI*, pp. 774–779, 2007.

[16] Y. Lecun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, 2015.

[17] G. Hinton, Y. Bengio, and Y. Lecun, "Breakthrough Deep Learning : machine learning algorithms based on," 2015.

[18] A. Skowron, H. Wang, A. Wojna, and J. Bazan, "Multimodal Classification : Case Studies," pp. 224–239.

[19] S. Mccann and D. G. Lowe, "Local Naive Bayes Nearest Neighbor for Image Classification_TR-2011-11_1.pdf," *IEEE Comput. Vis. Pattern Recognit.*, pp. 3650–3656, 2012.

[20] S. R. Kheradpisheh, M. Ghodrati, M. Ganjtabesh, and T. Masquelier, "Deep Networks Can Resemble Human Feed-forward Vision in Invariant Object Recognition," *Sci. Rep.*, vol. 6, no. August, 2016.

[21] D. Mus, "Audio Feature Extraction with Restricted Boltzmann Machines," *Liacs.Nl*, pp. 1–11.

[22] L. Wyse, "Audio Spectrogram Representations for Processing with Convolutional Neural Networks," vol. 1, no. 1, pp. 37–41, 2017.

[23] S. Zelikovitz and H. Hirsh, "Improving Text Classification with LSI Using Background Knowledge," *Learning*, 2001.

[24] Ş. Taşci and T. Güngör, "LDA-based keyword selection in text categorization," *2009 24th Int. Symp. Comput. Inf. Sci. Isc. 2009*, pp. 230–235, 2009.

[25] A. Ramachandran, S. Vasudevan, and V. Naganathan, "Deep Learning for Music Era Classification," pp. 1–8, 2016.

[26] C. D. Manning, P. Raghavan, and H. Schutze, "An Introduction to Information Retrieval," *Online*, no. c, p. 569, 2009.