
DLk-Sum: Image Collection Summarization Using Deep Learning based Architecture and k-Means Clustering

Anurag Singh , AV Subramanyam

Received: date / Accepted: date

Abstract Image Collection Summarization like video summarization is an important task that machines must be able to perform so as to remove the need for extra manual labor and interference to many tasks. Although image summarization is important it has received relatively less attention as compared to other standard tasks in computer vision and summarization in general. Within the domain of computer vision tasks such as classification, segmentation, object detection have been of keen interest to vision community. Also in domain of multimedia summarization of text, event and video remain hot areas for research. Unlike text, event or video there is no sure temporal sequence between the input data. The images in the data set that needs to be summarized need not be temporally co-related at all. In such scenarios it becomes extremely hard to summarize the datasets/ input data or and best images from a given set. In this paper we aim to build a general model for summarization of datasets and image collections using a deep learning architecture followed by K-means clustering.

Keywords Image Summarization, Auto Encoders , Generative Adversarial Networks, k-Means Clustering, Deep Learning

1 Introduction

With recent rise of data available to us and the rate at which data has been generated. It is also very much evident that there has been rise in its analysis. Organizations and companies wish to gather intelligence and knowledge out of that data. Out of the data generated and the various forms of media such as text, images and videos. One needs to build algorithms that can extract meaning from such varied forms of media. For understanding of text that is natural language a dedicated branch of computer science has emerged called as natural language

Anurag Singh
Netaji Subhas Institute of Technology, University of Delhi New Delhi,
E-mail: anurags.it@nsit.net.in

processing. Similarly, for analysis and understanding of images and videos and developing novel algorithms for tasks such as classification, segmentation, object detection, scene understanding is of primary interest to computer vision research community. As the internet penetration among the world increases there would be more and more communication round the globe and the percent of that data being generated in form of images would increase. In such cases there will be need to extensively analyze images and perform multiple tasks on them. Even the state of art algorithms cannot be trained on every image available as it would be computationally intractable. Therefore, summarization of such large corpus of data is going to be an inevitable challenge in the future. The key challenges are to rigorously define what a summary must and must not contain. How to identify redundancy and ensure that all the parts of source data find a representation in the summarized data. Also, how to qualitatively and quantitatively evaluate a summary i.e distinguish a good and bad summary. Summary of image corpus can be both qualitatively and quantitatively analyzed based on factor of relevance that is how relevant is a particular image to our task for which summary is being built and also the factor of diversity that all the images that are distinct are included in the summary i.e. cover all the aspects of data set and must not contain any redundancy [1]. The models for image collection summarization similar to any other automatic summarization problem can be categorized. Summarization techniques can be categorized into the following: simultaneous [1][2][3] and iterative[4]. We look at all the images simultaneously or their feature vectors at same time and pick some for summary. Or an iterative model which learns about the data set after completely iterating it multiple times. While the simultaneous models are good there is no efficient way to critically analyze large data all at once after certain point[5]. Thus challenging task is to create a model that summarizes by looking at small set of images at once in multiple iterations. Collections are diverse and the order might not have any temporal sequence or correlation within. A way to inherently identify how important a image is in set of images is critical. Immediate applications of automatic image corpus summarization[6] can be to summarize datasets and collection of images in settings where there is not enough communication band width or enough time to review the whole set of generated images manually. Recent applications also include finding the representative set of best images among a burst of images captured during a event. It is also a sub-domain of event summarization . A very important application is while training machine learning algorithms. While video summarization has been well explored for efficient browsing [7] [8], image corpus summarization has received far less attention [9]. We can train machine learning models for different applications using smaller and precise data sets built as a summary of huge data sets. Summarization of data set can help train models without trading-off much on accuracy as the diversity of data will be maintained. Iterative Summarization for images is a challenging task as at any point during the particular iteration even if some sort of temporal data is kept by the model it is not much likely that it can preserve and meaningful information for the model to exploit while it makes the decision to keep or reject the current image as to be part of the summary. Therefore, in this paper we try to explore the image summarization on the idea of looking at all the features at once. This although divides our model into two major components than being ended to end trained but it gives us added leverage to exploit and select images based on the representation of their feature vector maps. Also, it helps us consider all

the images in one go so the information from all the images is utilized in making the choice to keep or discard the image from the representative set. Therefore, the main contribution of the paper is to propose a novel simultaneous approach based image corpus summarization model based on a mix auto-encoder and GAN deep learning architecture to extract the features from the network and the cluster those features to find out the summary by picking up images closest to the centroids of the clusters.

2 Related Works

Summarization in case for videos has received much more recent attention and is most closely related to image summarization than text or event summarization. Lstm based models for supervised and unsupervised learning are being used to capture temporal relationships between video frames. In [8] the authors proposed a deep learning based architecture that can help to summarize different video lengths training in a unsupervised manner. They also propose the idea of using a encoder decoder classifier based network for the task of video summarization. The network consists of a auto encoder and a GAN which share a network component. The decoder of the auto encoder also acts as the generator of the GAN. Overall this combination allows the network to have a loss that aids in gradient flow to a fully connect multi-layer perceptron(MLP) feed forward scorer which scores the confidence of the membership of the frame in summary. The input to the scorer are the features extracted from Inception V3 model. The field of image collection summarization is an important field in information retrieval, machine learning and multimedia processing. The work done by J Camero et. al. [2] gives a multi-kernel clustering approach for the image summarization task. Many of the works use clustering in on way or another [3] [1]. Some other approaches involve graph methods[10] [11],similarity methods[12] or neural network based methods[4] such as self organizing maps. The authors in [1] consider scene summarization which deals with problem of concisely depicting a scene that is creating a visual summary for a given interest point.Images in a summary must not be totally identical to each other.This concept is referred as diversity or dispersion the property is also known as orthogonality. It can also be noted that authors in [1] give a concept of "likelihood" where image representing a set of images in summary must have similarity to that set of images. But it can be noted that the concept of diversity is not only for majority of images in original data set but also applies to minority. With the constraints of size and other parameters diversity and relevance must be maximized such that all the aspects i.e. both minority and majority must be included in summary. In 2013 Yang et. al. [6] formulate image summarization as an optimization problem. The authors apply a dictionary learning approach based of SIFT-Bag of Words model for creating the summary. In their NIPS 2014 paper Sebastian et. al. [9], pose a detailed analysis and try to solve problem of image summarization using sub modular functions.

In addition the authors in [9] also propose a novel evaluation metric which they name V-ROUGE based on recall inspired by ROUGE a evaluation metric extensively used in document summarization community. Recently deep learning based approach has been experimented. In [13] Subramanyam et. al propose a novel end to end deep learning based architecture for iterative summarization of

image corpus. They also proposed new techniques for evaluation of summaries quantitatively using Gini Coefficient and the idea of classification accuracy to test the information capturing capabilities of summary along with Reconstruction Error. They also proposed a brave new idea of considering task specific construction of summaries. The images in summary must be chosen keeping in mind the task that summary is going to be used for. For example a summary needed to train a cat classifier should keep that bias in mind to give more weight to images of cats so that all the hard examples needed to train the classifier are present in the summary. In this paper we also take inspiration from [13] to use a task specific loss in our objective function. The precision and recall metrics define quantitatively how good a summary only when they are provided with user annotated summary. This limits the scope to data sets where not much meta data is available. Secondly this also adds a issue to annotate the ground truth with relevance to a particular task for which summary is desired. As it is possible that which change of task the images needed in summary may change. Therefore a need for standard to compare different algorithms and approaches on particular dataset and fixed set of evaluation metrics needs to be established.

3 Proposed Model

Our network comprises of two main components a Deep learning network that consists of the features being extracted from the images using a encoder CNN (eCNN) and those features are then passed to decoder CNN(dCNN) which also acts a generator for the GAN. We use inception V3[14] for encoding the feature vectors. For the decoding and classification via the discriminator of the GAN(cCNN) which is nothing but a DC-GAN[15]. These feature embedding of images are then clustered using k means where number of clusters that is the $k = \sigma n$. Where the n is the number of images present in the whole dataset and the σ is the fraction images to selected in the summary. The image corresponding to the feature embedding nearest to the centroid of the cluster is selected as the representative image to that cluster to be selected in the summary.

3.1 Problem Formulation

Summarization can be approached in two different ways as a subset selection problem or as an optimization problem. Given a collection of n images $X = (X_1, X_2, \dots, X_n)$ we aim to find a subset S such that $S \subset X$ and $|S| < n$, while preserving the relevance and diversity.

3.2 Learning Framework

In our algorithm, we first extract features of the images $X = \{X_1, X_2, \dots, X_n\}$ using inception v3 [14]. Let these features be $x = \{x_t : t = 1 \dots n\}$. These feature vectors which will be generated after couple of epochs of training are then clustered using the k-means clustering and then the images corresponding to the feature vectors nearest to the centroids of the clusters are picked as summary. This then acts as a

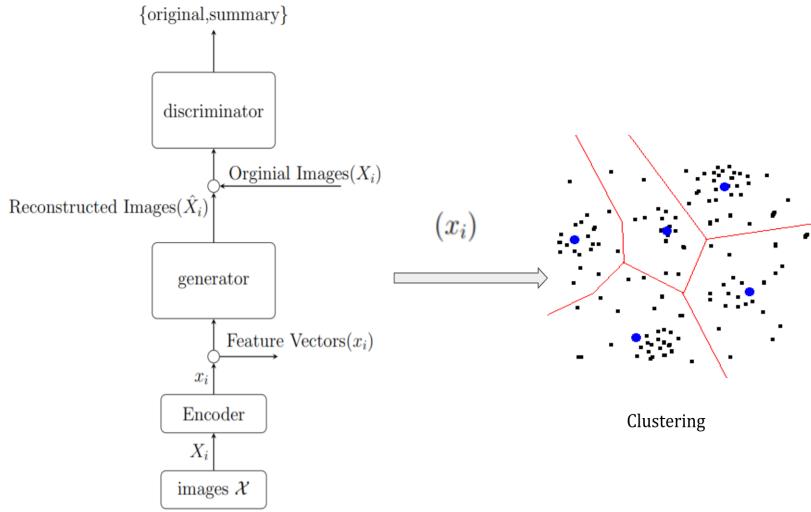


Fig. 1: Main components of our approach: The feature embedding come from the encoder (eCNN). Then each deep feature vector is de-convoluted by the decoder to produce a close estimation of the original image. The reconstructed the image collection is \hat{x} . The discriminator (cCNN) classifies \hat{x} as original or summary class. And decoder and classifier both form the generative adversarial network (GAN). Thus using this model we ensure that most of the important aspects of the image get embedded into the feature vectors. So as when clustering happens it is easier to identify and group similar images based on their feature vector representation and select a representative image for the summary

input to the generator which then reconstructs the image collection as a sequence of images. $\hat{X} = \{\hat{X}_1, \hat{X}_2, \dots, \hat{X}_n\}$.

The discriminator in the GAN is aimed to classify images as two distinct classes. Thus distinguish between images from X and \hat{X} as 'Original' and 'Summary'. The discriminator(cCNN) can be thought as a estimator of distance between X and \hat{X} and assigning different class labels to both of them if they are distinguishable. Therefore, it can be concluded that discriminator acts as a means to represent the error between the original image collection and reconstructed collection from the summary.GAN in the model is Similar to generative adversarial networks presented in DC-GANs and face-GANs[15].Generator and discriminator are trained adversarially until the discriminator is not able to discriminate between the re-collection and original collection. A important question to be addressed is that why are the images with feature vectors nearest to centroids being picked as a part of summary. Without loss of generality while summarizing the dataset in general we do not know what is the perfect length of the summary that should be kept which maximizes the diversity and has least redundancy. So say for a given $k = \sigma \times n$ we pick the image nearest to centroid. Now for a cluster the average distance of the centroid from every other image is going to be minimum. Therefore the centroid will serve as the best approximation for every image. Hence we pick

the image nearest to centroid as the representative image of that cluster. Here k is a hyper parameter that needs to be searched for by observing the average distance of elements in the cluster to their centroids [16][17].

3.3 Training the model

We discuss the different loss functions and training part of the algorithm in this section. The parameters of model are w_e, w_d, w_c for the encoder, decoder i.e. generator and classifier i.e. discriminator. The training of our model is defined by following losses Loss of GAN \mathcal{L}_{GAN} . Reconstruction loss $\mathcal{L}_{reconstruct}$. Similar to usual training GAN models in adversarial manner. The objective is iteratively achieved by:

1. for learning $\{w_e\}$, minimize $\mathcal{L}_{reconstruct}$
2. for learning $\{w_d\}$, minimize $\mathcal{L}_{GAN} + \mathcal{L}_{reconstruct}$
3. for learning $\{w_c\}$, maximize \mathcal{L}_{GAN}

3.4 Reconstruction Loss $\mathcal{L}_{reconstruct}$

$\mathcal{L}_{reconstruct}$ is used to make a summary that captures all relevant frames. It has been used in many of previous works in both image and video summarization [18][8][13]. If original all set of images can be reconstructed using the their deep features respectively it would mean that the network has been able to reduce the dimensionality of the problem without sacrificing much on the information present in the image. Then summary can be considered to have been generated on a tractable set of feature vector representations rather than images themselves. While keeping the encoded relevant information and making the clustering robust to negative examples where a mere pixel shift between a pair of similar images would have landed them far in hyper-dimension space while clustering based on their euclidean distances. $\mathcal{L}_{reconstruct} = \|\sum_{t=1}^n (\mathbf{X}_t - \hat{\mathbf{X}}_t)\|$ Where X_t is image in dataset, \hat{X}_t is reconstructed from the generator and n is number of images in dataset.

3.5 Loss of GAN \mathcal{L}_{GAN}

Similar to [19] we train the classifier i.e. discriminator such that it is able to distinguish between the 'original' X and 'summary' \hat{X} . The \mathcal{L}_{GAN} is thus defined as:

$$\mathcal{L}_{GAN} = \log(cCNN(X)) + \log(1 - cCNN(\hat{X})) \quad (1)$$

where $cCNN(\cdot)$ is the soft-max output of discriminator.

3.6 Task Specific Loss $\mathcal{L}_{task-specific}$

In [13] come up with a new approach for the task specific summarization and propose a method to account for the end goal for which the summary is going

to be used. In this paper we also use the task specific loss implemented for the classification task.

$$\mathcal{L}_{task-specific} = \frac{\mathcal{L}_{pre-trained}(\mathcal{X})}{\beta} \quad (2)$$

Here the standard classification task has been taken. So the images are passed from a trained network for classification task which also shares its weights with the encoder (eCNN) being the same network in case of our implementation (Inception V3)[14]. Therefore in the particular implementation of the task specific loss the eCNN part acts as a multi-task network that is both encoding of the image and its classification[20]. Multi-task learning has shown better results for training in general[21] and also for esoteric tasks such as person re-identification[22]. It is also to note that for a particular task the β is likely to control the degree to which the outliers are binned together with the inliers. Training for a particular task will ensure that the images that are labelled as belonging to same class have similar features. Now as for mentioned in the earlier discussion that a suitable value of k can be chosen as to further tune the number of outliers being present in the summary.

Algorithm 1 Training the model

```

1: function UPDATE PARAMS     $\triangleright$  where input is the feature vector sequence and output is
   learned parameters  $w_e, w_d, w_c$ 
2:   for max number of iterations do
3:      $X \leftarrow$  Mini Batch Of Images
4:      $x \leftarrow eCNN(X)$                        $\triangleright$  Encoding of the image into deep feature vector
5:      $\hat{X} \leftarrow dCNN(x)$                        $\triangleright$  Reconstruction
6:      $\{w_e\}^- = \nabla(\mathcal{L}_{reconstruct} + \mathcal{L}_{task-specific})$ 
7:      $\{w_g\}^- = \nabla(\mathcal{L}_{reconstruct} + \mathcal{L}_{GAN})$ 
8:      $\{w_d\}^+ = \nabla(\mathcal{L}_{GAN})$                        $\triangleright$  Maximization Update

```

4 Data set

The approach is evaluated using following datasets: CIFAR-10 [23], CIFAR-100 [23] Animals with attributes 2 (AwA2) [24], VOC2012 [25] and diversity 2016 [26]. CIFAR-10 consists of 60,000 32X32 tiny images belonging to 10 classes with same images per class. There are 50,000 training and 10,000 test images. CIFAR-100 consists of similar 60,000 32x32 tiny images with 600 images per class. The classes are divided into 20 super-classes with 5 classes per super-class. AwA2 is another data set used for classification purposes. It contains 37322 images of 50 animal classes. Visual Object Classes (VOC 2012) is another image classification data set with 20 classes and 11,530 images. While diversity 2016 contains the images with corresponding ground truth images for task of diversity in image retrieval. Images are ranked according to their importance within a class in ground truth annotations. There are 20821 images of multiple classes with each class containing 300 images. The classes correspond to events such as balloon festival, Buckingham guard change, Diwali or sports like surfing etc.

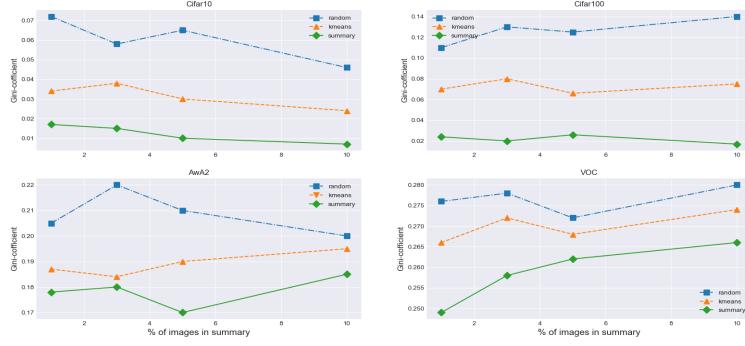


Fig. 2: The gini coefficient for different number of images in summary of multiple datasets. With images selected at random from data set, selected using k-means clustering and images selected in summary from the model

σ	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
precision	10	17.9	26.8	28.28	21.8	16.75	13.1	10.27	8.18	5.37
recall	6.58	9.7	28.5	37.1	50.2	61.1	69.8	81.1	90.8	100
F-score	8.1	13.8	27.3	32.1	30.4	26.3	22.2	18.1	15.0	10.2

Table 1: F-Scores vs σ for the diversity 2016 dataset

5 Results

Regarding evaluation metrics there have been multiple attempts to understand summaries both quantitatively and qualitatively. Still there exists a need for gold standard both in terms of data set with annotated summaries and evaluation of summaries generated for data sets with only meta data being image labels. Like previous works in video summarization [27] [8] where key frame annotations are given.

$$\text{Precision} = \frac{\text{number of images in summary and ground truth}}{\text{number of images in summary}} \quad (3)$$

$$\text{Recall} = \frac{\text{number of images in summary and ground truth}}{\text{number of images in ground truth}} \quad (4)$$

Precision and Recall can be used for evaluation of summaries of data sets with ground truths. precision is ratio of number of correct classifications to summary length and recall is ratio of number of correct classifications to size of ground truth. The F-score is harmonic mean of two. We also taking inspiration from the recent works in image summarization[13] where novel methods for evaluation for image corpus summarization were used when data without ground truth is summarized use Gini- Coefficient and Classification Accuracy for our evaluation purposes. To measure the diversity of our summary on data sets with only meta data available as image labels we use gini-index[28]. A metric often used in economics to define diversity of income levels in a country. In the F-scores are plotted vs the sparsity

Accuracy				
% of Data set used	100%	10%	30%	50%
CIFAR10	88.1%	78%	83.8%	85.1%
CIFAR100	66.2%	55.4%	60.6%	65.4%
VOC2012	80.2%	76.3%	78.6%	79.1%
AwA2	92.9%	89.5%	90.4%	91.9%

Table 2: Classification accuracy for different datasets

loss hyper-parameter σ and they peak out in range $0.3 < \sigma < 0.5$. The scores drop sharply outside this range as σ tends to 0 or 1. The precision value goes down with σ going from 0 to 1 because the size of summary increases and the recall goes up for same reasons. Now for the data sets only with labels as meta data we try to plot gini-indexes as a comparison between randomly picked images from data set, images picked corresponding to centers of K-means image clusters and gini index for images from summary are plotted in Fig:2 for different data sets. There is observed difference that our model is more diverse and shows lower gini index than the summary from clustering images directly. The summary is more diverse throughout the number of images present in summary with multiple data sets as lower the value of gini index more is diversity.

In order to evaluate whether the summary is a good representation of original, we perform the following experiment. We first fine-tune an inception v3 model on original datasets and compute the accuracy. We report these results in Table 2. Since our goal is to test the goodness of summary, we only run the model for few epochs attaining a decent accuracy, though stat of art accuracy may be achieved by using more epochs. Further, we repeat the experiment with summary and summary with data augmentation. We observe that the accuracy achieved using summary itself is good when compared to its original counterpart. In addition, with data augmentation, the accuracy boost up and is close to original case. These results are of great significance because the training using summary uses only 10% of original data, while the trade-off in accuracy is reasonable. Thus it is evident that the summary captures most of the aspects of the dataset.

5.1 Relevance and Diversity Trade-off

In table 3, we show the number of images of each class in the VOC2012 dataset, and their respective ratios to the total number of images in the dataset. In Table 3 the summary of the VOC2012 dataset with a $\sigma = 0.05$ and $\beta = 1.4$ is given. The value of β was chosen empirically as the one which gave the least reconstruction error and the highest classification accuracy, and thus, can be said to have a balance between the number of outliers and inliers in the summary. It is evident from the table that our model trains in a way so that the percentage of images selected from each class is inline with the number of images of that class present in the complete dataset. For example, from table 3 we see that 2.58% of the dataset contains images of the class bicycle, and in the summary 2.41% of images are of the class bicycle, which is quite close to the initial composition. Similarly, for the person class, 40.02% images in the dataset and 41.01% images in the summary are present. Thus, we can say that the relevant information in an image corpus

VOC full dataset	Class	bicycle	motorbike	bird	airplane	horse
	No. of images	390	395	669	604	421
	Ratio	0.0258	0.0262	0.0443	0.0400	0.0279
	No. of outliers	143	130	26	24	140
	Ratio of outliers	0.0421	0.0382	0.0076	0.0071	0.041
	Class	car	tv-monitor	train	chair	cat
	No. of images	700	376	470	491	946
	Ratio	0.0463	0.0249	0.0311	0.0325	0.0626
	No. of outliers	156	122	51	281	55
	Ratio of outliers	0.0459	0.0359	0.0150	0.0826	0.016
Summary $\sigma = 0.05$ $\beta = 1.4$	Class	cow	boat	dog	potted-plant	sheep
	No. of images	293	399	1069	234	303
	Ratio	0.0194	0.0264	0.0708	0.0155	0.0201
	No. of outliers	59	31	112	83	41
	Ratio of outliers	0.0174	0.0091	0.0329	0.0244	0.0121
	Class	bottle	sofa	bus	dining-table	person
	No. of images	342	337	360	261	6044
	Ratio	0.0226	0.0223	0.0238	0.0173	0.4002
	Number of outliers	149	215	39	169	1374
	Ratio of outliers	0.0438	0.0632	0.0115	0.0497	0.4041
Summary $\sigma = 0.05$ $\beta = 0.5$	Class	bicycle	motorbike	bird	airplane	horse
	No. of images	19	21	40	24	19
	Ratio	0.0241	0.0266	0.0506	0.0304	0.0241
	Class	car	tv-monitor	train	chair	cat
	No. of images	36	22	34	27	31
	Ratio	0.0456	0.0278	0.0430	0.0342	0.0392
	Class	cow	boat	dog	potted-plant	sheep
	No. of images	19	20	62	12	12
	Ratio	0.0241	0.0253	0.0785	0.0152	0.0152
	Class	bottle	sofa	bus	dining-table	person
	No. of images	14	16	16	22	324
	Ratio	0.0177	0.0203	0.0203	0.0278	0.4101

Table 3: The number of images of each class for the VOC full dataset (15104 images) and their proportion in dataset. Also the Outliers (3400 images) on loss threshold of 0.1, and their ratio. Similarly Number of images of a class in summary (790 images) at $\sigma = 0.05$ and $\beta = 1.4$ and their ratio to the total number of images in summary. Reconstruction Loss = 0.410. Classification Accuracy = 68.36. summary at $\sigma = 0.05$ and $\beta = 0.5$, and their ratio. Reconstruction Loss = 0.452. Classification Accuracy = 65.85

is maintained by our model by maintaining the composition of different classes in the summary. Thus, we can say that by minimizing the sparsity loss, regularize the composition of images of different classes in the summary and by minimizing the reconstruction loss, we ensure the most relevant images are included in this composition.

Outliers are images which are difficult to classify. There could be many reasons behind an image being an outlier, one being, its diversity. For example, consider an image containing a number of pets, say, a cat, a dog, a rabbit, and a horse. Now, no matter which one of the four labels is given to such an image, it will be an outlier, as the cross-entropy loss of such an image would be high because the probability that it belongs any of the remaining classes would be significant. Thus, we say that in order to ensure diversity, outliers must be present in the summary. In our paper, we show how task-specific loss is responsible for incorporating the outliers in the summary, and the affects of varying β on the number of outliers included in the summary.

For table 3, we calculate the number of outliers in each class of the VOC2012 dataset using a loss threshold of 0.1 i.e. all images having a cross-entropy loss greater than 0.1 were considered to be outliers. A truly diverse summary must have a composition similar to this table, mostly containing the outliers of different classes. the values in table 3 are calculated in the same way as table 3, except this time, β is taken to be a small value equal to 0.5. In the paper, we show that small values of β ensure the presence of a greater number of outliers in the summary. Here, we see that how small values of β lead to loss of relevant information from the summary by scaling up the affect of the task-specific loss, and leading the composition of the summary to be similar to the one only made up of outliers. For example, consider the class, dining-table. It's balanced composition, the one with the most appropriate number of outliers and inliers, in table 3, shows that 2.78% of the summary must comprise of the images of this class. But, when the β is lowered, its composition from table 3, shows that 4.27% of summary contains the images of this class. Thus, we see that by taking a significant affect of the task-specific with a lower value of β leads the composition of the summary to be similar to the composed of only outliers, as in table 3. Therefore, there is a trade-off between relevance and diversity in the summary.

5.2 t-SNE Visualization

We give the t-SNE visualization plots for different experiments conducted on the AWA2, VOC2012 and diversity 2016 data sets. For Diversity 2016, we generate image-embedded t-SNE plots for ground truth with top 20 ranked images from each class i.e. 1400 images in Figure 3(a) and top 50 ranked images from each class i.e. 3500 images in Figure 3(c). We also give the corresponding set of 1400 and 3500 images generated by our model in Figure 3(b) and Figure 3(d). For these figures, we have highlighted some of the images which were selected by our model, and were also present in the ground truth summary. Further, t-SNE scatter plots as in figure 4(a), 4(b), 4(a) and 4(b) show how the summary generated from our models in more sparse and avoids clusters, as compared to the randomly generated summary.

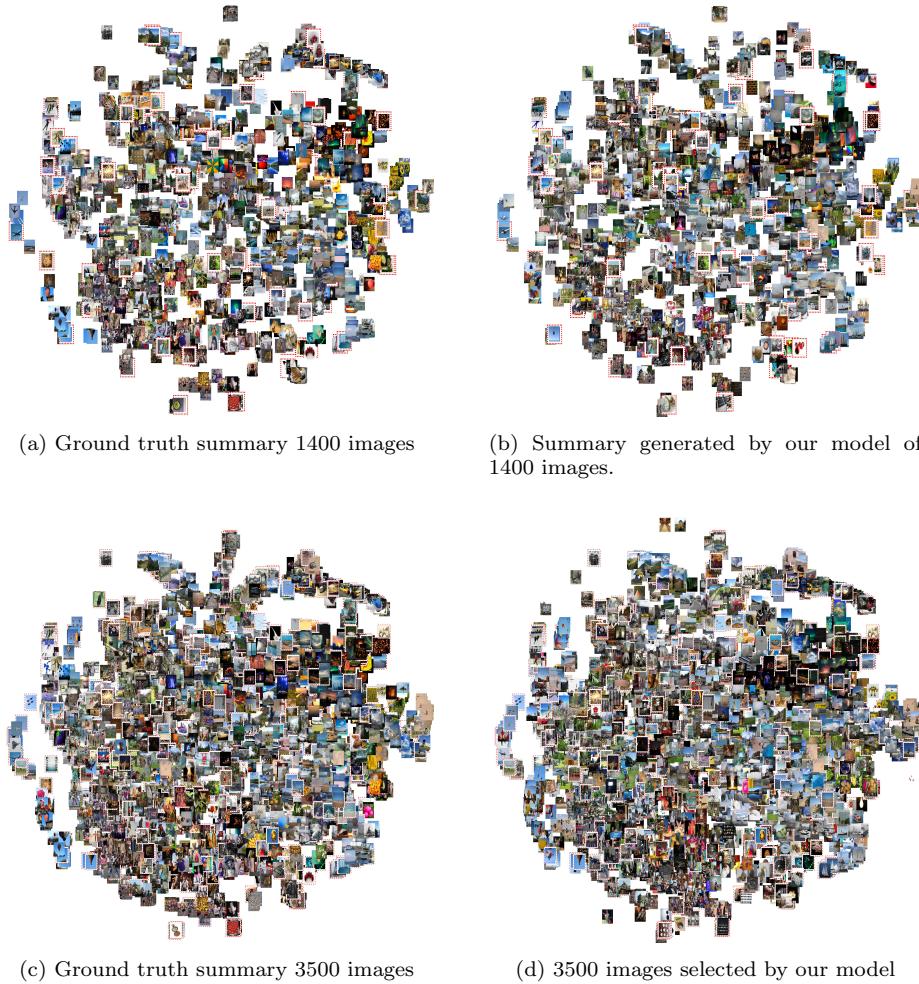


Fig. 3: Karaparthy style t-SNE plot for Diversity2016 of our model in comparison to ground truth. The red boundaries highlight the images which are common to both summary and ground truth. Top-50 and Top-20 images selected from each of 70 classes to make ground truth of 3500 and 1400 images respectively. Please zoom in for better visualization

5.3 Qualitative Visualization

For qualitative assessment of model proposed in this paper we did a comparison on diversity dataset where images were ranked according to their relevancy and diversity in a search result for a particular keyword which in turn was the class of those images. For every class we chose to keep the number of clusters 6. So the $K = 6 \times 70$ that is number of clusters in each class times the number of classes. In Figure 6 we show the top-6 results in the diversity dataset vs the 6 images selected from clustering by our model.

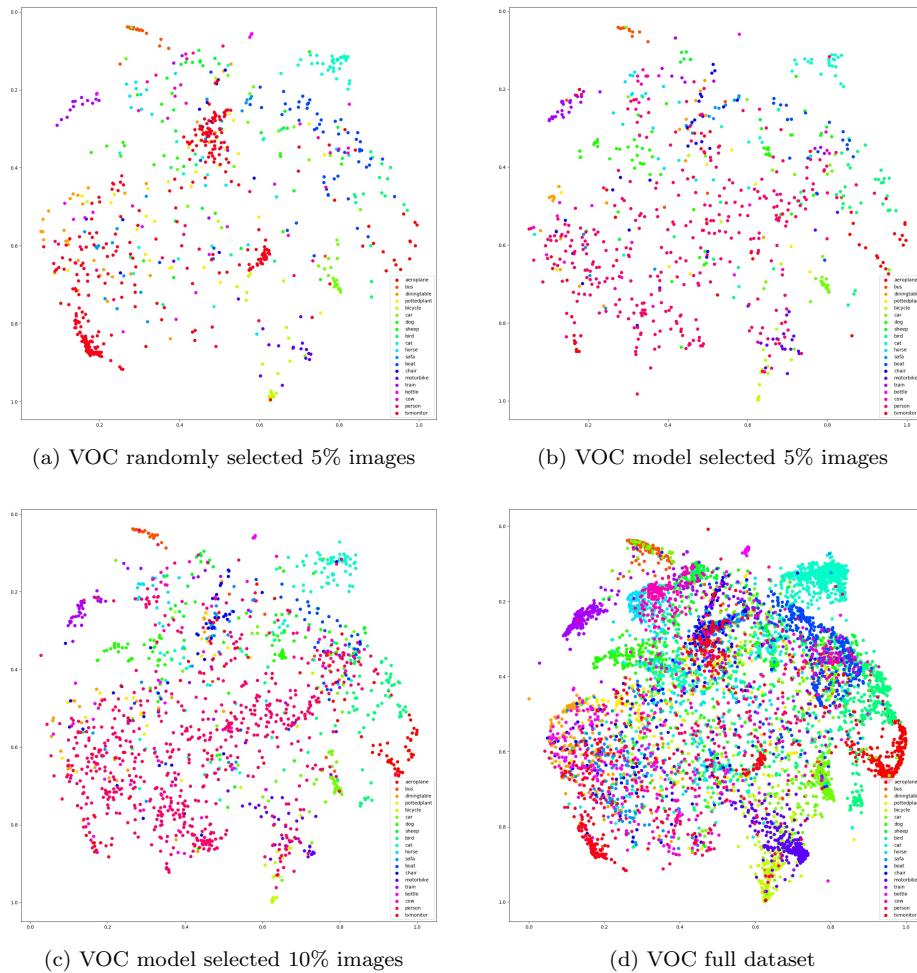


Fig. 4: The above figure talks about the t-SNE plots for VOC data set Fig:(a) is a randomly generated summary for VOC dataset at 5%. Fig(b) and Fig(c) is the summary for VOC dataset, 5% and 10% of images selected from dataset by our model respectively. Fig:(d) is t-SNE for full voc dataset. Different colors represent different classes. Please zoom in for better visualization

There is intersection between images in the ground truth and the images selected by summary. Due to lack of a standard dataset for task of generic summarization we thought of a novel method to qualitatively identify the efficacy of model. The qualitative visualization is done on a single image rather than a dataset where the image is broken into size of 32 pixels (tiny images). These tiny images are then fed forward to the model to generate feature vectors which are clustered by the model. Depending upon the length of the summary the number of tiny images are selected. In Figure 7 we use the famous lenna image commonly used in image processing research[29] to break it into tiny block images of size 32×32 and are

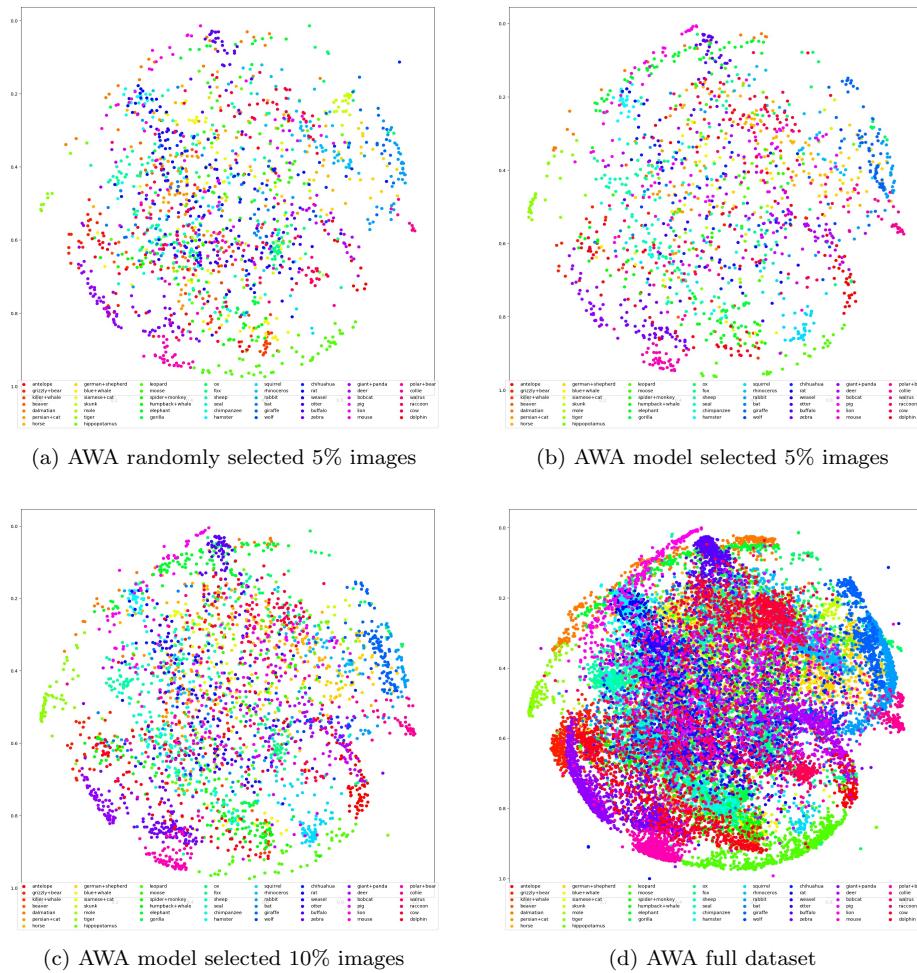


Fig. 5: The above figure talks about the t-SNE plots for AWA data set with Fig:(a) is a randomly generated summary for AWA dataset at 5%. Fig(b) and Fig(c) is the summary for dataset at 5% and 10% of images selected from dataset by our model. The Fig:(d) is t-SNE for full AwA dataset. Different colors represent different classes. Please zoom in for better visualization

clustered for different σ . Within the Figure 7 it can be observed that the model tends to discard the background and keeps most relevant and diverse information blocks as the size of summary continues to decrease. The eyes and the blocks containing hair are not so similar in their feature representations as compared to two adjacent blocks in the background. The trade-off between relevance and diversity as discussed suggests that for some values of σ there is minimum redundancy and high diversity both achieved. Decreasing the σ further will cause the important parts of the dataset/image to be discarded and keeping the σ too high would mean keeping redundancy. The efficacy can be observed by comparison with randomly

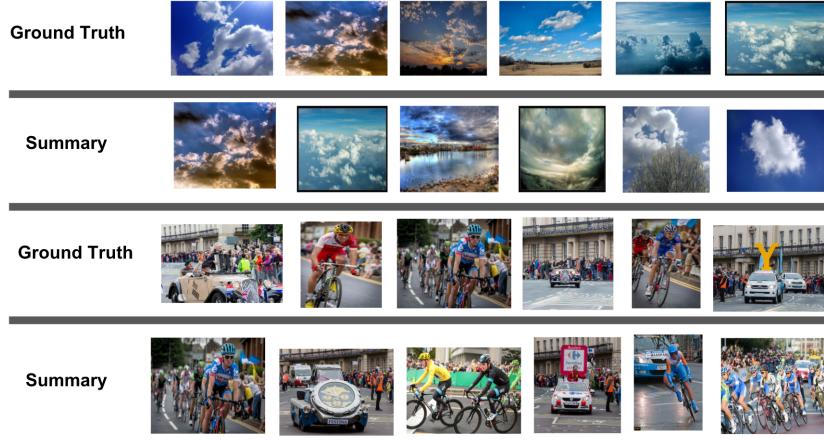


Fig. 6: Comparison of the images selected in our summary vs top images in the top-6 ground truth for diversity 2016 dataset for clouds and Tour-de France class.

picked blocks which do not consider the relevance and diversity of a block. While a randomly generated summary will be most effused because of the random nature but it is likely to neglect the relevance of a block absolutely which can be seen in case of all 10%, 30% and 50% of the blocks picked randomly.

6 Review of Task Specific Loss $\mathcal{L}_{task-specific}$

Task specific loss is important loss that needs to be present in objective function for the generation of task specific summary. Although many diversity regularization losses like Detrimental Point Processes or Repelling Regularization Losses[8] try to achieve diversity. But most of diversity regularization works in unsupervised fashion to help select images with high visual diversity by selecting points that are sufficiently distant in feature space representation. Repelling regularization works in a way that it repels selection of images which are near to each other in the feature space. But there should be a priority in terms of selection of images from feature space where images of multiple labels cluster together in the dataset. Since diversity regularization techniques mentioned above are unsupervised it is less likely that they would take care of such scenarios where a outlier of a particular class rests within the cluster of inliers of another class. For the use case of task specific summarization it is necessary that such outliers do get picked in general along with ensuring similar high visual diversity. As task specific loss leverages labels it is able to distinguish between images from class that are present within cluster of images of other class as outliers. It can be observed in the t-SNE visualization scatter plots for example in Figure 4(b) in the top right corner within selected images of a particular class i.e. cat we can find an outlier image of person being selected. Same can be observed in Figure 5(b) for example in bottom left

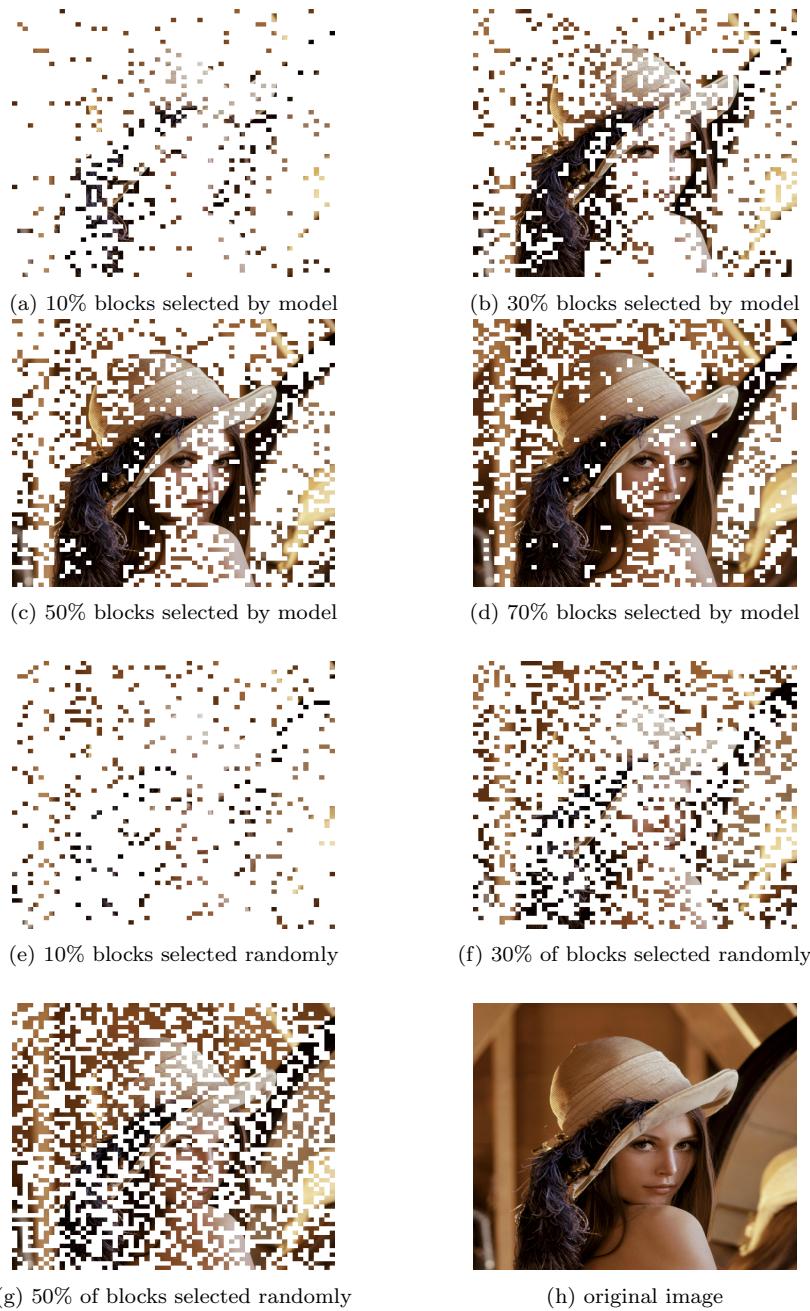


Fig. 7: Summary of the image for different values of σ and comparison with randomly generated summary

corner where images from giant panda are selected outliers from other classes are also selected.

7 Conclusions

In this work, we propose an unsupervised model to summarize a large collection of images. The classification results attained by training a deep network on summary only and on original dataset are close and show similar trend. Thus, model can also be used for a quick analysis of various models without needing to train on entire dataset. Moreover, image summarization can also be used to curate more precise data sets for given tasks. Summarization of data sets can be performed and accuracy achieved on summaries can be used to find out the small scale data sets for quicker training of models. In case the labels are not available, technique can be used to summarize and retains fraction of data, which can be relatively convenient to annotate. Further, one can perform different tasks on this data before scaling up the model as well other processing on the original data.

8 Future Works

There is a need for a standard in the terms of dataset specific summarization both in the terms of general summarization and task specific summarization where the human annotated summarized results are available as a ground truth reference for the solid comparison using F-score metrics. It will also help build better supervised algorithms and also to test the efficacy of the unsupervised models in a more rigorous manner. It was also observed that the idea of image collection summarization can help keep a strong foundation for image compression, saliency detection within images and attention models. The image collection summarization models can be fine tuned and used for the above mentioned tasks. For example similar idea to Figure 7 can perform saliency detection or attention model that it keeps only important parts in a image for small summary lengths i.e low values of σ . Also if the discarded information can be estimated using the remaining summary. The reconstruction would act as the compressed image from lossy compression.

References

1. Simon, I., Snavely, N., Seitz, S.M.: Scene summarization for online image collections. In: Proceedings of the IEEE International Conference on Computer Vision. (2007) 1–8
2. Camargo, J.E., González, F.A.: A multi-class kernel alignment method for image collection summarization. In: Iberoamerican Congress on Pattern Recognition, Springer (2009) 545–552
3. Stan, D., Sethi, I.K.: eid: a system for exploration of image databases. *Information processing & management* **39**(3) (2003) 335–361
4. Deng, D.: Content-based image collection summarization and comparison using self-organizing maps. *Pattern Recognition* **40**(2) (2007) 718–727
5. Chen, C.P., Zhang, C.Y.: Data-intensive applications, challenges, techniques and technologies: A survey on big data. *Information Sciences* **275** (2014) 314–347
6. Yang, C., Shen, J., Peng, J., Fan, J.: Image collection summarization via dictionary learning for sparse representation. *Pattern Recognition* **46**(3) (2013) 948–961

7. Khosla, A., Hamid, R., Lin, C.J., Sundaresan, N.: Large-scale video summarization using web-image priors. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2013) 2698–2705
8. Mahasseni, B., Lam, M., Todorovic, S.: Unsupervised video summarization with adversarial lstm networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 2982–2991
9. Tschiatschek, S., Iyer, R.K., Wei, H., Bilmes, J.A.: Learning mixtures of submodular functions for image collection summarization. In: Advances in neural information processing systems. (2014) 1413–1421
10. Cai, D., He, X., Li, Z., Ma, W.Y., Wen, J.R.: Hierarchical clustering of www image search results using visual, textual and link information. In: Proceedings of the 12th annual ACM international conference on Multimedia, ACM (2004) 952–959
11. Gao, B., Liu, T.Y., Qin, T., Zheng, X., Cheng, Q.S., Ma, W.Y.: Web image clustering by consistent utilization of visual features and surrounding texts. In: Proceedings of the 13th annual ACM international conference on Multimedia, ACM (2005) 112–121
12. Chen, J.Y., Bouman, C.A., Dalton, J.C.: Hierarchical browsing and search of large image databases. *IEEE transactions on Image Processing* **9**(3) (2000) 442–455
13. Virmani L, A.S., Subramanyam, A.: Summary Generation for Image Corpus using Generative Adversarial Networks. (2018)
14. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. *CoRR* **abs/1512.00567** (2015)
15. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in neural information processing systems. (2014) 2672–2680
16. Ketchen, D.J., Shook, C.L.: The application of cluster analysis in strategic management research: an analysis and critique. *Strategic management journal* **17**(6) (1996) 441–458
17. Thorndike, R.L.: Who belongs in the family? *Psychometrika* **18**(4) (1953) 267–276
18. Zhang, K., Chao, W.L., Sha, F., Grauman, K.: Video summarization with long short-term memory. In: Proceedings of the European conference on computer vision. (2016) 766–782
19. Larsen, A.B.L., Sønderby, S.K., Larochelle, H., Winther, O.: Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300* (2015)
20. Argyriou, A., Evgeniou, T., Pontil, M.: Multi-task feature learning. In: Advances in neural information processing systems. (2007) 41–48
21. Ruder, S.: An overview of multi-task learning in deep neural networks. *CoRR* **abs/1706.05098** (2017)
22. Chen, W., Chen, X., Zhang, J., Huang, K.: A multi-task deep network for person re-identification. *CoRR* **abs/1607.05369** (2016)
23. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images. (2009)
24. Xian, Y., Lampert, C.H., Schiele, B., Akata, Z.: Zero-shot learning - A comprehensive evaluation of the good, the bad and the ugly. *CoRR* **abs/1707.00600** (2017)
25. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International Journal of Computer Vision* **88**(2) (June 2010) 303–338
26. Ionescu, B., Gînscă, A.L., Boteanu, B., Lupu, M., Popescu, A., Müller, H.: Div150multi: A social image retrieval result diversification dataset with multi-topic queries. In: Proceedings of the 7th International Conference on Multimedia Systems. (2016) 46:1–46:6
27. Song, Y., Vallmitjana, J., Stent, A., Jaimes, A.: TvsuM: Summarizing web videos using titles. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2015) 5179–5187
28. Gini, C.: Variabilità e mutabilità. Reprinted in *Memorie di metodologica statistica* (Ed. Pizetti E, Salvemini, T). Rome: Libreria Eredi Virgilio Veschi (1912)
29. Fisher, Y.: Fractal image compression: theory and application. Springer Science & Business Media (2012)