

Sequential Person Recognition in Photo Albums with a Recurrent Network

Yao Li¹ Guosheng Lin^{2*} Bohan Zhuang¹ Lingqiao Liu¹ Chunhua Shen¹ Anton van den Hengel¹
¹The University of Adelaide ²Nanyang Technological University

Abstract

Recognizing the identities of people in everyday photos is still a very challenging problem for machine vision, due to issues such as non-frontal faces, changes in clothing, location and lighting. Recent studies have shown that rich relational information between people in the same photo can help in recognizing their identities. In this work, we propose to model the relational information between people as a sequence prediction task. At the core of our work is a novel recurrent network architecture, in which relational information between instances' labels and appearance are modeled jointly. In addition to relational cues, scene context is incorporated in our sequence prediction model with no additional cost. In this sense, our approach is a unified framework for modeling both contextual cues and visual appearance of person instances. Our model is trained end-to-end with a sequence of annotated instances in a photo as inputs, and a sequence of corresponding labels as targets. We demonstrate that this simple but elegant formulation achieves state-of-the-art performance on the newly released People In Photo Albums (PIPA) dataset.

1. Introduction

With the widespread adoption of digital cameras, the number of photos being taken has increased astronomically. The culture surrounding the use of these cameras means that a large proportion of these photos contain people. The overwhelming volume of these images is creating a demand for smart tools to organize photos containing people. One cornerstone step is to recognize each person in these everyday images. Previous work [1, 7, 8, 23, 29, 14, 12, 15] has shown that person recognition in such unconstrained settings remains a challenge problem for machine vision due to various factors, such as non-frontal faces, varying lighting and settings, and even just the variability in the appearance of a face over time.

To tackle these challenges, in addition to the appearance of the face, recent studies [1, 29, 14, 12] have shown that

*G. Lin's contribution was made when he was with The University of Adelaide. Correspondence should be addressed to C. Shen.

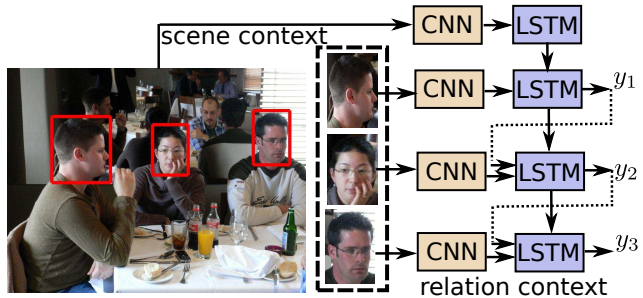


Figure 1. Our approach performs person recognition in photo albums as a sequence prediction task. Both contextual cues and instances' visual appearance are exploited in a unified framework.

contextual cues can help in recognizing peoples' identities in everyday photos. Other features of an individual, such as clothing [1, 7], may also provide valuable cues. The relationships between the person to be recognized and others, can also be a vital cue [23, 4, 12]. To take advantage of different relation cues, probabilistic graphical models have been widely exploited [1, 23, 4, 12].

In this work, we propose to model the rich relations between people in an image as a sequence prediction task. This is much motivated from the success of sequence prediction formulations in modeling relations between words in language problems [18, 22].

In our work, we propose a novel Recurrent Neural Network (RNN) architecture for the sequence prediction task (see Fig. 1), which consists of a Convolutional Neural Network (CNN) and a Long Short Term Memory (LSTM) network [10]. LSTMs have shown impressive performance on several sequence prediction problems such as image captioning [22], video description [6, 20], multi-label image classification [24], machine translation [18], *etc.* Initially we feed the LSTM with global image information provided by the CNN. At each subsequent step, the input to the LSTM is a joint embedding of the CNN feature representation of the current person instance and its predicted label at the last step. The LSTM then predicts the identity label for this person instance. In this sense, the number of steps in our sequence model is of variable length, depending on the number of annotated instances in the image.



Figure 2. Person recognition in photo albums.

Two sources of contextual cues are exploited in our model, including relation context and scene context (see Fig. 2). The relation context refers to relational information between multiple people in the same image (*e.g.*, some people are likely to appear together), which is naturally incorporated by our sequence prediction formulation. Based on the assumption that some people are more likely to appear in a certain scene, the scene itself can be used as a prior to indicate which identities tend to appear. This cue is exploited in our model by feeding the global image feature to our sequence prediction model at the initial step, which informs the system about the scene content. As we will show in the experiments, both contextual cues are critical to the methods ability to achieve state-of-the-art performance.

To our knowledge, this is the first approach to formulate person recognition in photo albums as a sequence prediction task. This simple but elegant approach (a) enables modeling visual appearance and contextual cues in the same framework, (b) handles a variable number of instances in an image and (c) is end-to-end trainable. We demonstrate that our model achieves state-of-the-art performance on the *People In Photo Albums (PIPA)* dataset [29], the benchmark dataset for person recognition in photo albums.

2. Related Work

Person recognition in photo albums. Person recognition in photo albums [1, 7, 23, 29, 14, 12] aims to recognize the identities of people in everyday photos. Intuitively, the face region can be an important cue for the task, however, it may not entirely be reliable, as in this unconstrained setting, people can have non-frontal, or even back views. This makes the problem much harder than the classical face recognition.

Recent studies on this topic have been boosted by the introduction of the PIPA dataset [29]. In this paper [29], the authors proposed a method which combines information from three sources, including the full body, poselets [3] and the DeepFace [19] model. Oh *et al.* [14] evaluated the

importance of different cues for the task, such as different body regions, scene and human attributes. More recently, Li *et al.* [12] proposed to incorporate contextual cues to the task, including group-level context and person-level context. However, the contextual models of [12] are treated as the post-processing steps after the classification result, we instead exploit contextual cues and visual appearance in a unified framework. This is particularly significant for the problem at hand because there are many possible identities for each detection, and non-frontal face cues can be extremely inconclusive. A unified framework means that the cues are used collectively to exploit all of the available information, and will succeed over the greedy approach when the face recognition result is ambiguous.

Our approach is also related to work on identifying people in group photos [8, 23, 4, 13], as our LSTM framework naturally handles multiple people in the same photo.

Modeling dependencies with RNNs. RNNs, and particularly LSTMs, have enjoyed much popularity recently in sequence modeling problems, largely due to their ability to model dependencies within sequences. For instance, LSTMs have been widely used in machine translation [18] and vision-to-language problems, such as image captioning [22, 27], video description [6, 20], and visual question answering [28, 26, 25].

In comparison with the sequence prediction model in machine translation [18], which contains both an encoder LSTM and decoder LSTM, we model both visual features and contextual cues using a single LSTM. In this sense, our model is closer to those used in image captioning [22] as the LSTM output is sent to a classification layer at each step (except the initial step), as is the case in image captioning. However, the main difference with respect to image captioning models is that we have a visual feature input to the LSTM at every step (not only the initial step).

Although the identities of a group of people in a photo are better described as a set than a sequence, there is an obvious dependence between them. It is this dependence that we seek to capture here using an RNN. Despite being very popular for the task, the RNN model is not sequence-specific, but rather, it can be employed to model a sequence by feeding the previous element of the sequence in as input. What the RNN actually generates is an output conditioned on its internal state and input. The output can be interpreted as a sequence, but it can equally be considered a set, or a variety of other types. Wang *et al.* [24], for instance, exploit LSTMs to model the dependencies between the multiple tags that different users apply to the same image (despite their being no natural order to the labels). Stewart *et al.* [17] similarly use an RNN to model the dependencies between detections in an image, which also have no natural order.

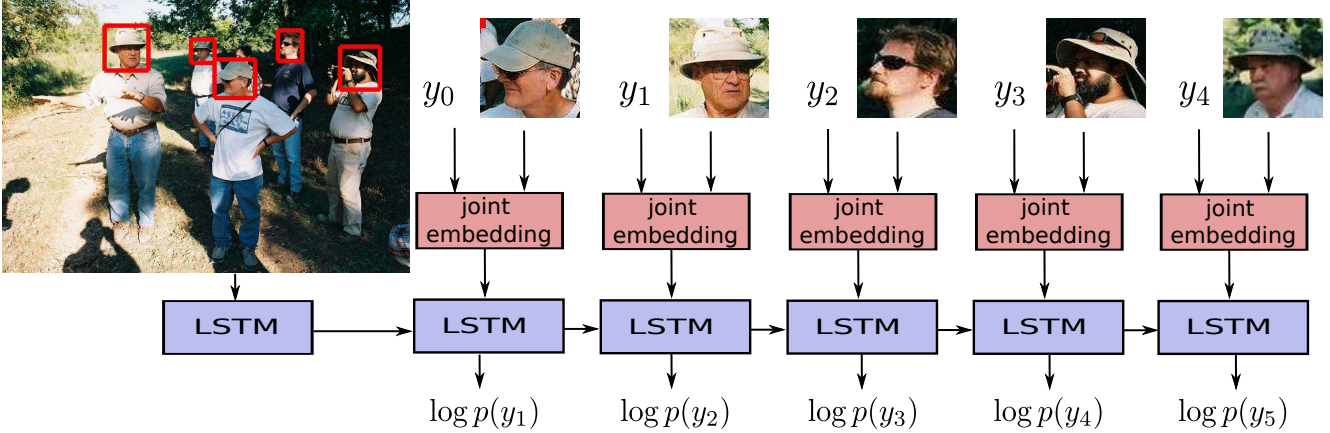


Figure 3. Our sequence prediction approach for recognizing people in photo albums. For an image which may contain multiple people, our approach predicts the identity of each person in a sequence using an LSTM-based framework. The initial state of the LSTM is informed by the scene context, and in each of the subsequent steps, the input to the LSTM is the joint embedding of the label of the last step and visual feature of the current instance (see Fig. 4 for details). The task of the LSTM is then to predict label of the current instance. In this way, the relationships between people are naturally incorporated in our framework. Note that y_0 is the label of an auxiliary identity.

3. Model

We propose a sequence prediction approach for recognizing people in photo albums. As depicted in Fig. 3, at each step (except the first one), we jointly embed the previous label and the current image of an identity, which is then served as the input to an LSTM. The LSTM then tries to predict the current correct label. Our work is largely motivated by the successful application of LSTMs in a range of sequence prediction tasks, such as image captioning [22] and machine translation [18].

More formally, at the training phase, a training sample consists of an image I , a set of annotated bounding boxes $\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_N\}$ of a human body region (e.g., head region) of N person instances in the image, and their corresponding labels $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$. Note that N varies across different images.

In our work, by treating both \mathcal{B} and \mathcal{Y} as sequences with some order (the order for \mathcal{B} and \mathcal{Y} must be same, so an instance in \mathcal{B} is matched with its label in \mathcal{Y}), we aim to look for a set of parameters θ^* which maximizes the log likelihood of producing the correct label sequence \mathcal{Y} given the input sequence \mathcal{B} and the global image I among all the training samples:

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{(\mathcal{B}, I, \mathcal{Y})} \log p(\mathcal{Y} | \mathcal{B}, I; \theta). \quad (1)$$

Intuitively, we can model the joint probability over all the labels $y_{1:N}$ using the chain factorization, that is,

$$\log p(\mathcal{Y} | \mathcal{B}, I; \theta) = \sum_{t=1}^N \log p(y_t | y_{1:t-1}, \mathbf{b}_{1:t}, I; \theta), \quad (2)$$

We assume that the current prediction of y_t does not depend on all instances in \mathcal{B} but only the previously seen instances $\mathbf{b}_{1:t-1}$ and the current instance \mathbf{b}_t .

Analogous to sequence prediction models in other tasks [22, 2], we model the conditional probability $p(y_t | y_{1:t-1}, \mathbf{b}_{1:t}, I; \theta)$ with a recurrent neural network, by introducing a hidden state vector \mathbf{h}_t , that is,

$$p(y_t | y_{1:t-1}, \mathbf{b}_{1:t}, I; \theta) = p(y_t | \mathbf{h}_t; \theta). \quad (3)$$

The hidden state vector \mathbf{h}_t has the following form:

$$\mathbf{h}_t = \begin{cases} f(I; \theta) & \text{if } t = 0, \\ f(\mathbf{h}_{t-1}, \mathbf{x}_t; \theta) & \text{otherwise.} \end{cases} \quad (4)$$

The \mathbf{x}_t in Eq. 4 is the novel part of our RNN architecture, which is a joint embedding of the previous label y_{t-1} and the current input instance \mathbf{b}_t (more details are provided in Sec. 3.1). For $f(\cdot)$ we opt for the LSTM component, which has shown state-of-the-art performance on sequence modeling tasks such as image captioning [22, 2].

As described in Eq. 4 and Fig. 3, at the initial step ($t = 0$), the input to the LSTM is the global image content I , which informs the network with scene context. For this purpose, we use features extracted from a pre-trained Convolutional Neural Network (CNN) to represent images. In subsequent steps the inputs are the current joint embedding \mathbf{x}_t and its previous hidden state \mathbf{h}_{t-1} ¹.

Let z_t denote the output of the LSTM at step t . We then add a fully-connected layer (W) with the softmax function

¹When $t = 1$, as there is no preceding identity labels, we add an auxiliary identity label y_0 . This is similar to the case in image captioning [22] where a special start token is used to represent the beginning of a sentence.

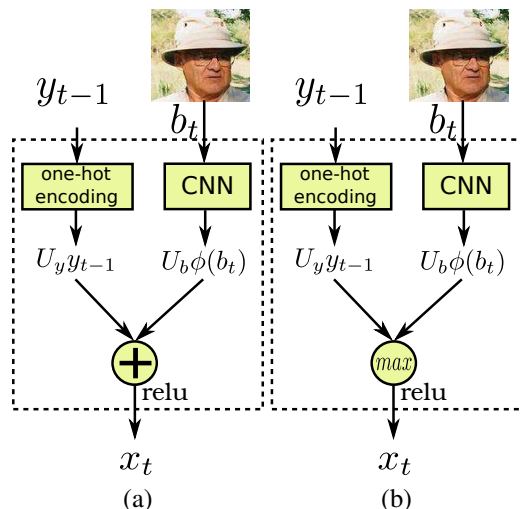


Figure 4. The joint embedding layer for generating the LSTM input \mathbf{x}_t . Either (a) addition or (b) an element-wise maximum can be used for joint embedding.

on the top to generate $p(\mathbf{y})$, the probability distribution over all the identity labels.

Our loss over all the steps is the sum of the negative log likelihood of the ground-truth identity label \mathbf{y}_t at each step:

$$L = - \sum_{t=1}^N \log p(\mathbf{y}_t). \quad (5)$$

The above loss can be minimized by the Back-propagation Through Time (BPTT) technique.

3.1. Joint embedding of instance feature and label

At any step t ($t > 0$), we argue that there are two sources of information which can help to predict the current label. The first source of information is the previous label \mathbf{y}_{t-1} (exploiting label co-occurrence information). The other is the appearance of current instance \mathbf{b}_t , which is a feature vector denoted by $\phi(\mathbf{b}_t)$. Therefore, to generate the LSTM input \mathbf{x}_t , we propose a joint embedding layer which combines valuable information from the two sources. Specifically, as shown in Fig. 4, after transforming labels into one-hot vectors, we define two embedding matrices \mathbf{U}_y and \mathbf{U}_b for encoding \mathbf{y}_{t-1} and $\phi(\mathbf{b}_t)$ respectively:

$$\mathbf{x}_t = \text{relu}(\mathbf{U}_y \mathbf{y}_{t-1} + \mathbf{U}_b \phi(\mathbf{b}_t)), \quad (6)$$

where *relu* stands for the Rectified Linear (ReLU) activation function. This is motivated by the label embedding formulation of [24].

An alternative to the addition above is to take the element-wise maximum with the ReLU activation, that is,

$$\mathbf{x}_t = \text{relu}(\max(\mathbf{U}_y \mathbf{y}_{t-1}, \mathbf{U}_b \phi(\mathbf{b}_t))). \quad (7)$$

The performance of these two formulations are analyzed in the experiment section (Sec. 4.2).

3.2. Training and inference

Random order training. In some sequence prediction models, such as image captioning, there is a natural order of the input sequence (*e.g.*, words in a sentence). For some tasks when the order is not obvious, the order is pre-defined based on some heuristic rules. For instance, in the human pose estimation work of [9], a tree based ordering of joints is used. The investigation of Vinyals *et al.* [21] has shown that for some simple problems most orderings perform equally well.

Our task differs from the sequence prediction tasks above, as people appearing in images don't have an inherent order. We thus opt for random orders at training time. To be more specific, for a training image, its annotated instances and their identity labels are randomly shuffled in the same order to generate an input sequence and a target sequence respectively. Therefore, the order of people in a training image varies in different epochs, which incorporates randomness to the training process.

Inference. At test time, to predict the identity label of a query instance, we generate multiple sequences for this instance, all with the query instance at the end and other instances in the image randomly ordered. The rationale behind this is that in order to take the advantage of the rich relational information between all people appearing in the same image, our sequence prediction model should “see” all other instances in the image before predicting the query instance. Fig. 5 provides a demonstration of the inference process. Note that the same process is done for every instance in a test image.

For each of the sequences for a query instance, we first feed the global image feature to generate initial state of the LSTM. The labels at subsequent steps are then predicted deterministically. More specifically, at step t ($t > 0$), the predicted label is the identity with the maximum output probability, *i.e.*, $\mathbf{y}_t = \text{argmax} p(\mathbf{y})$, which is then used as the input label for joint embedding with another randomly selected instance at the next step. This process stops after the query instance has been processed, which results in a probability distribution over identities for the query instance at the end of the sequence.

Given all the probability distributions from different sequences of the query instance, we take the element-wise maximum of the probability distributions and then the identity label with the maximum probability after this operation is assigned to the query instance.

The above inference process is applied to every instance

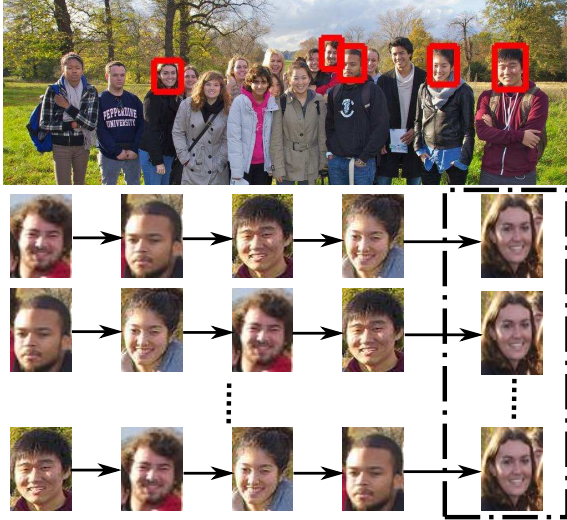


Figure 5. At test time, to predict the label of a query instance (e.g., the one in the dashed box), we generate multiple sequences for this instance, all with the instance at the end and other instances randomly ordered. The probability of an identity label for this instance is the maximum probability of this label in all sequences.

in a test image, i.e., every test instance will be used as a query instance for inference.

3.3. Discussion

Context information. Contextual cues beyond human appearance have been found valuable for person recognition in photo albums [14, 12]. There are two types of contextual cues exploited in our framework.

The first type of contextual cue is the relation context, which is the relational information between people in the same image, i.e., some people tend to appear together in the same image. In our work, this relational information is naturally handled by our sequence prediction formulation using an LSTM. Note that the relational information captured in our work is not the simple label co-occurrence between people, as our LSTM receives a joint embedding of both label and appearance as input (Fig. 4).

The second type of contextual cue is the scene context. This is based on the assumption that some identities appear more frequently in some certain scenes. Thus the scene context can be utilized as a prior for indicating which identities are likely to appear. In our work, this is done by feeding the global image feature to the LSTM at the initial step.

As we will show later in the experiments (Sec. 4.2), the above contextual cues are vital for improving the classification accuracy.

Multiple body regions. So far we only assume that \mathcal{B} is from a single annotated body part (e.g., head). Recent

Test split	#identities	#instances		#multi instances	
		$test_0$	$test_1$	$test_0$	$test_1$
Original	581	6442	6443	2802	2797
Album	581	6497	6388	2814	2751
Time	581	6440	6445	2591	2647
Day	199	2484	2485	792	744

Table 1. Statistics of the four test splits of the PIPA dataset. “#identities” denotes the number of identity labels, “#instances” refers to the number of all instances, and “#multi instances” is the number of instances from images contain multiple instances.

studies on person recognition in photo albums [29, 14, 12] have shown that recognition performance can be improved by fusing information from multiple body regions. In this work, we also extend our formulation to handle multiple body regions. One obvious solution is feature concatenation across multiple regions. We provide in-depth analysis of some feature fusion methods in the experiments section (Sec. 4.3).

4. Experiments

4.1. Experimental setup

Dataset and evaluation metric. The *People In Photo Albums (PIPA)* dataset [29] is adopted for evaluation of our approach as well as some baselines. The PIPA dataset is partitioned into train, validation, test, and leftover sets. The head region of each instance has been annotated in all sets (see Fig. 2). In previous work [29, 14, 12], the training set is used only for learning good feature representations for body regions. In the standard evaluation setting proposed in [29], the test set itself is split into two subsets, $test_0$ and $test_1$ with roughly the same number of instances. Given a recognition system that is trained on $test_0$, it is then evaluated on $test_1$, and vice versa.

Recently, in addition to the original test split proposed in [29], three more challenging splits are introduced in [14], including album, time and day splits. The album split ensures instances in $test_0$ and $test_1$ are from different albums, while time and day splits emphasize the temporal distance between $test_0$ and $test_1$ (different events, different days etc.). Generally speaking, the ranking of these four splits in the order of increasing difficulty is: original, album, time and day split. We provide an overview of the statistics of the four test splits in Table 1.

Following previous work [29, 14, 12], classification accuracy is used to evaluate the performance of our approach and some baselines. Specifically, the average classification accuracy over the $test_0$ and $test_1$ is reported. To have an in-depth understanding of our system, we also report the average accuracy on instances from images with multiple

Split	Method	Relation	Scene	Acc overall (%)	Acc multi (%)	Acc single (%)
Original	Appearance-only	–	–	75.43	77.93	73.51
	Ours-relation	✓	–	76.73	80.73	73.66
	Ours	✓	✓	81.75	84.85	79.36
Album	Appearance-only	–	–	68.31	72.00	65.52
	Ours-relation	✓	–	68.85	73.41	65.38
	Ours	✓	✓	74.21	78.22	71.16
Time	Appearance-only	–	–	57.19	59.79	55.39
	Ours-relation	✓	–	58.57	63.19	55.39
	Ours	✓	✓	63.73	67.17	61.37
Day	Appearance-only	–	–	36.37	36.72	36.23
	Ours-relation	✓	–	40.39	44.25	38.71
	Ours	✓	✓	42.75	47.25	40.74

Table 2. Classification accuracy(%) of two baselines as well as our full system under four settings on the PIPA test set. “Acc overall” refers to classification accuracy computed from all instances, whereas “Acc multi” (*resp.* “Acc single”) refers to accuracy of instances from images contain multiple (*resp.* single) instances. Head region is adopted for this analysis. Clearly, by modeling both relation and scene context, our full system has outperformed two baselines by a noticeable margin.

and single instances respectively.

Implementation details. Two body regions, head and upper body, are exploited in our system. Based on bounding box annotations of the head region which have been provided by the PIPA dataset, we estimate the annotations for the upper body region, similar to [14]. To learn feature representations, we fine-tune two VGG-16 networks [16] on the PIPA training set for these two regions respectively. On the test splits ($test_0$ and $test_1$), we extract CNN features from the $fc7$ layer of the fine-tuned networks for the two regions. The global image content I is the 4096-D $fc7$ feature extracted from the vanilla VGG-16 network, which is pre-trained on the ImageNet [5].

As each image may contain different numbers of instances, we unroll the LSTM to a fixed 22 steps (the maximum number of instances can appear on the PIPA test splits) during training. For images with instances less than 22, we pad labels with zeros and do not calculate loss on the padded labels.

We train all our weights, including the label embedding weight U_y , the image embedding weight U_b , the classification weights W , and weights in the LSTM, and using stochastic gradient descent with the Adam optimizer [11]. The initial learning rate is setted as 0.001 and decreased by 10 times after 20 epochs. We stop training after 80 epochs. We use 512 dimensions for the embeddings (Eq. 6 and Eq. 7) and the size of the LSTM memory.

4.2. Ablation study

To investigate the impact of different components of our approach, we consider the following two baselines as well as our full system.

1. “Appearance-only”: given the CNN network fine-tuned on a body region on the PIPA training set, we fine-tune the last fully-connected layer on the test instances on the either of the two test splits and evaluate on the other. In this sense, this setting is only based on the visual appearance of identities without any contextual information.
2. “Ours-relation”: We still use our LSTM-based framework to model the relation context between people, but we do not feed the LSTM with in global image content I at the initial step. In other words, the scene context is not exploited in this setting. Therefore, the only contextual cue exploited in this setting is the relation context.
3. “Ours”: Our full sequence prediction model with the global image content I fed at the initial step. In this sense, both relation and scene context are exploited.

Table. 2 depicts the performance of these three approaches under four different settings.

The importance of relation context. Comparing the overall performance of “Our-relation” with that of the “Appearance-only” (Table. 2), we observe that the former bypasses the later in all of the four settings. This reflects that our sequence prediction model has successfully taken advantage of the relation information between multiple people in the same photo.

When taking a closer look at the result of the two baselines in terms of “Acc multi”, it is clear that “Acc multi” has witnessed substantial increases in the “Our-relation” case, which contributes to the improvements in overall accuracy. We also observe that “Acc single” has stayed stable between the “Appearance-only” and “Our-relation”

Test split	Addition	Element-wise max
Original	81.51	81.75
Album	73.21	74.21
Time	62.97	63.73
Day	43.15	42.75

Table 3. Classification accuracy (%) of the two variants of joint embedding formulation in Fig. 4. Head region is used in this case.

cases. This is understandable because when there is only one instance in the image, there is no relation context to be exploited.

The importance of scene context. Comparing the performance of our full system with that of the “Our-relation” case (Table. 2), we observe that feeding the global image content to the our sequence prediction model consistently leads to further substantial improvements, resulting in the best performance in all settings.

More specifically, in all of the four settings, in terms of the overall accuracy, there is about a 3 ~ 4% accuracy gain by using our full system, compared with that of the “Our-relation” baseline which does not take advantage of the scene context. Similar improvements are also observed in both “Acc multi” and “Acc single”. This verifies that the scene context is valuable for person recognition.

Also the recent work of [14] has found that the scene contains useful information for person recognition, thus the observation in our work is consistent with [14]. In contrast with [14] in which the global image cue is analyzed independent of other cues, we are the first to incorporate different contextual cues, along with instances’ visual appearances into an end-to-end trainable framework.

Joint embedding analysis. As depicted in Fig. 4, we have proposed two formulations for the joint embedding of instance feature and label, in which information from the two sources is fused either by addition (Eq. 6) or element-wise max (Eq. 7). We hereby analyze the performance of these two variants on the PIPA test set.

As depicted in Table. 3, both embedding formulations have shown very close performance in the four settings, with the “Max” fusion slightly bypassing the “Addition” fusion in the three out of four settings. Therefore, we opt for the “Max” fusion method to report results in the following.

4.3. Region analysis

Recent studies on person recognition in photo albums [29, 14, 12] have shown different human body regions are complementary for achieving good performance. Thus, an ensemble of models built on different body regions are adopted in previous works. For instance, 107 poselet [3] classifiers are used by Zhang *et al.* [29].

Test split	Single region		Multiple region fusion		
	Head	Upper	Avg	Max	Concat
Original	81.75	79.92	84.93	84.07	82.86
Album	74.21	70.78	78.25	75.88	74.66
Time	63.73	58.80	66.43	65.63	63.62
Day	42.75	34.61	43.73	43.55	41.56

Table 4. Classification accuracy (%) using a single body region (columns 2-3), and their different fusions (columns 3-5).

Oh *et al.* [14] have analyzed the contribution of different body regions (*e.g.*, face, head, upper body, full body) to the recognition performance. Following the above works, we also extend our model to handle multiple body regions. In particular, the two body regions exploited in our work are the head and upper body regions.

Single region. As shown in Table. 4, the usage of the head region outperforms the upper body region in all settings, and the gap increases when the dataset becomes more challenging (*e.g.*, time or day setting). This is understandable because in the time or day setting, the appearance of the upper body of the same person can have significant changes (*e.g.*, change of clothing), which results in more failures when using features from the upper body region. In comparison, features from the head region are relatively stable.

Multiple region fusion. As different body regions have different levels of relative importance for the final recognition performance, the information from different regions should be fused in order to achieve better performance. In our work, we study three methods to fuse information from head and upper body regions in the our sequence prediction model.

1. “Concat”. This is our model trained with concatenated features from head and upper body regions.
2. “Avg”. Two models are trained with features from the head and upper body region respectively. The probability of a test instance is the average of probabilities from the two models.
3. “Max”. Same as the “Avg” case except that we take the maximum of probabilities rather than average.

The performance of the above three methods are presented in Table. 4. Clearly, the “Avg” fusion method shows the largest improvement over the performance of using a single head or upper body region. Therefore, we use the performance of “Avg” fusion method to compare with state-of-the-art approaches in the following.

Test split	Head		Upper body		Head + Upper body	
	[14]	Our	[14]	Our	[14]	Our
Original	76.42	81.75	75.07	79.92	83.63	84.93
Album	67.48	74.21	64.65	70.78	74.52	78.25
Time	57.05	63.73	50.90	58.80	63.98	66.43
Day	36.37	42.75	24.17	34.61	38.94	43.73

Table 5. Classification accuracy(%) of Oh *et al.* [14] and ours under the four different settings on the PIPA test set, using head region, upper body region and their fusion.



Figure 6. Our predictions using the head region on the PIPA test set. Correct predictions are denoted by red bounding boxes whereas wrong predictions are in green. Best viewed in color.

4.4. Comparison to state-of-the-arts

We now compare the performance of our approach to the state-of-the-art approaches in person recognition.

More specifically, we compare with [14] on head region, upper body region and their fusion (see Table. 5). In [14], the performance is achieved based on features extracted from fine-tuned CNNs without incorporating any contextual information. As shown in Table. 5, our sequence prediction model outperforms [14] by a reasonable margin in all of four settings. In the original setting, our result (84.93%) also outperforms that of Zhang *et al.* [29] (83.05%), although we just use two body regions. In comparison, features from 107 poselet regions are used in [29].

We are aware that by using better features learnt from external data, it is possible to achieve higher recognition performance. For instance, higher accuracy on the day setting is reported in [12], in which the features are extracted from the face region using a CNN trained for the face recognition task. However, in our work, we just use features fine-tuned on the PIPA training set. We are interested to show how our sequence prediction formulation models the contextual in-

formation (including the relation context and scene context) and visual cues in a unified framework for person recognition, which improves the recognition performance.

4.5. Visualization

We provide some visual examples of predicted instances in the PIPA test set by our sequential model in Fig. 6. As shown in Fig. 6, in most photos, our model correctly recognize the identities in the photo, including some challenging cases, such as non-frontal faces.

5. Conclusion

In this work, we have introduced a sequence prediction formulation for the task of person recognition in photo albums. The advantage of our approach is that it can model both the rich contextual information in the photo, and individual’s appearance in a unified framework. We trained our model end-to-end and witnessed a significant boost in the recognition performance, compared with baselines and state-of-the-approaches which do not exploit the contextual information for the task.

References

- [1] D. Anguelov, K. Lee, S. B. Gökürk, and B. Sumengen. Contextual identity recognition in personal photo albums. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2007.
- [2] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Proc. Adv. Neural Inf. Process. Syst.*, 2015.
- [3] L. D. Bourdev and J. Malik. Poselets: Body part detectors trained using 3d human pose annotations. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2009.
- [4] Q. Dai, P. Carr, L. Sigal, and D. Hoiem. Family member identification from photo collections. 2015.
- [5] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 248–255, 2009.
- [6] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, T. Darrell, and K. Saenko. Long-term recurrent convolutional networks for visual recognition and description. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015.
- [7] A. C. Gallagher and T. Chen. Clothing cosegmentation for recognizing people. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2008.
- [8] A. C. Gallagher and T. Chen. Understanding images of groups of people. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2009.
- [9] G. Gkioxari, A. Toshev, and N. Jaitly. Chained predictions using convolutional neural networks. In *Proc. Eur. Conf. Comp. Vis.*, 2016.
- [10] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8), 1997.
- [11] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proc. Int. Conf. Learn. Repr.*, 2015.
- [12] H. Li, J. Brandt, Z. Lin, X. Shen, and G. Hua. A multi-level contextual model for person recognition in photo albums. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016.
- [13] C. S. Mathialagan, A. C. Gallagher, and D. Batra. VIP: finding important people in images. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015.
- [14] S. J. Oh, R. Benenson, M. Fritz, and B. Schiele. Person recognition in personal photo collections. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2015.
- [15] S. J. Oh, R. Benenson, M. Fritz, and B. Schiele. Faceless person recognition: Privacy implications in social media. In *Proc. Eur. Conf. Comp. Vis.*, 2016.
- [16] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proc. Int. Conf. Learn. Repr.*, 2015.
- [17] R. Stewart, M. Andriluka, and A. Y. Ng. End-to-end people detection in crowded scenes. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016.
- [18] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Proc. Adv. Neural Inf. Process. Syst.*, 2014.
- [19] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2014.
- [20] S. Venugopalan, M. Rohrbach, J. Donahue, R. J. Mooney, T. Darrell, and K. Saenko. Sequence to sequence - video to text. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2015.
- [21] O. Vinyals, S. Bengio, and M. Kudlur. Order matters: Sequence to sequence for sets. In *Proc. Int. Conf. Learn. Repr.*, 2016.
- [22] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015.
- [23] G. Wang, A. C. Gallagher, J. Luo, and D. A. Forsyth. Seeing people in social context: Recognizing people and social relationships. In *Proc. Eur. Conf. Comp. Vis.*, 2010.
- [24] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu. CNN-RNN: A unified framework for multi-label image classification. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016.
- [25] P. Wang, Q. Wu, C. Shen, A. v. d. Hengel, and A. Dick. Fvqa: Fact-based visual question answering. *arXiv:1606.05433*, 2016.
- [26] Q. Wu, C. Shen, A. v. d. Hengel, P. Wang, and A. Dick. Image captioning and visual question answering based on attributes and their related external knowledge. *arXiv preprint arXiv:1603.02814*, 2016.
- [27] Q. Wu, C. Shen, L. Liu, A. R. Dick, and A. van den Hengel. What value do explicit high level concepts have in vision to language problems? In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016.
- [28] Q. Wu, P. Wang, C. Shen, A. van den Hengel, and A. R. Dick. Ask me anything: Free-form visual question answering based on knowledge from external sources. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016.
- [29] N. Zhang, M. Paluri, Y. Taigman, R. Fergus, and L. D. Bourdev. Beyond frontal faces: Improving person recognition using multiple cues. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015.