# Image Privacy Prediction Using Deep Neural Networks

Ashwini Tonge
Kansas State University
KS, USA
atonge@ksu.edu

Cornelia Caragea
University of Illinois at Chicago
IL, USA
cornelia@uic.edu

## ABSTRACT

Images today are increasingly shared online on social networking sites such as Facebook, Flickr, Foursquare, and Instagram. Image sharing occurs not only within a group of friends but also more and more outside a user's social circles for purposes of social discovery. Despite that current social networking sites allow users to change their privacy preferences, this is often a cumbersome task for the vast majority of users on the Web, who face difficulties in assigning and managing privacy settings. When these privacy settings are used inappropriately, online image sharing can potentially lead to unwanted disclosures and privacy violations. Thus, automatically predicting images' privacy to warn users about private or sensitive content before uploading these images on social networking sites has become a necessity in our current interconnected world.

In this paper, we explore learning models to automatically predict appropriate images' privacy as *private* or *public* using carefully identified image-specific features. We study deep visual semantic features that are derived from various layers of Convolutional Neural Networks (CNNs) as well as textual features such as user tags and deep tags generated from deep CNNs. Particularly, we extract deep (visual and tag) features from four pre-trained CNN architectures for object recognition, i.e., AlexNet, GoogLeNet, VGG-16, and ResNet, and compare their performance for image privacy prediction. Among all four networks, we observe that ResNet produces the best feature representations for this task. We also fine-tune the pre-trained CNN architectures on our privacy dataset and compare their performance with the models trained on pre-trained features. The results show that even though the overall performance obtained using the fine-tuned networks is comparable to that of pre-trained networks, the fine-tuned networks provide an improved performance for the private class as compared to models trained on the pre-trained features. Results of our experiments on a Flickr dataset of over thirty thousand images show that the learning models trained on features extracted from ResNet outperform the state-of-the-art models for image privacy prediction. We further investigate the combination of user tags and deep tags derived from CNN architectures using two settings: (1) SVM on the bag-of-tags features; and (2) text-based CNN. We compare these models with the models trained on ResNet visual features obtained for privacy prediction. Our results show that even though the models trained on the visual features perform better than those trained on the tag features, the combination of deep visual features with image tags shows improvements in performance over the individual feature sets. Our code, features, and the dataset used in experiments are available at https://github.com/ashwinitonge/deepprivate.git.

## CCS CONCEPTS

• **Security and privacy** → **Software and application security**;
• **Social network security and privacy**;

## KEYWORDS

Social networks, image analysis, image privacy prediction, deep learning.

## 1 INTRODUCTION

Online image sharing through social networking sites such as Facebook, Flickr, and Instagram is on the rise, and so is the sharing of private or sensitive images, which can lead to potential threats to users' privacy when inappropriate privacy settings are used in these platforms. Many users quickly share private images of themselves and their family and friends, without carefully thinking about the consequences of unwanted disclosure and privacy violations [Ahern et al. 2007; Zerr et al. 2012b]. For example, it is common now to take photos at cocktail parties and share them on social networking sites without much hesitation. The smartphones facilitate the sharing of photos virtually at any time with people all around the world. These photos can potentially reveal a user's personal and social habits and may be used in the detriment of the photos' owner.

Gross and Acquisti [2005] analyzed more than 4,000 Carnegie Mellon University students' Facebook profiles and outlined potential threats to privacy. The authors found that users often provide personal information generously on social networking sites, but they rarely change default privacy settings, which could jeopardize their privacy. In a parallel study, Lipford et al. [2008] showed that, although current social networking sites allow users to change their privacy preferences, the vast majority of users on the Web face difficulties in assigning and managing privacy settings. Interestingly, Orekondy et al. [2017] showed that, even when users change their privacy settings to comply with their personal privacy preference, they often misjudge the private information in images, which fails to enforce their own privacy preferences. Not surprising, employers these days often perform background checks for their future employees using social networks and about 8% of companies have
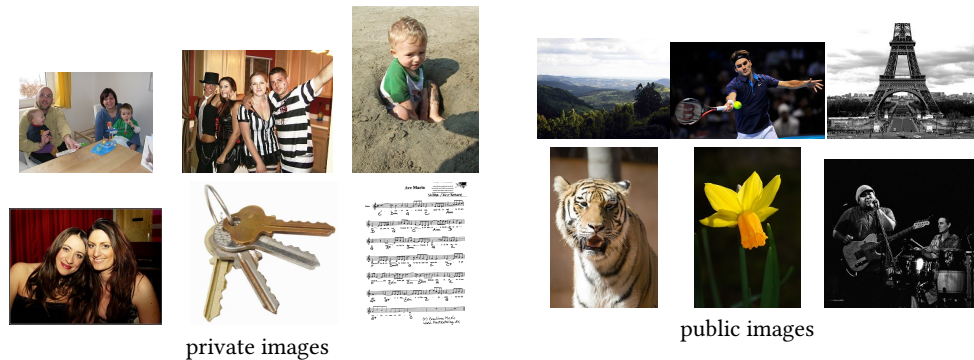
Figure 1: Examples of images manually identified as private (left) and public (right).

already fired employees due to their inappropriate social media content [Waters and Ackerman 2011]. A study carried out by the Pew Research center reported that 11% of users of social networks regret the posted content [Madden 2012]. The Director of the AI Research at Facebook, Yann LeCun [2017] urges the development of a digital assistant to warn people about private or sensitive content before embarrassing photos are shared with everyone on social networks.

Identifying private or sensitive content from images is inherently difficult because images' privacy is dependent on the owners' personality traits and their level of awareness towards privacy. Still, images' privacy is not purely subjective, but generic patterns of privacy exist. Consider, for example, the images shown in Figure 1, which are manually annotated and consistently rated as *private* and *public* by multiple annotators in a study conducted by Zerr et al. [2012b,a]. Notice that the presence of people generally pinpoints to private images, although this is not always true. For example, an image of a musical band in concert is considered to be public. Similarly, images with no people in them could be private, e.g., images with door keys, music notes, legal documents, or someone's art are considered to be private. Indeed, Laxton et al. [2008] described a "tele-duplication attack" that allows an adversary to create a physical key duplicate simply from an image.

Researchers showed that generic patterns of images' privacy can be automatically identified when a large set of images are considered for analysis and investigated binary prediction models based on user tags and image content features such as SIFT (Scale Invariant Feature Transform) and RGB (Red Green Blue) [Squicciarini et al. 2014, 2017a; Zerr et al. 2012b]. More recently, several studies [Tonge and Caragea 2015, 2016, 2018; Tran et al. 2016] started to explore privacy frameworks that leverage the benefits of Convolutional Neural Networks (CNNs) for object recognition since, intuitively, the objects present in images significantly impact images' privacy (as can be seen from Figure 1). However, these studies used only the AlexNet architecture of CNNs on small dataset sizes. To date, many deep CNN architectures have been developed and achieve state-of-the-art performance on object recognition. These CNNs include GoogLeNet [Szegedy et al. 2014], VGG-16 [Simonyan and Zisserman 2014], and ResNet [He et al. 2016a] (in addition to AlexNet [Krizhevsky et al. 2012]). Towards this end, in this paper,

we present an extensive study to carefully identify the CNN architectures and features derived from these CNNs that can adequately predict the class of an image as *private* or *public*. Our research is motivated by the fact that increasingly, online users' privacy is routinely compromised by using social and content sharing applications [Zheleva and Getoor 2009]. Our models can help users to better manage their participation in online image sharing sites by identifying the sensitive content from the images so that it becomes easier for regular users to control the amount of personal information that they share through these images.

Our contributions are as follows:

- We study deep visual semantic features and deep image tags derived from CNN architectures pre-trained on the ImageNet dataset and use them in conjunction with Support Vector Machine (SVM) classifiers for image privacy prediction. Specifically, we extract deep features from four successful (pre-trained) CNN architectures for object recognition, AlexNet, GoogLeNet, VGG-16, and ResNet and compare their performance on the task of privacy prediction. Through carefully designed experiments, we find that ResNet produces the best feature representations for privacy prediction compared with the other CNNs.

- We fine-tune the pre-trained CNN architectures on our privacy dataset and use the softmax function to predict the images' privacy as *public* or *private*. We compare the fine-tuned CNNs with the SVM models obtained on the features derived from the pre-trained CNNs and show that, although the overall performance obtained by the fine-tuned CNNs is comparable to that of SVM models, the fine-tuned networks provide improved recall for the private class as compared to the SVM models trained on the pre-trained features.

- We show that the best feature representation produced by ResNet outperforms several baselines for image privacy prediction that consider CNN-based models and SVM models trained on traditional visual features such as SIFT and global GIST descriptor.

- Next, we investigate the combination of user tags and deep tags derived from CNNs in two settings: (1) using SVM on the bag-of-tags features; and (2) applying the text CNN [Kim 2014] on the combination of user tags and deep tags for privacy prediction using the softmax function. We compare

these models with the models trained on the most promising visual features extracted from ResNet (obtained from our study) for privacy prediction. Our results show that the models trained on the visual features perform better than those trained on the tag features.

- Finally, we explore the combination of deep visual features with image tags and show further improvement in performance over the individual sets of features.

The rest of the paper is organized as follows. We summarize prior work in Section 2. In Section 3, we describe the problem statement in details. Section 4 describes the image features obtained from various CNNs for privacy prediction, whereas in Section 5, we provide details about the dataset that we use to evaluate our models. In Section 6, we present the experiments and describe the experimental setting and results. We finish our analysis in Section 7, where we provide a brief discussion of our main findings, interesting applications of our work, future directions, and conclude the paper.

## 2 RELATED WORK

Emerging privacy violations in social networks have started to attract various researchers to this field [Zheleva and Getoor 2009]. Researchers also provided public awareness of privacy risks associated with images shared online [Henne et al. 2013; Xu et al. 2015]. Along this line, several works are carried out to study users' privacy concerns in social networks, privacy decisions about sharing resources, and the risk associated with them [Ghazinour et al. 2013; Gross and Acquisti 2005; Ilia et al. 2015; Krishnamurthy and Wills 2008; Parra-Arnau et al. 2014; Parra-Arnau et al. 2012; Simpson 2008].

Moreover, several works on privacy analysis examined privacy decisions and considerations in mobile and online photo sharing [Ahern et al. 2007; Besmer and Lipford 2009; Gross and Acquisti 2005; Jones and O'Neill 2011]. For example, Ahern et al. [2007] explored critical aspects of privacy such as users' consideration for privacy decisions, content and context based patterns of privacy decisions, and how different users adjust their privacy decisions and behavior towards personal information disclosure. The authors concluded that applications that could support and influence users' privacy decision-making process should be developed. Jones and O'Neill [2011] reinforced the role of privacy-relevant image concepts. For instance, the authors determined that people are more reluctant to share photos capturing social relationships than photos taken for functional purposes; certain settings such as work, bars, concerts cause users to share less. Besmer and Lipford [2009] mentioned that users want to regain control over their shared content, but meanwhile, they feel that configuring proper privacy settings for each image is a burden.

More recent and related to our line of work are the automated image privacy approaches that have been explored along four lines of research: *social group based approaches*, in which users' profiles are used to partition the friends' lists into multiple groups or circles, and the friends from the same circle are assumed to share similar privacy preferences; *location-based approaches*, in which location contexts are used to control the location-based privacy disclosures; *tag-based approaches*, in which tags are used for privacy setting

recommendations; and *visual-based approaches*, in which the visual content of images is leveraged for privacy prediction.

*Social group based approaches.* Several works emerged to provide the automated privacy decisions for images shared online based on the social groups or circles [Bonneau et al. 2009a,b; Christin et al. 2013; Danezis 2009; Fang and LeFevre 2010; Joshi and Zhang 2009; Kepez and Yolum 2016; Klemperer et al. 2012; Mannan and van Oorschot 2008; Pesce et al. 2012; Petkos et al. 2015; Squicciarini et al. 2012, 2015, 2009; Watson et al. 2015; Yuan et al. 2017; Zerr et al. 2012b]. For example, Christin et al. [2013] proposed an approach to share content with the users within privacy bubbles. Privacy bubbles represent the private sphere of the users and the access to the content is provided by the bubble creator to people within the bubble. Bonneau et al. [2009b] introduced the notion of privacy suites which recommend users a set of privacy settings that "expert" users or the trusted friends have already established so that ordinary users can either directly accept a setting or perform minor modifications only. Fang and LeFevre [2010] developed a privacy assistant to help users grant privileges to their friends. The approach takes as input the privacy preferences for the selected friends and then, using these labels, constructs a classifier to assign privacy labels to the rest of the (unlabeled) friends based on their profiles. Danezis [2009] generated privacy settings based on the policy that the information produced within the social circle should remain in that circle itself. Along these lines, Adu-Oppong et al. [2008] obtained privacy settings by forming clusters of friends by partitioning a user's friends' list. Yuan et al. [2017] proposed an approach for context-dependent and privacy-aware photo sharing. This approach uses the semantics of the photo and the requester's contextual information in order to define whether an access to the photo will be granted or not at a certain context. These social group based approaches mostly considered the user trustworthiness, but ignored the image content sensitiveness, and thus, they may not necessarily provide appropriate privacy settings for online images as the privacy preferences might change according to sensitiveness of the image content.

*Location-based approaches.* These approaches [Baokar 2016; Bilogrevic et al. 2016; Choi et al. 2017; Fisher et al. 2012; Freudiger et al. 2012; Friedland and Sommer 2010; Olejnik et al. 2017; Ravichandran et al. 2009; Shokri et al. 2011; Toch 2014; Yuan et al. 2017; Zhao et al. 2014] leverage geo-tags, visual landmarks and other location contexts to control the location-based privacy disclosures. The geo-tags can be provided manually via social tagging or by adding location information automatically through the digital cameras or smart-phones having GPS. The location can also be inferred by identifying places from the shared images through the computer vision techniques.

*Tag-based approaches.* Previous work in the context of tag-based access control policies and privacy prediction for images [Apostolova and Demner-Fushman 2009; De Choudhury et al. 2009; Klemperer et al. 2012; Kurtan and Yolum 2018; Mannan and van Oorschot 2008; Pesce et al. 2012; Ra et al. 2013; Squicciarini et al. 2012, 2015, 2017b; Vyas et al. 2009; Yeung et al. 2009; Zerr et al. 2012b] showed initial success in tying user tags with access control rules. For example, Squicciarini et al. [2012, 2017b], Zerr et al. [2012b], and Vyas

et al. [2009] explored learning models for image privacy prediction using user tags and found that user tags are informative for predicting images' privacy. Moreover, Squicciarini et al. [2015] proposed an Adaptive Privacy Policy Prediction framework to help users control access for their shared images. The authors investigated social context, image content, and metadata as potential indicators of privacy preferences. Klemperer et al. [2012] studied whether the user annotated tags help to create and maintain access-control policies more intuitively. However, the scarcity of tags for many online images [Sundaram et al. 2012] and the dimensions of user tags precluded an accurate analysis of images' privacy. Hence, in our previous work, [Tonge and Caragea 2015, 2016, 2018; Tonge et al. 2018a,b], we explored automatic image tagging and showed that the predicted tags combined with user tags can improve the overall privacy prediction performance.

*Visual-based approaches.* Several works used visual features derived from the images' content and showed that they are informative for predicting images' privacy settings [Buschek et al. 2015; Dufaux and Ebrahimi 2008; Hu et al. 2016; Kuang et al. 2017; Nakashima et al. 2011, 2012, 2016; Orekondy et al. 2018; Shamma and Uddin 2014; Squicciarini et al. 2014, 2017a; Tonge and Caragea 2015, 2016, 2018; Tran et al. 2016; von Zezschwitz et al. 2016; Wu et al. 2018; Yu et al. 2017a, 2018; Yuan et al. 2018; Zerr et al. 2012b; Zhang et al. 2005]. For example, Buschek et al. [2015] presented an approach to assigning privacy to shared images using metadata (location, time, shot details) and visual features (faces, colors, edges). Zerr et al. [2012b] proposed privacy-aware image classification and learned classifiers on Flickr images. The authors considered image tags and visual features such as color histograms, faces, edge-direction coherence, and Scale Invariant Feature Transform (SIFT) for the privacy classification task. SIFT as well as GIST are among the most widely used traditional features for image analysis in computer vision. SIFT [Lowe 2004] detects scale, rotation, and translation invariant key-points of objects in images and extracts a pool of visual features, which are represented as a "bag-of-visual-words." GIST [Oliva and Torralba 2001] encodes global descriptors for images and extracts a set of perceptual dimensions (naturalness, openness, roughness, expansion, and ruggedness) that represent the dominant spatial structure of the scene. Squicciarini et al. [2014, 2017a] performed an in-depth analysis of image privacy classification using Flickr images and found that SIFT and image tags work best for predicting privacy of users' images.

Recently, the computer vision community has shifted towards convolutional neural networks (CNNs) for tasks such as object detection [Sermanet et al. 2014, 2013] and semantic segmentation [Farabet et al. 2013]. CNNs have acquired state-of-the-art results on ImageNet for object recognition [Russakovsky et al. 2015] using supervised learning [Krizhevsky et al. 2012]. Given the recent success of CNNs, several researchers [Kuang et al. 2017; Tonge and Caragea 2015, 2016, 2018; Tran et al. 2016; Yu et al. 2017a, 2018] showed promising privacy prediction results compared with visual features such as SIFT and GIST. Yu et al. [2017b] adopted CNNs to achieve semantic image segmentation and also learned object-privacy relatedness to identify privacy-sensitive objects.

Using CNNs, some works started to explore personalized privacy prediction models [Orekondy et al. 2017; Spyromitros-Xioufis

et al. 2016; Zhong et al. 2017]. For example, Spyromitros-Xioufis et al. [2016] used features extracted from CNNs to provide personalized image privacy classification. Zhong et al. [2017] proposed a Group-Based Personalized Model for image privacy classification in online social media sites that learns a set of archetypical privacy models (groups) and associates a given user with one of these groups. Orekondy et al. [2017] defined a set of privacy attributes, which were first predicted from the image content and then used these attributes in combination with users' preferences to estimate personalized privacy risk. Although there is evidence that individuals' sharing behavior is unique, Zhong et al. [2017] argued that personalized models generally require large amounts of user data to learn reliable models, and are time and space consuming to train and store models for each user, while taking into account possible sudden changes of users' sharing activities and privacy preferences. Orekondy et al. [2017] tried to resolve some of these limitations by clustering users' privacy profiles and training a single classifier that maps the target user into one of these clusters to estimate the personalized privacy score. However, the users' privacy profiles are obtained using a set of attributes. which are defined based on the Personally Identifiable Information [McCallister et al. 2010], the US Privacy Act of 1974 and official online social network rules, instead of collecting opinions about sensitive content from the actual users of social networking sites. Hence, the definition of sensitive content may not meet a user's actual needs, which limits their applicability in a real-world usage scenario [Li et al. 2018]. In this context, it is worth mentioning that CNNs were also used in another body of privacy related work such as multi-party privacy conflict detection [Zhong et al. 2018] and automatic redaction of sensitive image content [Orekondy et al. 2018].

The image representations using visual features and tags are pivotal in above privacy prediction works. In this paper, we aim to study "deep" features derived from CNNs, by abstracting out users' privacy preferences and sharing behavior. Precisely, our goal is to identify a set of "deep" features that have the highest discriminative power for image privacy prediction and to flag images that contain private or sensitive content before they are shared on social networking sites. To our knowledge, this is the first study to provide a detailed analysis of various CNN architectures for privacy prediction. Our comprehensive set of experiments can provide the community with evidence about the best CNN architecture and features for the image privacy prediction task, especially since the results obtained outperformed other complex approaches, on a large dataset of more than 30, 000 images.

## 3 PROBLEM STATEMENT

Our goal is to accurately identify private or sensitive content from images before they are shared on social networking sites. Precisely, given an image, we aim to learn models to classify the image into one of the two classes: *private* or *public*, based on generic patterns of privacy. Private images belong to the private sphere (e.g., self-portraits, family, friends, someone's home) or contain information that one would not share with everyone else (e.g., private documents). Public images capture content that can be seen by everyone without incurring privacy violations. To achieve our goal, we extract a variety of features from several CNNs and identify those

features that have the highest discriminative power for image privacy prediction.

As the privacy of an image can be determined by the presence of one or more objects described by the visual content and the description associated with it in the form of tags, we consider both visual features and image tags for our analysis. For the purpose of this study, we did not consider other contextual information about images (e.g., personal information about the image owner or the owner social network activities, which may or may not be available or easily accessible) since our goal is to predict the privacy of an image solely from the image's content itself. We rely on the assumption that, although privacy is a subjective matter, generic patterns of images' privacy exist that can be extracted from the images' visual content and textual tags.

We describe the feature representations considered for our analysis in the next section.

## 4 IMAGE ENCODINGS

In this section, we provide details on visual content encodings and tag content encodings derived from various CNNs (pre-trained and fine-tuned) to carefully identify the most informative feature representations for image privacy prediction. Particularly, we explore four CNN architectures, AlexNet [Krizhevsky et al. 2012], GoogLeNet [Szegedy et al. 2014], VGG-16 [Simonyan and Zisserman 2014], and ResNet [He et al. 2016a] to derive features for all images in our dataset, which are labeled as private or public. The choice of these architectures is motivated by their good performance on the large scale ImageNet object recognition challenge [Russakovsky et al. 2015]. We also leverage a text-based CNN architecture used for sentence classification [Kim 2014] and apply it to images' textual tags for privacy prediction.

### 4.1 Preliminary: Convolutional Neural Networks

CNN is a type of feed-forward artificial neural network which is inspired by the organization of the animal visual cortex. The learning units in the network are called neurons. These neurons learn to convert input data, i.e., a picture of a dog into its corresponding label, i.e., "dog" through automated image recognition. The bottom layers of a CNN consist of interleaved convolution and pooling layers, and the top layers consist of fully-connected (fc) layers, and a probability (prob) layer obtained by applying the softmax function to the input from the previous fc layer, which represents the probability distribution over the available categories for an input image. As we ascend through an architecture, the network acquires: (1) lower layers features (color blobs, lines, corners); (2) middle layer features (textures resulted from a combination of lower layers); and (3) higher (deeper) layer features (high-level image content like objects obtained by combining middle layers). Since online images may contain multiple objects, we consider features extracted from deeper layers as they help to encode the objects precisely.

A CNN exploits the 2D topology of image data, in particular, *local connectivity* through convolution layers, performs *weight sharing* to handle very high-dimensional input data, and can deal with more *abstract or global information* through pooling layers. Each unit within a convolution layer receives a small region of its input at location $l$, denoted $\mathbf{r}_l(\mathbf{x})$ (a.k.a. *receptive field*), and applies a non-linear function to it. More precisely, given an input image $\mathbf{x}$, a unit that is responsible for region $l$ computes $\sigma(\mathbf{W} \cdot \mathbf{r}_l(\mathbf{x}) + \mathbf{b})$, where $\mathbf{W}$ and $\mathbf{b}$ represent the matrix of weights and the vector of biases, respectively, and $\sigma$ is a non-linear function such as the sigmoid activation or rectified linear activation function. $\mathbf{W}$ and $\mathbf{b}$ are learned during training and are shared by all units in a convolution layer. Each unit within a pooling layer receives a small region from the previous convolution layer and performs average or max-pooling to obtain more abstract features. During training, layers in CNNs are responsible for a forward pass and a backward pass. The forward pass takes inputs and generates the outputs. The backward pass takes gradients with respect to the output and computes the gradient with respect to the parameters and to the inputs, which are consecutively back-propagated to the previous layers [Jia et al. 2014].
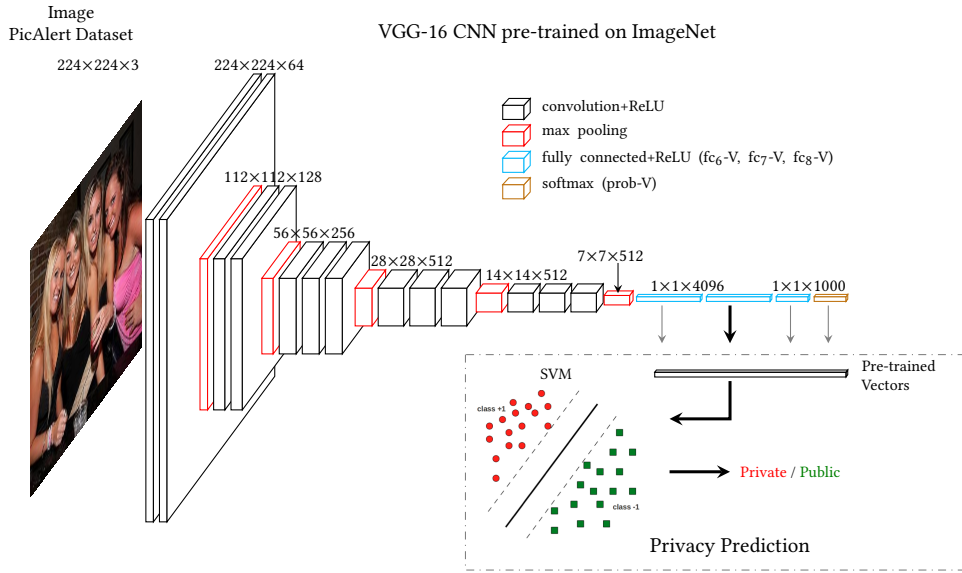
### 4.2 Features Derived Through Pre-Trained CNNs

We describe a diverse set of features derived from CNN architectures pre-trained on the ILSVRC-2012 object classification subset of the ImageNet dataset that contains 1000 object categories and 1.2 million images [Russakovsky et al. 2015]. We consider powerful features obtained from various fully-connected layers of a CNN that are generated by the previous convolutional layers, and use them to learn a decision function whose sign represents the class (*private* or *public*) assigned to an input image $\mathbf{x}$. The activations of the fully connected layers capture the complete object contained in the region of interest. Hence, we use the activations of the fully-connected layers of a CNN as a feature vector. For image encoding, we also use the probability (prob) layer obtained by applying the softmax function to the output of the (last) fully-connected layer. We extract features from the four pre-trained CNNs as follows.

The ***AlexNet*** architecture implements an eight-layer network; the first five layers of AlexNet are convolutional, and the remaining three layers are fully-connected. We extract features from the three fully-connected layers, which are referred as $fc_6$-A, $fc_7$-A, and $fc_8$-A, and from the output layer denoted as "prob-A." The dimensions of $fc_6$-A, $fc_7$-A, $fc_8$-A, and prob-A are 4096, 4096, 1000, and 1000, respectively.

The ***GoogLeNet*** architecture implements a 22 layer deep network with Inception architecture. The architecture is a combination of all layers with their output filter bank concatenated so as to form input for the next stage. We extract features from the last two layers named as "$loss_3$-G/classifier" (InnerProduct layer) and the output layer denoted as "prob-G." The dimension of $loss_3$-G and prob-G is 1000.

The ***VGG-16*** architecture implements a 16 layer deep network; a stack of convolutional layers with a very small receptive filed: $3 \times 3$ followed by fully-connected layers. The architecture contains 13 convolutional layers and 3 fully-connected layers. The number of channels of the convolutional layers starts from 64 in the first layer and then increases by a factor of 2 after each max-pooling layers until it reaches 512. We refer to features extracted from the fully-connected layers as $fc_6$-V, $fc_7$-V, $fc_8$-V, and the output layer as

**Figure 2: Image encoding using pre-trained CNN: (1) We employ a CNN (e.g. VGG-16) pre-trained on the ImageNet object dataset. (2) We derive high-level features from the image's visual content using fully connected layers (fc$_6$-V, fc$_7$-V, and fc$_8$-V) and probability layer (softmax) of the pre-trained network.**

"prob-V." The dimensions of fc$_6$-V, fc$_7$-V, fc$_8$-V, and prob-V are 4096, 4096, 1000, and 1000, respectively.

The **ResNet** (or Residual network) alleviates the vanishing gradient problem by introducing short paths to carry gradient throughout the extent of very deep networks and allows the construction of deeper architectures. A residual unit with an identity mapping is defined as:

$$X^{l+1} = X^l + \mathcal{F}(X^l)$$

where $X^l$ is the input and $X^{l+1}$ is the output of the residual unit; $\mathcal{F}$ is a residual function, e.g., a stack of two $3 \times 3$ convolution layers in [He et al. 2016a]. The main idea of the residual learning is to learn the additive residual function $\mathcal{F}$ with respect to $X^l$ [He et al. 2016b]. Intuitively, ResNets can be explained by considering residual functions as paths through which information can propagate easily. This interprets as ResNets learn more complex feature representations which are combined with the shallower descriptions obtained from previous layers. We refer to features extracted from the fully-connected layer as fc-R and the output layer as "prob-R." The dimension of fc-R and prob-R is 1000.

The feature extraction using the pre-trained network for an input image from our dataset is shown in Figure 2. In the figure, we show VGG-16 as the pre-trained network for illustrating the feature extraction.

### 4.3 Fine-tuned CNN

For this type of encoding, models trained on a large dataset (e.g., the ImageNet dataset) are fine-tuned using a smaller dataset (e.g., the privacy-labeled dataset). Fine-tuning a network is a procedure based on the concept of transfer learning [Bengio 2012; Donahue et al. 2013]. This strategy fine-tunes the weights of the pre-trained network by continuing the back-propagation on the small dataset,
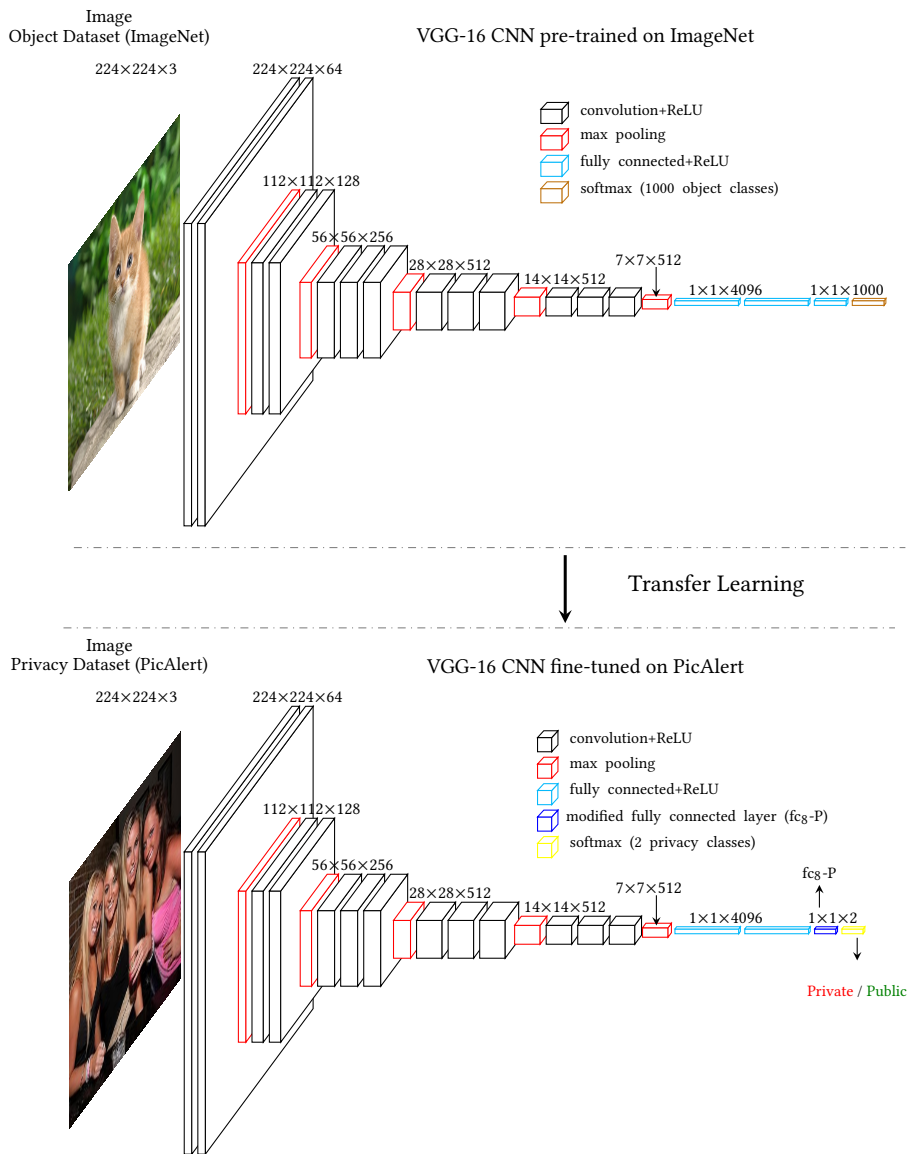
i.e., privacy dataset in our scenario. The features become more dataset-specific after fine-tuning, and hence, are distinct from the features obtained from the pre-trained CNN. We modify the last fully-connected layer of all four network architectures, AlexNet, GoogLeNet, VGG-16, and ResNet by changing the output units from 1000 (object categories) to 2 (with respect to privacy classes) (e.g., changing fc$_8$ with 1000 output units to fc$_8$-P with 2 output units). We initialize the weights of all the layers of this modified architectures with the weights of the respective layers obtained from the pre-trained networks. We train the network by iterating through all the layers of the networks using the privacy data. We use the softmax function to predict the privacy of an image. Precisely, we use the probability distribution over 2 privacy classes for the input image obtained by applying the softmax function over the modified last fully-connected layer (e.g. fc$_8$-P in VGG-16) of the fine-tuned networks (See Figure 3, second network, blue rectangle). The conditional probability distribution over 2 privacy classes can be defined using a softmax function as given below:

$$P(y = P_r | \mathbf{z}) = \frac{exp(z_{P_r})}{exp(z_{P_u}) + exp(z_{P_r})}, P(y = P_u | \mathbf{z}) = \frac{exp(z_{P_u})}{exp(z_{P_u}) + exp(z_{P_r})}$$

where, in our case, $\mathbf{z}$ is the output of the modified last fully-connected layer (e.g., the fc$_8$-P layer of VGG-16) and $P_r$ and $P_u$ denote *private* and *public* class, respectively. The fine-tuning process using VGG-16 is shown in Figure 3.

### 4.4 Image Tags (Bag-of-Tags model)

Prior works on privacy prediction [Squicciarini et al. 2014, 2017b; Tonge and Caragea 2015, 2016; Zerr et al. 2012b] found that the tags associated with images are indicative of their sensitive content. Tags are also crucial for image-related applications such as indexing, sharing, searching, content detection and social discovery [Bischoff

**Figure 3: Image encoding using fine-tuned CNN: (1) We modify the last fully-connected layer of the pre-trained network (top network) by changing the output units from** $1000$ **(object categories) to** $2$ **(privacy classes). (2) To train the modified network (bottom network) on privacy dataset, we first adopt weights of all the layers of the pre-trained network as initial weights and then iterate through all the layers using privacy data. (3) To make a prediction for an input image (privacy dataset), we use the probability distribution over** $2$ **privacy classes (softmax layer, yellow rectangle) for the input image obtained by applying the softmax function over the last modified fully-connected layer (fc$_8$-P, bottom network) of the fine-tuned network.**

et al. 2008; Gao et al. 2011; Hollenstein and Purves 2010; Tang et al. 2009]. Since not all images on social networking sites have user tags or the set of user tags is very sparse [Sundaram et al. 2012], we use an automatic technique to annotate images with tags based on their visual content as described in our previous work [Tonge and Caragea 2015, 2016]. Precisely, we predict top $k$ object categories from the probability distribution extracted from a pre-trained CNN. These top $k$ categories are images' deep tags, used to describe an image. For example, we obtain deep tags such as "Maillot," "Wig,"

"Brassiere," "Bra," "Miniskirt" for the picture in Figure 4 (note that only top 5 deep tags are shown in the figure). Note that the deep tags give some description about the image, but still some relevant tags such as "people" and "women" are not included since the 1000 object categories of the ImageNet dataset do not contain these tags. Images on social networking sites also give additional information about them through the tags assigned by the user. We call these tags "User Tags." Examples of user tags for the image in Figure 4 are: "Birthday Party," "Night Life," "People," etc. For user tags, we
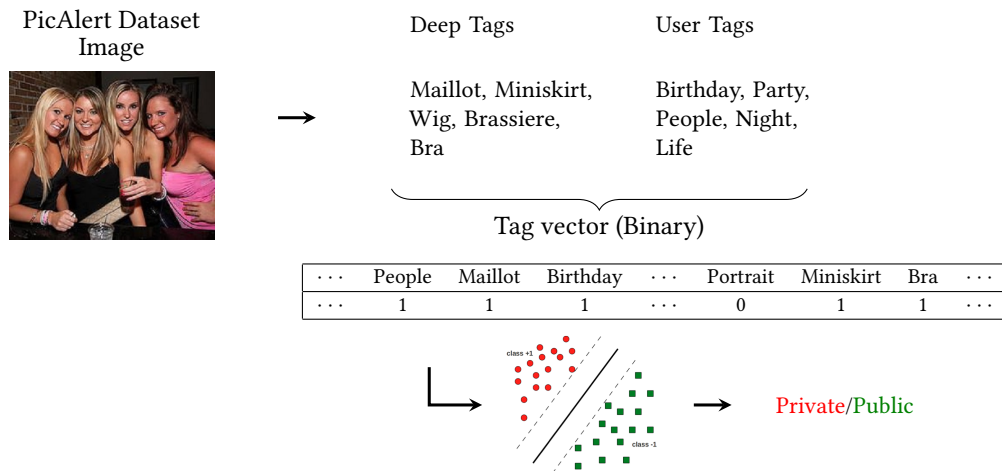
**Figure 4: Image encoding using tag features: We encode the combination of user tags and deep tags using binary vector representation, showing presence and absence of tags from tag vocabulary $V$. We set $1$ if a tag is present in the tag set or $0$ otherwise. We refer this model as Bag-of-Tags (BoT) model.**
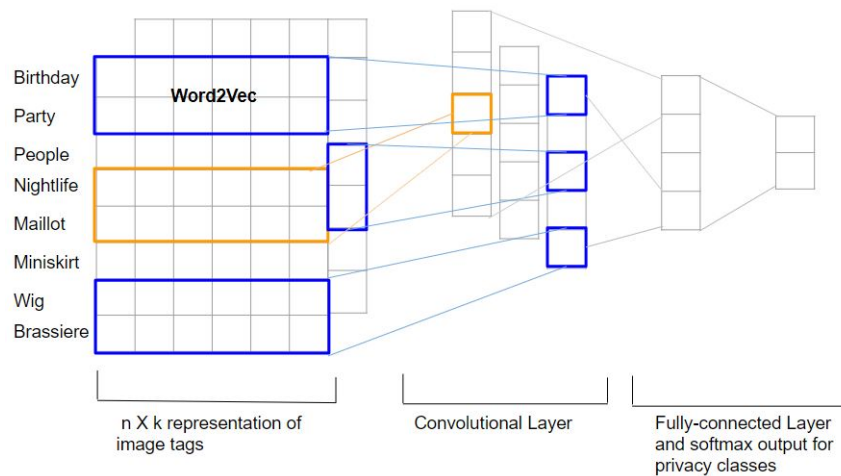


**Figure 5: Tag CNN architecture to classify an image as public or private using image tags.**

remove special characters and numbers from the user tags, as they do not provide any information with respect to privacy.

We combine deep tags and user tags and generate a binary vector representation for the tag set of an image, illustrating presence or absence of tags from tag vocabulary $V$. Particularly, we create a vector of size $|V|$, wherein, for all tags in the tag set, we set 1 on the position of the tag in the vocabulary ($V$) and 0 otherwise. We refer to this model as a Bag-of-Tags (BoT) model and show it's pictorial representation in Figure 4.

### 4.5 Tag CNN

CNN based models have achieved exceptional results for various NLP tasks such as semantic parsing [Yih et al. 2014], search query retrieval, sentence modeling [Kalchbrenner et al. 2014], sentence classification [Kim 2014], and other traditional NLP tasks [Collobert

et al. 2011]. Kim [2014] developed a CNN architecture for sentence level classification task. A sentence contains keywords in the form of objects, subjects, and verbs that help in the classification task. Image tags are nothing but keywords that are used to describe an image. Thus, for privacy prediction, we employ a CNN architecture that has proven adequate for sentence classification [Kim 2014].

The CNN architecture by Kim [2014] shown in Figure 5 is a slight variant of the CNN architecture of Collobert et al. [2011]. This architecture contains one layer of convolution on top of word vectors obtained from an unsupervised neural language model. The first layer embeds words (tags in our case) into the word vectors. The word vectors are first initialized with the word vectors that were trained on 100 billion words of Google News, given by Le and Mikolov [2014]. Words that are not present in the set of pre-trained words are initialized randomly. These word vectors are

then fine-tuned on the tags from the privacy dataset. The next layer performs convolutions on the embedded word vectors using multiple filter sizes of 3, 4 and 5, where we use 128 filters from each size and produce a tag feature representation. A max-pooling operation [Collobert et al. 2011] over a feature map is applied to take the maximum value of the features to capture the most important feature of each feature map. These features are passed to a fully connected softmax layer to obtain the probability distribution over privacy labels. An illustration of the Tag CNN model can be seen in Figure 5.

## 5  DATASET

We evaluated our approach on a subset of $32,000$ Flickr images sampled from the PicAlert dataset, made available by Zerr et al. [2012b,a]. PicAlert consists of Flickr images on various subjects, which are manually labeled as *public* or *private* by external viewers. The dataset contains photos uploaded on Flickr during the period from January to April 2010. The data have been labeled by six teams providing a total of 81 users of ages between 10 and 59 years. One of the teams included graduate computer science students working together at a research center, whereas other teams contained users of social platforms. Users were instructed to consider that their camera has taken these pictures and to mark them as "private," "public," or "undecidable." The guideline to select the label is given as private images belong to the private sphere (like self-portraits, family, friends, someone's home) or contain information that one would not share with everyone else (such as private documents). The remaining images are labeled as public. In case no decision could be made, the image was marked as undecidable. Each image was shown to at least two different users. In the event of disagreement, the photos were presented to additional users. We only consider images that are labeled as public or private.

For all experiments, our $32,000$ images dataset is split into train and test sets of $27,000$ and $5,000$ images, respectively. Each experiment is repeated five times with a different train/test split (obtained using five different random seeds), and the final results are averaged across the five runs. The public and private images are in the ratio of 3:1 in both train and test sets.

## 6  EXPERIMENTS, RESULTS AND OBSERVATIONS

In this section, we perform a broad spectrum of experiments to evaluate features extracted from various deep architectures in order to understand which architecture can capture the complex privacy characteristics and help to distinguish between privacy classes. We first choose the machine learning classifier between generative models, ensemble methods, and discriminative algorithms for privacy prediction. Then, we use the chosen classifier to examine the visual features extracted from all four deep architectures AlexNet, GoogLeNet, VGG-16, and ResNet pre-trained on object data. We further investigate these architectures by fine-tuning them on the privacy data. Next, we compare the performance of models trained on the highest performing features with that of the state-of-the-art models and baseline approaches for privacy prediction. Additionally, we show the performance of the deep tags obtained through all four pre-trained networks and also study the combination of

deep tags and user tags in details for privacy prediction. We show the tag performance in two settings: (1) Bag-of-Tags models and (2) Tag CNN. We analyze the most promising features derived from both visual and tag encodings for privacy classification. We also provide a detailed analysis of the most informative tags for privacy prediction. Finally, we show the performance of the models trained on the fusion of visual and most informative tag features.
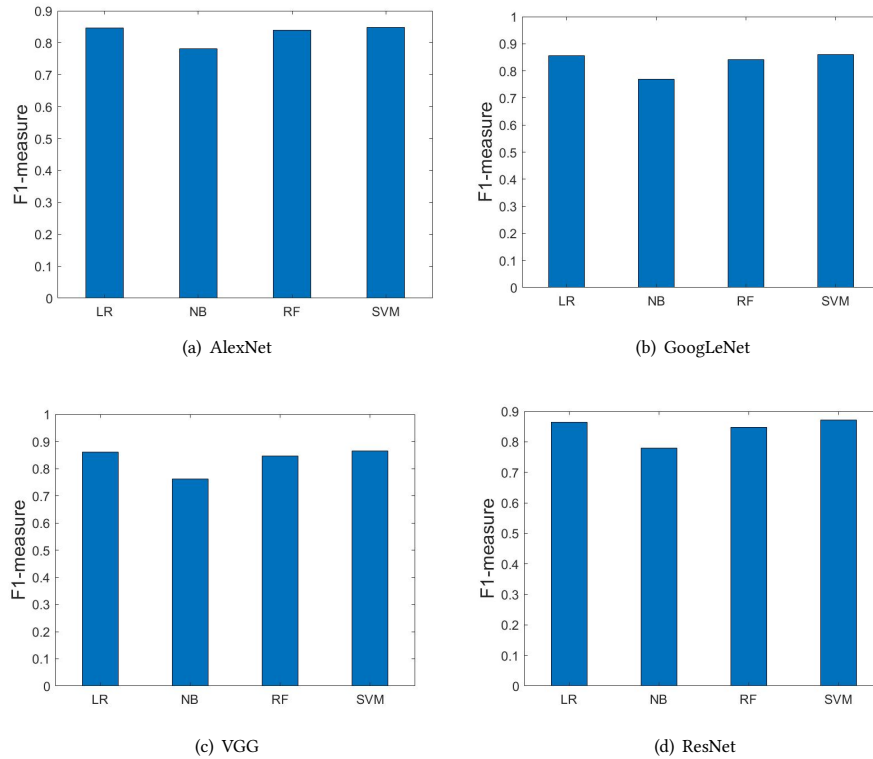
### 6.1  Classification Experiments for Features Derived From Pre-Trained CNNs

We first determine the classifier that works best with the features derived from the pre-trained CNNs. We study the performance of the features using the following classification algorithms: Naive Bayes (NB), Random Forest (RF), Logistic Regression (LR) and Support Vector Machine (SVM). NB is a generative model, whereas RF is an ensemble method using decision trees, and SVM and LR are discriminative algorithms. We evaluate the performance of these classifiers using the features derived from the last fully-connected layer of all the architectures, i.e., $fc_8$-A of AlexNet, $loss_3$-G of GoogLeNet, $fc_8$-V of VGG-16, and fc-R of ResNet. Figure 6 shows the performance of these classifiers in terms of F1-measure for all four architectures. From the figure, we notice that almost all the classifiers perform similarly except NB which performs worse. For example, for Alexnet, with NB we get an F1-measure of 0.781, whereas SVM obtains an F1-measure of 0.849. We can also observe that, generally, SVM and LR perform better than RF. For example, for ResNet, using SVM, we get an F1-measure of 0.872, whereas with RF we get an F1-measure of 0.848. SVM and LR perform comparably to each other for almost all the architectures except for ResNet. For ResNet, we obtain F1-measure of 0.872 and 0.865 using SVM and LR, respectively. The results of SVM over the LR classifier are statistically significant for p-values < 0.05. Thus, we chose to use SVM with the features derived from pre-trained CNNs for all of our next experiments.

To evaluate the proposed features, we used the SVM Weka implementation and chose the hyper-parameters that gave the best performance using 10-fold cross-validation on the training set. We experimented with $C = \{0.001, 0.01, 1.0, \cdots, 10.0\}$, kernels: Polynomial and RBF, the $\gamma$ parameter in RBF, and the degree $d$ of a polynomial. Hyper-parameters shown in all subsequent tables follow the format: "R/P,C,$\gamma/d$" where "R" denotes "RBF" and "P" denotes "Polynomial" kernel.

### 6.2  The Impact of the CNN Architecture on the Privacy Prediction

In this experiment, we aim to determine which architecture performs best for privacy prediction by investigating the performance of privacy prediction models based on visual semantic features extracted from all four architectures, AlexNet, GoogLeNet, VGG-16, and ResNet pre-trained on object data of ImageNet. We extract deep visual features: (1) $fc_6$-A, $fc_7$-A, $fc_8$-A and "prob-A" from AlexNet, (2) $loss_3$-G and "prob-G" from GoogLeNet, (3) $fc_6$-V, $fc_7$-V, $fc_8$-V and "prob-V" from VGG-16, and (4) fc-R and "prob-R" from ResNet. For AlexNet and GoogLeNet, we used the pre-trained networks that come with the CAFFE open-source framework for CNNs [Jia et al. 2014]. For VGG-16, we used an improved version of pre-trained

(a) AlexNet

(b) GoogLeNet

(c) VGG

(d) ResNet

**Figure 6: Performance of various classifiers (LR, NB, RF, SVM) using the features derived from all four architectures AlexNet, GoogLeNet, VGG, and ResNet.**

models presented by the VGG-16 team in the ILSVRC-2014 competition [Simonyan and Zisserman 2014]. For ResNet, we use the ResNet pre-trained models of 101 layers given by He et al. [2016a].
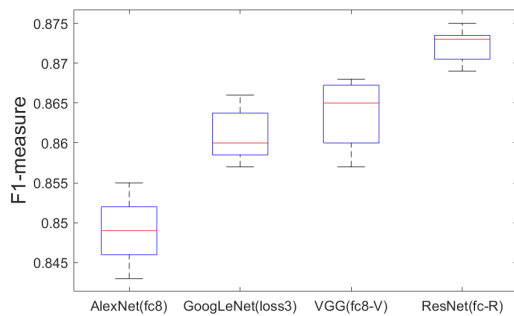
Table 1 shows the performance (Accuracy, F1-measure, Precision, Recall) of SVMs trained on the features extracted from all four pre-trained networks. From the table, we can observe that the models trained on the features extracted from ResNet consistently yield the best performance. For example, ResNet achieves an F1-measure of 0.872 as compared with 0.849, 0.861, 0.864 achieved by AlexNet, GoogLeNet, and VGG-16, respectively. These results suggest that the deep Residual Networks have more representational abilities compared to the other networks, and are more effective for predicting appropriate privacy classes of images. Additionally, ResNets are substantially deeper than their "plain" counterparts, which allows extracting various image-specific features that are beneficial for learning images' privacy characteristics better. Since privacy involves understanding the complicated relationship between the objects present in images, the features derived from ResNet prove to be more adequate than the features obtained by simply stacking convolutional layers. In Table 1, we also show the class-specific privacy prediction performance in order to identify which features characterize the private class effectively as sharing private images on the Web with everyone is not desirable. Interestingly, we found that the model trained on features obtained from ResNet provides improved F1-measure, precision, and recall for the private class.

Precisely, F1-measure for the private class improves from 0.661 (for AlexNet) to 0.717 (for ResNet), yielding an improvement of 6%. Similarly, for precision and recall, we obtain an increase of 4% and 7%, respectively, using ResNet features over the AlexNet features.

From Table 1, we also notice that the overall best performance (shown in orange and blue color) obtained for each network is higher than $\approx 85\%$ in terms of all compared measures (overall - Accuracy, F1-measure, precision and recall). Note that a naive baseline which classifies every image as "public" obtains an accuracy of 75%. Additionally, analyzing the results obtained by the VGG-16 features, we notice that as we ascend the fully-connected layers of the VGG-16 network from $fc_6$-V to $fc_8$-V, the F1-measure improves from 0.837 to 0.864 (see Table 1). Similarly, for AlexNet, the F1-measure improves from 0.82 (for $fc_6$-A) to 0.849 (for $fc_8$-A). This shows that the high-level object interpretations obtained through the last fully-connected layer helped to derive better privacy characteristics. Moreover, it is worth noting that "prob" features perform worse than the features extracted from the fully-connected layers (on all architectures). For example, prob-G obtains an F1-measure of 0.815, whereas $loss_3$-G achieves an F1-measure of 0.861. One possible explanation could be that squashing the values at the previous layer (e.g., $loss_3$-G in GoogleNet) through the softmax function, which yields the "prob" layer, produces a non-linearity that is less useful for SVM compared to the untransformed values. We also experimented with a combination of features, e.g., $fc_7$-A concatenated

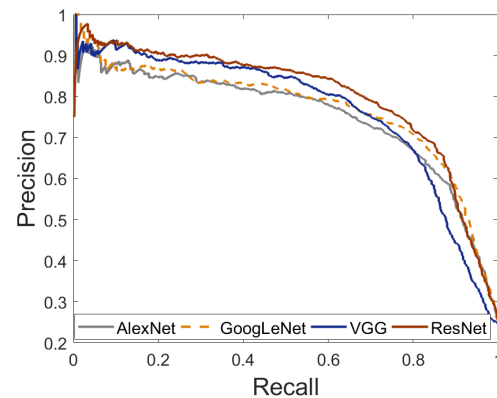| Features | H-Param | Acc % | Overall | | | Private | | | Public | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | F1 | Prec | Re | F1 | Prec | Re | F1 | Prec | Re |
| AlexNet | | | | | | | | | | | |
| $fc_6$-A | R,1.0,0.05 | 82.29 | 0.82 | 0.819 | 0.823 | 0.613 | 0.639 | 0.591 | 0.885 | 0.875 | 0.895 |
| $fc_7$-A | R,2.0,0.01 | 82.97 | 0.827 | 0.825 | 0.83 | 0.627 | 0.656 | 0.602 | 0.889 | 0.878 | 0.901 |
| $fc_8$-A | R,1.0,0.05 | *85.51* | *0.849* | *0.849* | *0.855* | *0.661* | *0.746* | *0.595* | *0.908* | *0.881* | *0.936* |
| prob-A | R,5.0,1.0 | 82.76 | 0.815 | 0.816 | 0.828 | 0.568 | 0.704 | 0.477 | 0.892 | 0.851 | 0.937 |
| GoogLeNet | | | | | | | | | | | |
| $loss_3$-G | P,0.001,2.0 | *86.42* | *0.861* | *0.86* | *0.864* | *0.695* | *0.746* | *0.652* | *0.913* | *0.895* | *0.93* |
| prob-G | R,50.0,0.05 | 82.66 | 0.815 | 0.816 | 0.827 | 0.573 | 0.694 | 0.488 | 0.891 | 0.853 | 0.933 |
| VGG-16 | | | | | | | | | | | |
| $fc_6$-V | R,1.0,0.01 | 83.85 | 0.837 | 0.836 | 0.839 | 0.652 | 0.67 | 0.636 | 0.895 | 0.888 | 0.902 |
| $fc_7$-V | R,2.0,0.01 | 84.43 | 0.843 | 0.842 | 0.844 | 0.663 | 0.684 | 0.644 | 0.899 | 0.891 | 0.907 |
| $fc_8$-V | R,2.0,0.05 | *86.72* | *0.864* | *0.863* | *0.867* | *0.7* | *0.758* | *0.65* | *0.915* | *0.895* | *0.935* |
| prob-V | R,2.0,0.05 | 81.72 | 0.801 | 0.804 | 0.817 | 0.528 | 0.687 | 0.429 | 0.887 | 0.84 | 0.939 |
| ResNet | | | | | | | | | | | |
| fc-R | R,1.0,0.05 | **87.58** | **0.872** | **0.872** | **0.876** | **0.717** | **0.783** | **0.662** | **0.92** | **0.899** | **0.943** |
| prob-R | R,2.0,0.05 | 80.6 | 0.784 | 0.789 | 0.806 | 0.473 | 0.67 | 0.366 | 0.881 | 0.826 | 0.943 |

**Table 1: Comparison of SVMs trained on features extracted from pre-trained architectures AlexNet, GoogLeNet, VGG-16 and ResNet. The best performance is shown in bold and blue color. The best performance for each network is shown in italics and orange color.**



**Figure 7: Box plot of F1-measure (overall) obtained for the best-performing features derived from each CNN over five splits.**



**Figure 8: Precision-recall curves for the private class obtained using features extracted from all four architectures AlexNet ($fc_8$), GoogLeNet ($loss_3$), VGG-16 ($fc_8$-V) and ResNet (fc-R).**

with $fc_8$-A, but we did not obtain a significant improvement over the individual ($fc_7$-A and $fc_8$-A) features.

We also analyze the performance by showing the box plots of F1-measure in Figure 7, obtained for the most promising features of all the architectures over the five random splits of the dataset. The figure indicates that the model trained on ResNet features is statistically significantly better than the models that are trained on the

features derived from the other architectures. We further compare features derived through all the architectures using precision-recall curves given in Figure 8. The curves show again that features derived from ResNet perform better than the features obtained from the other architectures, for a recall ranging from 0.5 to 0.8. For example, for a recall of 0.7, we achieve a precision of 0.75, 0.8, 0.8 and 0.85 for AlexNet, GoogLeNet, VGG-16, and ResNet, respectively.

| Features | H-Param | Acc % | Overall | | | Private | | | Public | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | F1 | Prec | Re | F1 | Prec | Re | F1 | Prec | Re |
| Fine-tuned AlexNet | | | | | | | | | | | |
| ft-A | fc | 85.01 | 0.846 | 0.845 | 0.851 | 0.657 | 0.723 | 0.606 | 0.904 | 0.883 | 0.926 |
| ft-A | fc-all | 85.14 | 0.849 | 0.847 | 0.852 | 0.669 | 0.713 | 0.632 | 0.904 | 0.889 | 0.92 |
| ft-A | all | 85.07 | 0.848 | 0.847 | 0.851 | 0.67 | 0.707 | *0.638* | 0.904 | 0.89 | 0.917 |
| Pre-trained AlexNet | | | | | | | | | | | |
| $fc_8$-A | R,1,0.05 | 85.51 | 0.849 | 0.849 | 0.855 | 0.661 | *0.746* | 0.595 | 0.908 | 0.881 | 0.936 |
| Fine-tuned GoogLeNet | | | | | | | | | | | |
| ft-G | fc | 86.27 | 0.86 | 0.859 | 0.863 | 0.694 | 0.74 | 0.653 | 0.911 | 0.895 | 0.928 |
| ft-G | all | 86.77 | 0.867 | 0.867 | 0.868 | *0.717* | 0.732 | *0.705* | 0.914 | 0.909 | 0.919 |
| Pre-trained GoogLeNet | | | | | | | | | | | |
| $loss_3$-G | P,0.001,2 | 86.42 | 0.861 | 0.86 | 0.864 | 0.695 | 0.746 | 0.652 | 0.913 | 0.895 | 0.930 |
| Fine-tuned VGG-16 | | | | | | | | | | | |
| ft-V | fc | 86.74 | 0.864 | 0.865 | 0.869 | 0.695 | *0.782* | 0.631 | 0.916 | 0.891 | *0.944* |
| ft-V | fc-all | 86.92 | 0.869 | 0.87 | 0.869 | **0.722** | 0.73 | **0.717** | 0.914 | **0.912** | 0.917 |
| ft-V | all | 86.76 | 0.867 | 0.867 | 0.868 | 0.718 | 0.729 | 0.709 | 0.913 | 0.91 | 0.917 |
| Pre-trained VGG-16 | | | | | | | | | | | |
| $fc_8$-V | R,2,0.05 | 86.72 | 0.864 | 0.863 | 0.867 | 0.700 | 0.758 | 0.65 | 0.915 | 0.895 | 0.935 |
| Fine-tuned ResNet | | | | | | | | | | | |
| ft-R | fc | 87.23 | 0.87 | 0.869 | 0.873 | 0.717 | 0.759 | *0.68* | 0.918 | 0.903 | 0.932 |
| ft-R | all | 86.19 | 0.856 | 0.856 | 0.863 | 0.672 | 0.776 | 0.594 | 0.913 | 0.881 | **0.946** |
| Pre-trained ResNet | | | | | | | | | | | |
| fc-R | R,1,0.05 | **87.58** | **0.872** | **0.872** | **0.876** | 0.717 | **0.783** | 0.662 | **0.92** | 0.899 | 0.943 |

**Table 2: Fine-tuned networks vs. Pre-trained networks. The best performance is shown in bold and blue color. The performance measures that achieve a better performance after fine-tuning a CNN over pre-trained features are shown in italics and orange color.**

## 6.3 Fine-Tuned Networks vs. Pre-Trained Networks

Previous works showed that the features transferred from the network pre-trained on the object dataset to the privacy data achieved a good performance [Tran et al. 2016]. Moreover, many other works used "transfer learning" to get more dataset specific features [Bengio 2012; Donahue et al. 2013]. Thus, we determine the performance of fine-tuned networks on the privacy dataset. We compare fine-tuned networks of all four architectures with the deep features obtained from pre-trained networks. We refer the fine-tuned networks of AlexNet, GoogLeNet, VGG-16, and ResNet as "ft-A," "ft-G," "ft-V," and "ft-R" respectively. For fine-tuning, we used the same CNN architectures pre-trained on the object dataset, and employed in previous experiments. To fine-tune the networks, we experiment with the three types of settings: (1) fine-tune the last fully-connected layer (that has two output units corresponding to 2 privacy classes) with higher learning rates as compared to the learning rates of the rest of the layers of the networks (0.001 vs. 0.0001), referred as "fc." (2) fine-tune all the fully-connected layers of the networks with higher learning rates and convolutional layers are learned with smaller learning rates. We refer to this setting as "fc-all." (3)

fine-tune all layers with the same learning rates and denoted as "all." Note that since ResNet and GoogLeNet have only one fully-connected layer, we report the performance obtained only using "fc," and "all" settings. The very low learning rate avoids substantial learning of the pre-trained layers. In other words, due to a very low learning rate (0.0001), pre-trained layers learn very slowly as compared to the layers that have a higher learning rate (0.001) to learn the required weights for privacy data.

Table 2 shows the performance comparison of the models obtained by fine-tuning architectures on privacy data and the models trained on the features derived from the pre-trained networks. We notice that we get mostly similar results when we fine-tune pre-trained models on our privacy dataset as compared to the models trained on the features derived from the pre-trained architectures. However, we get improved recall for the private class when we fine-tune the networks on the privacy dataset. For example, the fine-tuned VGG-16 network gets an improvement of 6.7% in the recall for the private class (see ft-V, fc-all setting vs. $fc_8$-V) over the models trained on the features extracted from the pre-trained VGG-16. The performance measures that achieve a better performance after fine-tuning a CNN over pre-trained features are shown in italics and orange color for each network. We notice that the

| Features | H-Param | Acc % | Overall | | | Private | | | Public | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | F1 | Prec | Re | F1 | Prec | Re | F1 | Prec | Re |
| Highest performing CNN architecture | | | | | | | | | | | |
| fc-R | R,1.0,0.05 | **87.58** | **0.872** | **0.872** | **0.876** | **0.717** | **0.783** | **0.662** | **0.92** | 0.899 | **0.943** |
| #1 PCNH framework | | | | | | | | | | | |
| PCNH | – | 83.13 | 0.824 | 0.823 | 0.831 | 0.624 | 0.704 | 0.561 | 0.891 | 0.863 | 0.921 |
| #2 AlexNet Deep Features | | | | | | | | | | | |
| $fc_8$-A | R,1.0,0.05 | 85.51 | 0.849 | 0.849 | 0.855 | 0.661 | 0.746 | 0.595 | 0.908 | 0.881 | 0.936 |
| #3 SIFT & GIST models | | | | | | | | | | | |
| SIFT | P,1.0,2.0 | 77.31 | 0.674 | 0.598 | 0.773 | 0.002 | 0.058 | 0.001 | 0.87 | 0.772 | 0.995 |
| GIST | R,0.001,0.5 | 77.33 | 0.674 | 0.598 | 0.773 | 0.002 | 0.058 | 0.001 | 0.87 | 0.772 | 0.995 |
| SIFT & GIST | R,0.05,0.5 | 72.67 | 0.704 | 0.691 | 0.727 | 0.27 | 0.343 | 0.223 | 0.832 | 0.793 | 0.874 |
| #4 Rule-based models | | | | | | | | | | | |
| Rule-1 | – | 77.35 | 0.683 | 0.694 | 0.672 | 0.509 | 0.47 | 0.556 | 0.853 | 0.875 | 0.832 |
| Rule-2 | – | 77.93 | 0.673 | 0.704 | 0.644 | 0.458 | 0.373 | 0.593 | 0.897 | **0.914** | 0.88 |

**Table 3: Highest performing visual features (fc-R) vs. Prior works.**

fine-tuned VGG gives the best performance for the F1-measure and recall of the private class (shown in bold and blue color). However, the models trained on the features derived from the pre-trained ResNet yield the best overall performance (shown in bold and blue color). Thus, we compare the models trained on fc-R features with prior privacy prediction approaches in the next subsection.

## 6.4 ResNet Features-Based Models vs. Prior Works

We compare the performance of the state-of-the-art works on privacy prediction, as detailed below, with the models trained using ResNet features, i.e., fc-R.
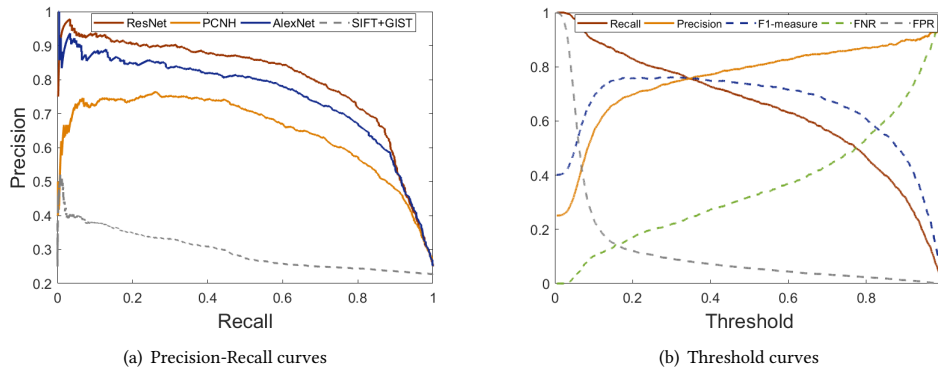
**1. PCNH privacy framework** [Tran et al. 2016]: This framework combines features obtained from two architectures: one that extracts convolutional features (size = 24, referred as Convolutional CNN), and another that extracts object features (size = 24, referred as Object CNN). The Convolutional CNN contains two convolutional layers and three fully-connected layers of size 512, 512, 24, respectively. On the other hand, the object CNN is an extension of AlexNet architecture that appends three fully-connected layers of size 512, 512, and 24, at the end of the last fully-connected layer of AlexNet and form a deep network of 11 layers. The two CNNs are connected at the output layer. The PCNH framework is first trained on the ImageNet dataset and then fine-tuned on a small privacy dataset.

**2. AlexNet features** [Tonge and Caragea 2015, 2016, 2018]: We consider the model trained on the features extracted from the last fully-connected layer of AlexNet, i.e., $fc_8$-A as another baseline, since in our previous works we achieved a good performance using these features for privacy prediction.

**3. SIFT & GIST** [Squicciarini et al. 2014, 2017a; Zerr et al. 2012b]: We also consider classifiers trained on the best performing features between SIFT, GIST, and their combination as our baselines. Our choice of these features is motivated by their good performance over other visual features such as colors, patterns, and edge directions in prior works [Squicciarini et al. 2014; Zerr et al. 2012b]. For SIFT, we construct a vocabulary of 128 visual words for our experiments. We tried different numbers of visual words such as 500, 1000, etc., but we did not get a significant improvement over the 128 visual words. For a given image, GIST is computed by first convolving the image with 32 Gabor filters at 4 scale and 8 orientations, which produces 32 feature maps; second, dividing the feature map into a $4 \times 4$ grid and averaging feature values of each cell; and third, concatenating these 16 averaged values for 32 feature maps, which results in a feature vector of 512 ($16 \times 32$) length.

**3. Rule-based classifiers**: We also compare the performance of models trained on ResNet features fc-R with two rule-based classifiers which predict an image as *private* if it contains persons. Otherwise, the image is classified as *public*. For the first rule-based classifier, we detect front and profile faces by using Viola-Jones algorithm [Viola and Jones 2001]. For the second rule-based classifier, we consider user tags such as "women," "men," "people." Recall that these tags are not present in the set of 1, 000 categories of the ILSVRC-2012 subset of the ImageNet dataset, and hence, we restrict to user tags only. If an image contains one of these tags or detects a face, we consider it as "private," otherwise "public."

Table 3 compares the performance of models trained on fc-R features (the highest performing features obtained from our previous experiments) with the performance obtained by prior works. As can be seen from the table, the deep features extracted from the pre-trained ResNet achieve the highest performance, and hence,

(a) Precision-Recall curves

(b) Threshold curves

**Figure 9: Precision-Recall and Threshold curves for the private class obtained using ResNet features (fc-R) and prior works.**

are able to learn the privacy characteristics better than the prior works with respect to both the classes. Precisely, using fc-R features, F1-measure improves from 0.824 obtained by PCNH framework to 0.872 obtained by fc-R features, providing an overall improvement of 5%. Moreover, for the private class, fc-R features yield an improvement of 9.8% in F1-measure over the more sophisticated PCNH framework (from 0.624, PCNH to 0.717, fc-R features).

One possible explanation could be that the object CNN of PCNH framework is formed by appending more fully-connected layers to the AlexNet architecture and the increase in the number of complex non-linear layers (fully-connected layers) introduces more parameters to learn. At the same time, with a relatively small amount of training data (PicAlert vs. ImageNet), the object CNN model can over-fit. On the other hand, as images' privacy greatly depends on the objects in images, we believe that the low-level features controlling the distinct attributes of the objects (e.g., edges of swimsuit vs. short pants) obtained through the convolutional layers can better approximate the privacy function compared with adding more non-linear layers (as in PCNH). This is justified by the results, where the network with more convolutional layers, i.e., ResNet achieves a better performance as compared to the network with more fully-connected layers, i.e., PCNH. Additionally, even though PCNH attempted to capture convolutional features using Convolutional CNN, both CNN (convolutional and object) vary in their discriminative power and thus obtaining an optimal unification of convolutional CNN and object CNN is difficult. Moreover, PCNH is required to first train on ImageNet and then fine-tune on the PicAlert dataset. Training a deep network such as PCNH two times significantly increases the processing power and time. On the other hand, through our experiments, we found that the features derived from the state-of-the-art ResNet model can reduce the overhead of re-training and achieve a better performance for privacy prediction.

As discussed before, the models trained on ResNet features outperform those trained on AlexNet features. Interestingly, the best performing baseline among all corresponds to the SVM trained on the deep features extracted from the AlexNet architecture. For example, the SVM trained on the AlexNet features (fc$_8$-A) yields an F1-measure of 0.849 as compared with the F1-measure of 0.824 achieved by the PCNH framework. We hypothesize that this is due to the model complexity and the small size of the privacy dataset

used to train the PCNH framework. For example, merging two CNNs (as in PCNH) that vary in depth, width, and optimization algorithm can become very complex and thus the framework potentially has more local minima, that may not yield the best possible results. Additionally, unlike Tran et al. [2016], that used 800 images in their evaluation, we evaluate the models on a large set of images (32000), containing a large variety of image subjects. The features derived from the various layers of the state-of-the-art AlexNet reduce the overhead of training the complex structure and still achieve a good performance for privacy prediction.

Another interesting aspect to note is that, although we showed earlier that the fine-tuned network (in this case VGG-16) does not show a significant improvement over the ResNet pre-trained features (see Table 2), our fine-tuning approach yields better results compared to the PCNH framework. For example, fine-tuned VGG-16 (ft-V) achieves an F1-measure of 0.869 whereas PCNH achieves an F1-measure of 0.824 (see Tables 2 and 3). The possible reasons could be that we use a larger privacy dataset to fine-tune a simpler architecture, unlike PCNH that merges two convolutional neural networks. Additionally, we fine-tune the state-of-the-art VGG-16 model presented by Simonyan and Zisserman [2014], contrary to PCNH that required estimating optimal network parameters to train the merged architecture on the ImageNet dataset.

As expected, we can see from Table 3 that the baseline models trained on SIFT/GIST and the rule-based models are the lowest performing models. For example, the fc-R based models achieve improvement in performance as high as 17% over SIFT/GIST models. With a paired T-test, the improvements over the prior approaches for F1-measure are statistically significant for p-values < 0.05. It is also interesting to note that rules based on facial features exhibit better performance than SIFT and GIST and suggest that feature representing persons are helpful to predict private images. However, fc-R features outperform the rule-based models based on facial features by more than 10% in terms of all measures.

We further analyze fc-R features and compare their performance with the prior works through precision-recall curves shown in Figure 9 (a). As can be seen from the figure, the SVM trained on ResNet features achieve a precision of ≈ 0.8 for recall values up to 0.8, and after that, the precision drops steadily.

| Features | H-Param | Acc % | Overall | | | Private | | | Public | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | F1 | Prec | Re | F1 | Prec | Re | F1 | Prec | Re |
| Best performing CNN architecture | | | | | | | | | | | |
| fc-R | R,1.0,0.05 | **87.58** | **0.872** | **0.872** | **0.876** | **0.717** | **0.783** | 0.662 | **0.92** | 0.899 | **0.943** |
| #1 User Tags (BoT) | | | | | | | | | | | |
| UT | R,2.0,0.05 | 78.63 | 0.777 | 0.772 | 0.786 | 0.496 | 0.565 | 0.442 | 0.865 | 0.837 | 0.894 |
| #2 Deep Tags (BoT) | | | | | | | | | | | |
| DT-A | R,1.0,0.1 | 83.34 | 0.825 | 0.824 | 0.833 | 0.601 | 0.699 | 0.529 | 0.895 | 0.863 | 0.929 |
| DT-G | R,1.0,0.05 | 83.59 | 0.828 | 0.827 | 0.836 | 0.606 | 0.699 | 0.534 | 0.896 | 0.866 | 0.929 |
| DT-V | P,1.0,1.0 | 83.42 | 0.826 | 0.825 | 0.834 | 0.607 | 0.698 | 0.537 | 0.895 | 0.865 | 0.927 |
| DT-R | P,1.0,1.0 | 83.78 | 0.833 | 0.831 | 0.838 | 0.631 | 0.688 | 0.584 | 0.896 | 0.876 | 0.917 |
| #3 User Tags & Deep Tags | | | | | | | | | | | |
| UT+DT-R (BoT) | R,1.0,0.05 | 84.33 | 0.84 | 0.839 | 0.843 | 0.67 | 0.709 | 0.636 | 0.897 | 0.882 | 0.913 |
| Tag CNN | – | 85.13 | 0.855 | 0.855 | 0.854 | 0.706 | 0.700 | **0.712** | 0.901 | **0.903** | 0.898 |

**Table 4: Visual features vs. Tag features.**

The performance measures shown in previous experiments are calculated using a classification threshold of 0.5. In order to see how the performance measures vary for different classification thresholds, we plot the threshold curve and show this in Figure 9 (b). From the figure, we can see that the precision increases from $\approx 0.68$ to $\approx 0.97$ at a slower rate with the classification threshold. The recall slowly decreases to 0.8 for a threshold value of $\approx 0.4$, and the F1-measure remains comparatively constant until $\approx 0.75$. At a threshold of $\approx 0.4$, we get equal precision and recall of $\approx 0.78$, which corresponds to the breakeven point. In the figure, we also show the false negative rate and false positive rate, so that depending on a user's need (high precision or high recall), the classifier can run at the desired threshold. Also, to reduce the number of content-sensitive images shared with everyone on the Web, lower false negative (FN) rates are desired. From Figure 9 (b), we can see that we achieve lower FN rates up to $\approx 0.4$ for the threshold values up to 0.8.

## 6.5 Best Performing Visual Features vs. Tag Features

Image tags provide relevant cues for privacy-aware image retrieval [Zerr et al. 2012b] and can become an essential tool for surfacing the hidden content of the deep Web without exposing sensitive details. Additionally, previous works showed that user tags performed better or on par compared with visual features [Squicciarini et al. 2014; Tonge and Caragea 2015, 2016, 2018; Zerr et al. 2012b]. For example, in our previous work [Tonge and Caragea 2015, 2016, 2018], we showed that the combination of user tags and deep tags derived from AlexNet performs comparably to the AlexNet based visual features. Hence, in this experiment, we compare the performance of fc-R features with the tag features. For deep tags, we follow the same approach as in our previous work [Tonge and Caragea 2015, 2016, 2018] and consider the top $k = 10$ object labels since $k = 10$ worked best. "DT-A," "DT-G," "DT-V," and "DT-R" denote deep tags generated by AlexNet, GoogleNet, VGG-16, and ResNet, respectively. Deep tags are generated using the probability distribution
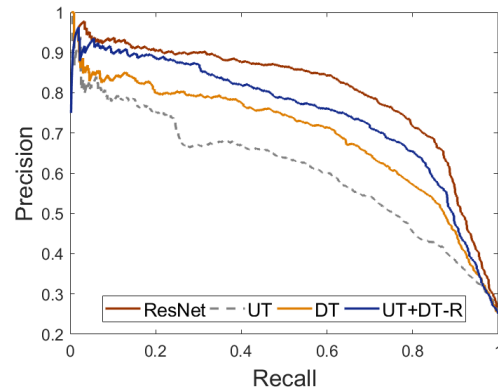


**Figure 10: Precision-Recall curves for the private class obtained using visual features (fc-R) and tag features as user tags (UT), deep tags (DT-R), the combination of user tags and deep tags (UT + DT-R).**

over 1,000 object categories for the input image obtained by applying the softmax function over the last fully-connected layer of the respective CNN.

Table 4 compares the performance obtained using models trained on fc-R features with the performance of models trained on the tag features. We consider tag features as: (1) user tags (UT); (2) deep tags (DT) obtained from all architectures; (3) the combination of user tags and best performing deep tag features using Bag-of-Tags (BoT) model; and (4) Tag CNN applied to the combination of user and deep tags. As can be seen from the table, the visual features extracted from ResNet outperform the user tags and deep tags independently as well as their combination. The models trained on fc-R features achieve an improvement of 2% over the CNN trained on the combination of user tags and deep tags (Tag CNN). Additionally, the models trained on fc-R features yield an increase of 9.5% in the F1-measure over the user tags alone and an increase of 4% over the
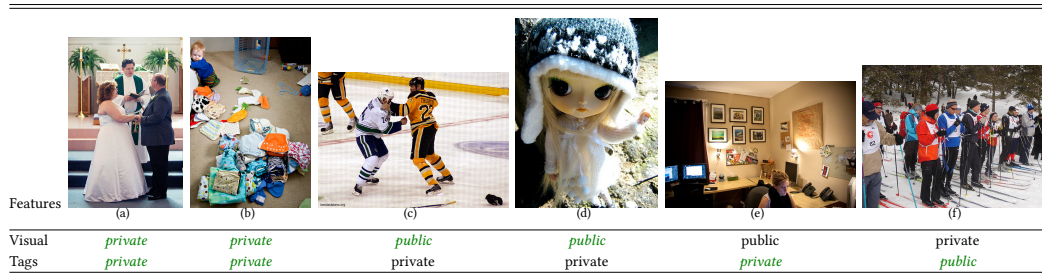
| Features | (a) | (b) | (c) | (d) | (e) | (f) |
|---|---|---|---|---|---|---|
| Visual | *private* | *private* | *public* | *public* | public | private |
| Tags | *private* | *private* | private | private | *private* | *public* |

**Figure 11: Privacy predictions obtained by image content encodings.**

| Rank 1-10 | Rank 11-20 | Rank 21-30 | Rank 31-40 | Rank 41-50 |
|---|---|---|---|---|
| **people** | pyjama | maillot | **promontory** | jersey |
| wig | jammies | **girl** | t-shirt | mole |
| **portrait** | sweatshirt | suit of clothes | foreland | groin |
| bow-tie | **outdoor** | ice lolly | **headland** | bulwark |
| neck brace | **lakeside** | suit | bandeau | seawall |
| **groom** | **lakeshore** | lollipop | miniskirt | **seacoast** |
| **bridegroom** | sun blocker | two-piece | breakwater | **indoor** |
| laboratory coat | sunscreen | tank suit | **vale** | stethoscope |
| hair spray | sunglasses | bikini | hand blower | **valley** |
| shower cap | military uniform | swimming cap | **jetty** | **head** |

**Table 5: Top** 50 **highly informative tags. We use the combination of deep tags and user tags (DT+UT) to calculate the information gain. User tags are shown in bold.**

best performing deep tags, i.e., DT-R (among the deep tags of the four architectures).

From Table 4, we also observe that the Tag CNN performs better than the Bag-of-Tags model (DT-R+UT), yielding an improvement of 3.0% in the F1-measure of private class. Additionally, even though the visual features (fc-R) yield overall a better performance than the tag features, for the private class, the F1-measure (0.717) of the visual features (fc-R) is comparable to the F1-measure (0.706) of the Tag CNN. It is also interesting to note that the Visual CNN (fc-R) achieves an increase of 8% in the precision (private class) over the Tag CNN whereas the Tag CNN obtains an improved recall (private class) of 5% over the Visual CNN.

In order to see how precision varies for different recall values, we also show the precision-recall curves for the visual and tag features in Figure 10. To avoid clutter we show the precision-recall curves for deep tags derived through ResNet and the combination of user tags and deep tags (DT-R) using BoT model. From the curves, we can see that the ResNet visual features perform better than the tag features, for a wide range of recall values from 0.3 to 0.8.

We further analyze both the type of image encodings (visual & tag) by examining the privacy predictions obtained for anecdotal examples using both the encodings.

*6.5.1 Anecdotal Examples:* In order to understand the quality of predictions obtained by visual and tag features, we show privacy predictions for some samples obtained by both type of features. Figure 11 shows the predictions obtained using SVM models trained on the visual features and those trained on the combination of user tags and deep tags. Correct predictions are shown in italics

and green in color. We can see that for images (a) and (b), the models trained on image tags (UT+DT) and visual features provide correct predictions. The tags such as "groom," "bride," "wedding," "photography" describe the picture (a) adequately, and hence, using these tags appropriate predictions are obtained. Similarly, visual features identify the required objects, and a relationship between the objects and provide an accurate prediction for these images. Consider now examples (c) and (d). For these images, visual features capture the required objects to make accurate predictions, whereas, image tags such as "bruins," "fight," of image (c) and "cute," "wool," "bonnet" of image (d) do not provide adequate information about the picture and hence, yield an incorrect prediction. However, tags such as "hockey," "sports" for image (c) and "toy," "doll" for image (d) would have helped to make an appropriate prediction. We also show some examples, (e) and (f), for which visual features fail to predict correct privacy classes. Particularly, for image (f), we notice that visual features capture object information that identifies the image as private. On the other hand, the image tags such as "festival" and "sport" (describing the scene) provide additional information (over the object information) that helps the tag-based classifier to identify the picture as public.

Next, we provide the detailed analysis of image tags with respect to privacy.

*6.5.2 Analysis of Image Tags with Respect to Privacy Classes:* We provide an analysis of the deep tags (capturing the visual content of the image) and user tags to learn their correlation with the private and public classes. First, we rank user tags and deep tags based on their information gain on the train set. Table 5 shows

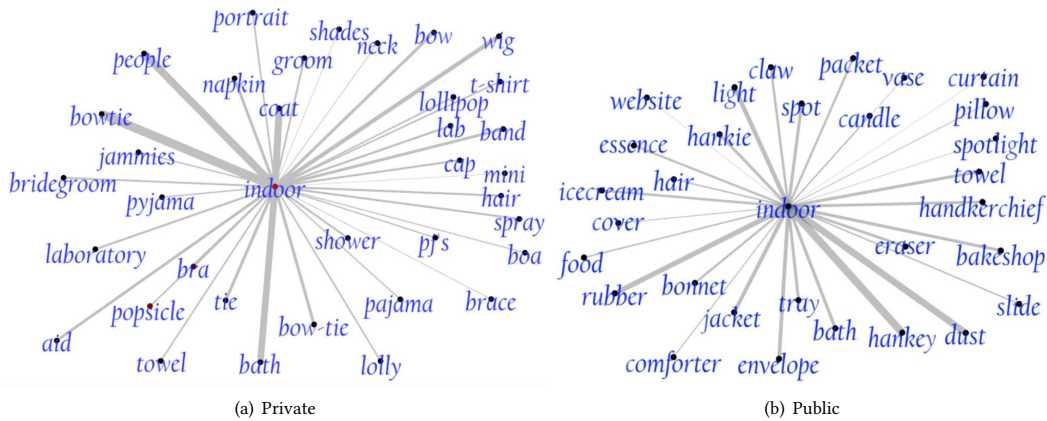**Figure 12: High frequency tag clouds with respect to public and private images.**



**Figure 13: Tag association graph.**

top 50 tags with high information gain. From the table, we observe that the tags such as "maillot," "two-piece," "sandbar" provide high correlation to the privacy classes. We also notice that deep tags (objects) contribute to a significant section of top 50 highly informative tags. Secondly, we rank both the tags (user and deep tags) based on their frequency in public and private classes. We show 50 most frequent tags for each privacy class using word clouds in Figure 12. The tags with larger word size depict a higher frequency of the tag. We notice that tags such as "indoor," "people," "portrait" occur more frequently in the private class, whereas tags such as "outdoor," "lakeside," "fountain," occur more frequently in the public class.

We also observe that some informative tags overlap in both public and private clouds (See Figure 12, e.g., "indoor"). Thus, we analyze other tags that co-occur with the overlapping tags to further discriminate between their association with the public and private classes. To inspect the overlapping tags, we create two graphs with respect to public and private classes. For the public graph, we consider each tag as a node in the graph and draw an edge between

the two nodes if both the tags belong to the same public image. Likewise, we construct another graph using private images. Figure 13 shows portions of both public and private graphs for "indoor" tag. To reduce the complexity of visualization, we only display nodes with stronger edges that have the co-occurrence greater than a certain threshold. Note that stronger edges (edges with higher width) represent the high co-occurrence coefficient between two nodes (in our case, tags). From the graphs, we observe that the overlapping tag "indoor" tends to have different highly co-occurring tags for public and private classes. For example, the "indoor" tag shows high co-occurrence with tags such as "people," "bath," "coat," "bowtie," "bra" (tags describing private class) in the private graph. On the other hand, in the public graph, the tag shows high co-occurrence with "dust," "light," "hankey," "bakeshop," "rubber," and so forth (tags describing public class). Even though some tags in the graph have comparatively low co-occurrence, the tags occurring in the private graph tend to associate with the private class whereas the tags from the public graph are more inclined towards the public class.
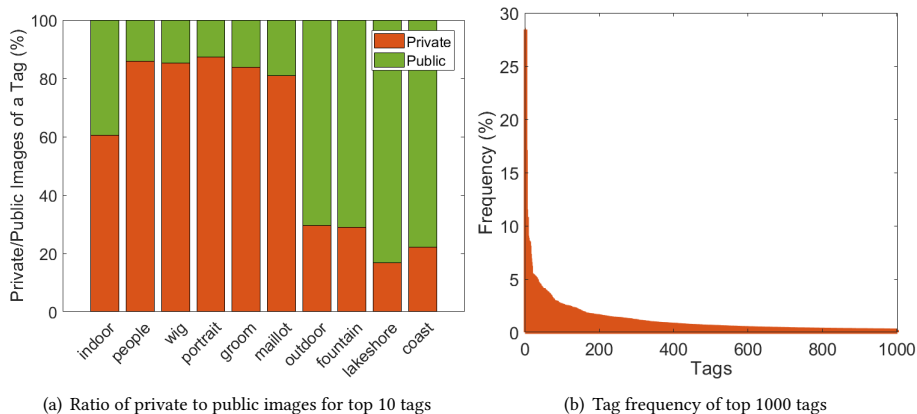
(a) Ratio of private to public images for top 10 tags

(b) Tag frequency of top 1000 tags

**Figure 14: Analysis of top frequently occurring tags.**

| Features | Overall | | | | Private | | | Public | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Acc % | F1 | Prec | Re | F1 | Prec | Re | F1 | Prec | Re |
| fc-R | 87.58 | 0.872 | 0.872 | 0.876 | 0.717 | 0.783 | 0.662 | 0.92 | 0.899 | **0.943** |
| fc-R+UT | **88.29** | **0.881** | **0.88** | **0.883** | **0.753** | **0.799** | **0.713** | **0.923** | **0.907** | 0.94 |

**Table 6: Results for the combination of visual and tag features.**

We further analyze the privacy differences of top 10 private and public image subjects. We consider "outdoor," "indoor," "fountain," "lakeshore," and "coast" for the public class. On the other hand, we consider "indoor," "people," "wig," "portrait," "outdoor," "groom," and "maillot" for the private class. Note that since images may have various tags associated with them, an image can be counted towards more than one tag. Given that the dataset contains three times more public images than private images (3 : 1 public to private ratio), we count 3 for each private image as opposed to the public class where we count 1 for each public image for a fair comparison. The ratio of private to public content for a specific tag is shown in Figure 14 (a). For example, out of the total images that possess the "indoor" tag, 60% images are of private class. From the figure, we observe that tags except for "indoor" show a significant difference in the inclination towards public and private classes. We also plot the frequency of top 1000 tags normalized by the dataset size in Figure 14 (b). The plot shows that the top 200 tags befall in 3% − 30% of the dataset with very few tags occurring in around 20% of the dataset. We also observe that most of the tags lie below 3% of the dataset showing the variation in the images' subjects and complexity of the dataset which justifies the fact that increasing the number of images increases the challenges of the problem statement.

## 6.6 Fusion of Visual and Tag Features for Image Privacy Prediction

Visual encoding and tag encoding capture different aspects of images. Thus, we add the top 350 correlated tags to the visual features fc-R and evaluate their performance for privacy prediction. We experiment with the number of top correlated tags = $\{10, 20, \cdots, 50, 100, \cdots, 500, 1000, 5000, 10000\}$. However, we get the best results with the top 350 correlated tags. Table 6 shows the

results obtained using SVMs trained on fc-R and the combination of fc-R with the top 350 correlated user tags (fc-R+tag). The results reveal that adding the highly correlated tags improves the privacy prediction performance. Precisely, we get a significant improvement of 4% on F1-measure of private class over the performance obtained using visual features fc-R. Note that, in our previous works [Tonge and Caragea 2015, 2016, 2018] and Experiment 6.5 (where we compare visual and tag features), we described visual content using tags (deep tags) and combined with the user tag to achieve a better performance. However, the combination of user tags and deep tags (combining one type of encoding) yields a lower performance as compared to the combination of user tags and fc-R features (combining two types of encodings). Precisely, the combination of user tags (UT) and fc-R features yields an improvement of 5% in the F1-measure of private class (refer Tables 4 and 6) over the combination of user tags and deep tags.

## 7 CONCLUSION

In this paper, we provide a comprehensive study of the deep features derived from various CNN architectures of increasing depth to discover the best features that can provide an accurate privacy prediction for online images. Specifically, we explored features obtained from various layers of the pre-trained CNNs such as AlexNet, GoogLeNet, VGG-16, and ResNet and used them with SVM classifiers to predict an image's privacy as *private* or *public*. We also fine-tuned these architectures on a privacy dataset. The study reveals that the SVM models trained on features derived from ResNet perform better than the models trained on the features derived from AlexNet, GoogLeNet, and VGG-16. We found that the overall performance obtained using models trained on the features derived through pre-trained networks is comparable to the fine-tuned

architectures. However, fine-tuned networks provide improved performance for the private class as compared to the models trained on pre-trained features. The results show remarkable improvements in the performance of image privacy prediction as compared to the models trained on CNN-based and traditional baseline features. Additionally, models trained on the deep features outperform rule-based models that classify images as private if they contain people. We also investigate the combination of user tags and deep tags derived from CNN architectures in two settings: (1) using SVM on the bag-of-tags features; and (2) applying the text CNN over these tags. We thoroughly compare these models with the models trained on the highest performing visual features obtained for privacy prediction. We further provide a detailed analysis of tags that gives insights for the most informative tags for privacy predictions. We finally show that the combination of deep visual features with these informative tags yields improvement in the performance over the individual sets of features (visual and tag).

The result of our classification task is expected to aid other very practical applications. For example, a law enforcement agent who needs to review digital evidence on a suspected equipment to detect sensitive content in images and videos, e.g., child pornography. The learning models developed here can be used to filter or narrow down the number of images and videos having sensitive or private content before other more sophisticated approaches can be applied to the data. Consider another example, images today are often stored in the cloud (e.g., Dropbox or iCloud) as a form of file backup to prevent their loss from physical damages and they are vulnerable to unwanted exposure when the storage provider is compromised. Our work can alert users before uploading their private (or sensitive) images to the cloud systems to control the amount of personal information (eg. social security number) shared through images.

In the future, using this study, an architecture can be developed, that will incorporate other contextual information about images such as personal information about the image owner, owner's privacy preferences or the owner social network activities, in addition to the visual content of the image. Another interesting direction is to extend these CNN architectures to describe and localize the sensitive content in private images.

## ACKNOWLEDGMENTS

## REFERENCES

Fabeah Adu-Oppong, Casey K. Gardiner, Apu Kapadia, and Patrick P. Tsang. 2008. Social Circles: Tackling Privacy in Social Networks. In *Symposium on Usable Privacy andSecurity (SOUPS)*.

Shane Ahern, Dean Eckles, Nathaniel S. Good, Simon King, Mor Naaman, and Rahul Nair. 2007. Over-exposed?: Privacy Patterns and Considerations in Online and Mobile Photo Sharing. In *Proceedings of the SIGCHI Conference (CHI '07)*. ACM, New York, NY, USA, 357–366.

Emilia Apostolova and Dina Demner-Fushman. 2009. Towards Automatic Image Region Annotation - Image Region Textual Coreference Resolution. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers*. Association for Computational Linguistics, 41–44.

Arjun Baokar. 2016. A Contextually-Aware, Privacy-Preserving Android Permission Model. In *Technical Report No. UCB/EECS-2016-69, University of California, Berkeley*.

Yoshua Bengio. 2012. Deep Learning of Representations for Unsupervised and Transfer Learning.. In *ICML Unsupervised and Transfer Learning (JMLR Proceedings)*, Isabelle Guyon, Gideon Dror, Vincent Lemaire, Graham W. Taylor, and Daniel L. Silver

(Eds.), Vol. 27. JMLR.org, 17–36. http://dblp.uni-trier.de/db/journals/jmlr/jmlrp27.html#Bengio12

Andrew Besmer and Heather Lipford. 2009. Tagged photos: concerns, perceptions, and protections. In *CHI '09: 27th international conference extended abstracts on Human factors in computing systems*. ACM, New York, NY, USA, 4585–4590.

Igor Bilogrevic, Kévin Huguenin, Berker Agir, Murtuza Jadliwala, Maria Gazaki, and Jean-Pierre Hubaux. 2016. A machine-learning based approach to privacy-aware information-sharing in mobile social networks. *Pervasive and Mobile Computing* 25 (2016), 125–142.

Kerstin Bischoff, Claudiu S. Firan, Wolfgang Nejdl, and Raluca Paiu. 2008. Can All Tags Be Used for Search?. In *Proceedings of the 17th ACM Conference on Information and Knowledge Management (CIKM '08)*. ACM, New York, NY, USA, 193–202.

Joseph Bonneau, Jonathan Anderson, and Luke Church. 2009a. Privacy Suites: Shared Privacy for Social Networks. In *Proceedings of the 5th Symposium on Usable Privacy and Security (SOUPS '09)*. ACM, New York, NY, USA, Article 30, 1 pages.

Joseph Bonneau, Jonathan Anderson, and George Danezis. 2009b. Prying Data out of a Social Network. In *Proceedings of the 2009 International Conference on Advances in Social Network Analysis and Mining (ASONAM '09)*. IEEE Computer Society, Washington, DC, USA, 249–254.

Daniel Buschek, Moritz Bader, Emanuel von Zezschwitz, and Alexander De Luca. 2015. Automatic Privacy Classification of Personal Photos. In *Human Computer Interaction INTERACT 2015*, Julio Abascal, Simone Barbosa, Mirko Fetter, Tom Gross, Philippe Palanque, and Marco Winckler (Eds.). Vol. 9297. 428–435.

Jaeyoung Choi, Martha Larson, Xinchao Li, Kevin Li, Gerald Friedland, and Alan Hanjalic. 2017. The Geo-Privacy Bonus of Popular Photo Enhancements. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval (ICMR '17)*. ACM, New York, NY, USA, 84–92. https://doi.org/10.1145/3078971.3080543

Delphine Christin, Pablo SáNchez LóPez, Andreas Reinhardt, Matthias Hollick, and Michaela Kauer. 2013. Share with Strangers: Privacy Bubbles As User-centered Privacy Control for Mobile Content Sharing Applications. *Inf. Secur. Tech. Rep.* 17, 3 (Feb. 2013), 105–116.

Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural Language Processing (Almost) from Scratch. *J. Mach. Learn. Res.* 12 (Nov. 2011), 2493–2537.

George Danezis. 2009. Inferring Privacy Policies for Social Networking Services. In *Proceedings of the 2nd ACM Workshop on Security and Artificial Intelligence (AISec '09)*. ACM, New York, NY, USA, 5–10. https://doi.org/10.1145/1654988.1654991

Munmun De Choudhury, Hari Sundaram, Yu-Ru Lin, Ajita John, and Doree Duncan Seligmann. 2009. Connecting Content to Community in Social Media via Image Content, User Tags and User Communication. In *Proceedings of the 2009 IEEE International Conference on Multimedia and Expo (ICME'09)*. IEEE Press, Piscataway, NJ, USA, 1238–1241.

Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. 2013. DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. *CoRR* abs/1310.1531 (2013).

F. Dufaux and T. Ebrahimi. 2008. Scrambling for Privacy Protection in Video Surveillance Systems. *IEEE Trans. Cir. and Sys. for Video Technol.* 18, 8 (Aug. 2008), 1168–1174. https://doi.org/10.1109/TCSVT.2008.928225

Lujun Fang and Kristen LeFevre. 2010. Privacy Wizards for Social Networking Sites. In *Proceedings of the 19th International Conference on World Wide Web (WWW '10)*. ACM, New York, NY, USA, 351–360.

Clement Farabet, Camille Couprie, Laurent Najman, and Yann LeCun. 2013. Learning Hierarchical Features for Scene Labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (August 2013).

Drew Fisher, Leah Dorner, and David Wagner. 2012. Short Paper: Location Privacy: User Behavior in the Field. In *Proceedings of the Second ACM Workshop on Security and Privacy in Smartphones and Mobile Devices (SPSM '12)*. ACM, New York, NY, USA, 51–56. https://doi.org/10.1145/2381934.2381945

Julien Freudiger, Reza Shokri, and Jean-Pierre Hubaux. 2012. Evaluating the Privacy Risk of Location-Based Services. In *Financial Cryptography and Data Security*, George Danezis (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 31–46.

Gerald Friedland and Robin Sommer. 2010. Cybercasing the Joint: On the Privacy Implications of Geo-Tagging. In *HotSec*. USENIX Association.

Yue Gao, Meng Wang, Huanbo Luan, Jialie Shen, Shuicheng Yan, and Dacheng Tao. 2011. Tag-based Social Image Search with Visual-text Joint Hypergraph Learning. In *Proceedings of the 19th ACM International Conference on Multimedia (MM '11)*. ACM, New York, NY, USA, 1517–1520. https://doi.org/10.1145/2072298.2072054

Kambiz Ghazinour, Stan Matwin, and Marina Sokolova. 2013. Monitoring and Recommending Privacy Settings in Social Networks. In *Proceedings of the Joint EDBT/ICDT 2013 Workshops (EDBT '13)*. ACM, New York, NY, USA, 164–168.

Ralph Gross and Alessandro Acquisti. 2005. Information Revelation and Privacy in Online Social Networks. In *Proceedings of the 2005 ACM Workshop on Privacy in the Electronic Society*. 71–80.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016a. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. 770–778.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016b. Identity Mappings in Deep Residual Networks. In *ECCV*.

Benjamin Henne, Christian Szongott, and Matthew Smith. 2013. SnapMe if You Can: Privacy Threats of Other Peoples' Geo-tagged Media and What We Can Do About It (WiSec '13). 12. https://doi.org/10.1145/2462096.2462113

Livia Hollenstein and Ross Purves. 2010. Exploring place through user-generated content: Using Flickr tags to describe city cores. J. Spatial Information Science 1, 1 (2010), 21–48. https://doi.org/10.5311/JOSIS.2010.1.3

D. Hu, F. Chen, X. Wu, and Z. Zhao. 2016. A Framework of Privacy Decision Recommendation for Image Sharing in Online Social Networks. In 2016 IEEE First International Conference on Data Science in Cyberspace (DSC). 243–251.

Panagiotis Ilia, Iasonas Polakis, Elias Athanasopoulos, Federico Maggi, and Sotiris Ioannidis. 2015. Face/Off: Preventing Privacy Leakage From Photos in Social Networks. In Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security (CCS '15). ACM, New York, NY, USA, 781–792. https://doi.org/10.1145/2810103.2813603

Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. 2014. Caffe: Convolutional Architecture for Fast Feature Embedding. In Proceedings of the ACM International Conference on Multimedia. 675–678.

Simon Jones and Eamonn O'Neill. 2011. Contextual dynamics of group-based sharing decisions. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '11). ACM, 1777–1786.

James B. D. Joshi and Tao Zhang (Eds.). 2009. The 5th International Conference on Collaborative Computing: Networking, Applications and Worksharing, CollaborateCom 2009, Washington DC, USA, November 11-14, 2009. ICST / IEEE.

Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A Convolutional Neural Network for Modelling Sentences. CoRR abs/1404.2188 (2014).

Berkant Kepez and Pinar Yolum. 2016. Learning Privacy Rules Cooperatively in Online Social Networks. In Proceedings of the 1st International Workshop on AI for Privacy and Security (PrAISe '16). ACM, New York, NY, USA, Article 3, 4 pages.

Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In EMNLP. ACL, 1746–1751.

Peter F. Klemperer, Yuan Liang, Michelle L. Mazurek, Manya Sleeper, Blase Ur, Lujo Bauer, Lorrie Faith Cranor, Nitin Gupta, and Michael K. Reiter. 2012. Tag, you can see it! Using tags for access control in photo sharing. In CHI 2012: Conference on Human Factors in Computing Systems. ACM.

Balachander Krishnamurthy and Craig E. Wills. 2008. Characterizing Privacy in Online Social Networks. In Proceedings of the First Workshop on Online Social Networks (WOSN '08). ACM, New York, NY, USA, 37–42.

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In Advances in Neural Information Processing Systems 25, F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger (Eds.). Curran Associates, Inc., 1097–1105.

Zhenzhong Kuang, Zongmin Li, Dan Lin, and Jianping Fan. 2017. Automatic Privacy Prediction to Accelerate Social Image Sharing. In Third IEEE International Conference on Multimedia Big Data, BigMM 2017, Laguna Hills, CA, USA, April 19-21, 2017. 197–200. https://doi.org/10.1109/BigMM.2017.70

Abdurrahman Can Kurtan and Pinar Yolum. 2018. PELTE: Privacy Estimation of Images from Tags. In Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '18). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1989–1991. http://dl.acm.org/citation.cfm?id=3237383.3238047

Benjamin Laxton, Kai Wang, and Stefan Savage. 2008. Reconsidering Physical Key Secrecy: Teleduplication via Optical Decoding. In Proceedings of the 15th ACM Conference on Computer and Communications Security (CCS '08). ACM, 469–478.

Quoc V. Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. In ICML (JMLR Workshop and Conference Proceedings), Vol. 32. JMLR.org, 1188–1196.

Yann LeCun. 2017. Facebook Envisions AI That Keeps You From Uploading Embarrassing Pics. https://www.wired.com/2014/12/fb/all/1. (2017). [Online; accessed 12-April-2017].

Yifang Li, Wyatt Troutman, Bart P. Knijnenburg, and Kelly Caine. 2018. Human Perceptions of Sensitive Content in Photos. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops.

Heather Richter Lipford, Andrew Besmer, and Jason Watson. 2008. Understanding Privacy Settings in Facebook with an Audience View. In Proceedings of the 1st Conference on Usability, Psychology, and Security (UPSEC'08). 2:1–2:8.

David G. Lowe. 2004. Distinctive Image Features from Scale-Invariant Keypoints. IJCV 60, 2 (Nov. 2004), 91–110.

Mary Madden. 2012. Privacy management on social media sites. http://www.pewinternet.org/2012/02/24/privacy-management-on-social-media-sites. (2012). [Online; accessed 12-November-2017].

Mohammad Mannan and Paul C. van Oorschot. 2008. Privacy-enhanced Sharing of Personal Content on the Web. In Proceedings of the 17th International Conference on World Wide Web. 487–496.

Erika McCallister, Timothy Grance, and Karen A. Scarfone. 2010. SP 800-122. Guide to Protecting the Confidentiality of Personally Identifiable Information (PII). Technical Report. Gaithersburg, MD, United States.

Yuta Nakashima, Noboru Babaguchi, and Jianping Fan. 2011. Automatic generation of privacy-protected videos using background estimation. In Proceedings of the 2011

IEEE International Conference on Multimedia and Expo, ICME 2011, 11-15 July, 2011, Barcelona, Catalonia, Spain. 1–6. https://doi.org/10.1109/ICME.2011.6011955

Yuta Nakashima, Noboru Babaguchi, and Jianping Fan. 2012. Intended human object detection for automatically protecting privacy in mobile video surveillance. Multimedia Syst. 18, 2 (2012), 157–173. https://doi.org/10.1007/s00530-011-0244-y

Yuta Nakashima, Noboru Babaguchi, and Jianping Fan. 2016. Privacy Protection for Social Video via Background Estimation and CRF-Based Videographer's Intention Modeling. IEICE Transactions 99-D, 4 (2016), 1221–1233.

Katarzyna Olejnik, Italo Dacosta, Joana Soares Machado, Kévin Huguenin, Mohammad Emtiyaz Khan, and Jean-Pierre Hubaux. 2017. SmarPer: Context-Aware and Automatic Runtime-Permissions for Mobile Devices. In 38th IEEE Symposium on Security and Privacy (S&P). IEEE, San Jose, CA, United States, 1058–1076. https://doi.org/10.1109/SP.2017.25

Aude Oliva and Antonio Torralba. 2001. Modeling the Shape of the Scene: A Holistic Representation of the Spatial Envelope. IJCV 42, 3 (May 2001), 145–175.

Tribhuvanesh Orekondy, Mario Fritz, and Bernt Schiele. 2018. Connecting Pixels to Privacy and Utility: Automatic Redaction of Private Information in Images. In Conference on Computer Vision and Pattern Recognition (CVPR).

Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. 2017. Towards a Visual Privacy Advisor: Understanding and Predicting Privacy Risks in Images. In IEEE International Conference on Computer Vision, ICCV 2017. 3706–3715.

J. Parra-Arnau, A. Perego, E. Ferrari, J. Forné, and D. Rebollo-Monedero. 2014. Privacy-Preserving Enhanced Collaborative Tagging. IEEE Transactions on Knowledge and Data Engineering 26, 1 (Jan 2014), 180–193.

Javier Parra-Arnau, David Rebollo-Monedero, Jordi Forné, Jose L. Muñoz, and Oscar Esparza. 2012. Optimal tag suppression for privacy protection in the semantic Web. Data Knowl. Eng. 81-82 (2012), 46–66.

João Paulo Pesce, Diego Las Casas, Gustavo Rauber, and Virgílio Almeida. 2012. Privacy Attacks in Social Media Using Photo Tagging Networks: A Case Study with Facebook. In Proceedings of the 1st Workshop on Privacy and Security in Online Social Media (PSOSM '12). ACM, New York, NY, USA, Article 4, 8 pages.

Georgios Petkos, Symeon Papadopoulos, and Yiannis Kompatsiaris. 2015. Social Circle Discovery in Ego-Networks by Mining the Latent Structure of User Connections and Profile Attributes. In Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015 (ASONAM '15). ACM, New York, NY, USA, 880–887.

Moo-Ryong Ra, Ramesh Govindan, and Antonio Ortega. 2013. P3: Toward Privacy-preserving Photo Sharing. In Proceedings of the 10th USENIX Conference on Networked Systems Design and Implementation (nsdi'13). USENIX Association, Berkeley, CA, USA, 515–528. http://dl.acm.org/citation.cfm?id=2482626.2482675

Ramprasad Ravichandran, Michael Benisch, Patrick Gage Kelley, and Norman M. Sadeh. 2009. Capturing Social Networking Privacy Preferences:. Springer Berlin Heidelberg, Berlin, Heidelberg, 1–18.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. ImageNet Large Scale Visual Recognition Challenge. IJCV (April 2015), 1–42.

Pierre Sermanet, David Eigen, Xiang Zhang, Michael Mathieu, Rob Fergus, and Yann LeCun. 2014. OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks. In ICLR 2014. CBLS.

Pierre Sermanet, Koray Kavukcuoglu, Soumith Chintala, and Yann Lecun. 2013. Pedestrian Detection with Unsupervised Multi-stage Feature Learning. In Proceedings of the 2013 IEEE Conference on Computer Vision and Pattern Recognition. 3626–3633.

S. Shamma and M. Y. S. Uddin. 2014. Towards privacy-aware photo sharing using mobile phones. In 8th International Conference on Electrical and Computer Engineering. 449–452. https://doi.org/10.1109/ICECE.2014.7026919

Reza Shokri, George Theodorakopoulos, Jean-Yves Le Boudec, and Jean-Pierre Hubaux. 2011. Quantifying Location Privacy. In IEEE Symposium on Security and Privacy. IEEE Computer Society, 247–262.

Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR abs/1409.1556 (2014).

Andrew Simpson. 2008. On the Need for User-defined Fine-grained Access Control Policies for Social Networking Applications. In Proceedings of the Workshop on Security in Opportunistic and SOCial Networks (SOSOC '08). ACM, New York, NY, USA, Article 1, 8 pages. https://doi.org/10.1145/1461469.1461470

Eleftherios Spyromitros-Xioufis, Symeon Papadopoulos, Adrian Popescu, and Yiannis Kompatsiaris. 2016. Personalized Privacy-aware Image Classification. In Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval (ICMR '16). ACM, New York, NY, USA, 71–78.

Anna Squicciarini, Cornelia Caragea, and Rahul Balakavi. 2014. Analyzing Images' Privacy for the Modern Web. In Proceedings of the 25th ACM Conference on Hypertext and Social Media (HT '14). ACM, New York, NY, USA, 136–147.

Anna Squicciarini, Cornelia Caragea, and Rahul Balakavi. 2017a. Toward Automated Online Photo Privacy. ACM Trans. Web 11, 1, Article 2 (April 2017), 29 pages.

Anna Squicciarini, D. Lin, S. Karumanchi, and N. DeSisto. 2012. Automatic social group organization and privacy management. In 8th International Conference on Collaborative Computing: Networking, Applications and Worksharing (CollaborateCom).

89–96.

Anna Squicciarini, Dan Lin, Smitha Sundareswaran, and Joshua Wede. 2015. Privacy Policy Inference of User-Uploaded Images on Content Sharing Sites. *IEEE Trans. Knowl. Data Eng.* 27, 1 (2015), 193–206.

Anna Squicciarini, Andrea Novelli, Dan Lin, Cornelia Caragea, and Haoti Zhong. 2017b. From Tag to Protect: A Tag-Driven PolicyRecommender System for Image Sharing. In *PST '17*.

Anna Squicciarini, Mohamed Shehab, and Federica Paci. 2009. Collective Privacy Management in Social Networks. In *Proceedings of the 18th International Conference on World Wide Web (WWW '09)*. ACM, New York, NY, USA, 521–530.

H. Sundaram, L. Xie, M. De Choudhury, Y.R. Lin, and A. Natsev. 2012. Multimedia Semantics: Interactions Between Content and Community. *Proc. IEEE* 100, 9 (2012), 2737–2758.

Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2014. Going Deeper with Convolutions. *CoRR* abs/1409.4842 (2014).

Jinhui Tang, Shuicheng Yan, Richang Hong, Guo-Jun Qi, and Tat-Seng Chua. 2009. Inferring Semantic Concepts from Community-contributed Images and Noisy Tags. In *Proceedings of the 17th ACM International Conference on Multimedia (MM '09)*. ACM, New York, NY, USA, 223–232. https://doi.org/10.1145/1631272.1631305

Eran Toch. 2014. Crowdsourcing Privacy Preferences in Context-aware Applications. *Personal Ubiquitous Comput.* 18, 1 (Jan. 2014), 129–141. https://doi.org/10.1007/s00779-012-0632-0

Ashwini Tonge and Cornelia Caragea. 2015. Privacy Prediction of Images Shared on Social Media Sites Using Deep Features. *CoRR* abs/1510.08583 (2015).

Ashwini Tonge and Cornelia Caragea. 2016. Image Privacy Prediction Using Deep Features. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*. 4266–4267.

Ashwini Tonge and Cornelia Caragea. 2018. On the Use of "Deep" Features for Online Image Sharing. In *The Web Conference Companion*.

Ashwini Tonge, Cornelia Caragea, and Anna Squicciarini. 2018a. Privacy-Aware Tag Recommendation for Image Sharing. In *Proceedings of the 29th on Hypertext and Social Media (HT '18)*. ACM, New York, NY, USA, 52–56.

Ashwini Tonge, Cornelia Caragea, and Anna Squicciarini. 2018b. Uncovering Scene Context for Predicting Privacy of Online Shared Images. In *AAAI' 18*.

Lam Tran, Deguang Kong, Hongxia Jin, and Ji Liu. 2016. Privacy-CNH: A Framework to Detect Photo Privacy with Convolutional Neural Network Using Hierarchical Features. In *Proceedings of the Thirtieth AAAI Conference*. 1317–1323.

Paul Viola and Michael Jones. 2001. Robust Real-time Object Detection. In *IJCV*.

Emanuel von Zezschwitz, Sigrid Ebbinghaus, Heinrich Hussmann, and Alexander De Luca. 2016. You Can'T Watch This!: Privacy-Respectful Photo Browsing on Smartphones. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI '16)*. ACM, New York, NY, USA, 4320–4324. https://doi.org/10.1145/2858036.2858120

Nitya Vyas, Anna Squicciarini, Chih-Cheng Chang, and Danfeng Yao. 2009. Towards automatic privacy management in Web 2.0 with semantic analysis on annotations. In *CollaborateCom*. 1–10.

Susan Waters and James Ackerman. 2011. Exploring Privacy Management on Facebook: Motivations and Perceived Consequences of Voluntary Disclosure. *Journal of Computer-Mediated Communication* 17, 1 (2011), 101–115.

Jason Watson, Heather Richter Lipford, and Andrew Besmer. 2015. Mapping User Preference to Privacy Default Settings. *ACM Trans. Comput.-Hum. Interact.* 22, 6, Article 32 (Nov. 2015), 20 pages. https://doi.org/10.1145/2811257

Zhenyu Wu, Zhangyang Wang, Zhaowen Wang, and Hailin Jin. 2018. Towards Privacy-Preserving Visual Recognition via Adversarial Training: A Pilot Study. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XVI (Lecture Notes in Computer Science)*, Vittorio Ferrari,

Martial Hebert, Cristian Sminchisescu, and Yair Weiss (Eds.), Vol. 11220. Springer, 627–645. https://doi.org/10.1007/978-3-030-01270-0_37

Haitao Xu, Haining Wang, and Angelos Stavrou. 2015. Privacy Risk Assessment on Online Photos.. In *RAID*, Herbert Bos, Fabian Monrose, and Gregory Blanc (Eds.), Vol. 9404. 427–447.

C.M.A. Yeung, L. Kagal, N. Gibbins, and N. Shadbolt. 2009. Providing access control to online photo albums based on tags and linked data. *Social Semantic Web: Where Web* 2 (2009).

Wen-tau Yih, Xiaodong He, and Christopher Meek. 2014. Semantic Parsing for Single-Relation Question Answering. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Baltimore, Maryland, 643–648. http://www.aclweb.org/anthology/P14-2105

Jun Yu, Zhenzhong Kuang, Zhou Yu, Dan Lin, and Jianping Fan. 2017a. Privacy Setting Recommendation for Image Sharing. In *16th IEEE International Conference on Machine Learning and Applications, ICMLA 2017, Cancun, Mexico, December 18-21, 2017*. 726–730. https://doi.org/10.1109/ICMLA.2017.00-73

Jun Yu, Zhenzhong Kuang, Baopeng Zhang, Wei Zhang, Dan Lin, and Jianping Fan. 2018. Leveraging Content Sensitiveness and User Trustworthiness to Recommend Fine-Grained Privacy Settings for Social Image Sharing. *IEEE Trans. Information Forensics and Security* 13, 5 (2018), 1317–1332. https://doi.org/10.1109/TIFS.2017.2787986

Jun Yu, Baopeng Zhang, Zhengzhong Kuang, Dan Lin, and Jianping Fan. 2017b. iPrivacy: Image Privacy Protection by Identifying Sensitive Objects via Deep Multi-Task Learning. *IEEE Trans. Information Forensics and Security* 12, 5 (2017), 1005–1016.

Lin Yuan, Joel Regis Theytaz, and Touradj Ebrahimi. 2017. Context-Dependent Privacy-Aware Photo Sharing based on Machine Learning. *Proc. of 32nd International Conference on ICT Systems Security and Privacy Protection (IFIP SEC)* (2017).

X. Yuan, X. Wang, C. Wang, Anna Squicciarini, and K. Ren. 2018. Towards Privacy-Preserving and Practical Image-Centric Social Discovery. *IEEE Transactions on Dependable and Secure Computing* 15, 5 (Sept 2018), 868–882.

Sergej Zerr, Stefan Siersdorfer, Jonathon Hare, and Elena Demidova. 2012b. Privacy-aware image classification and search. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*. ACM, NY, USA.

Sergej Zerr, Stefan Siersdorfer, and Jonathon S. Hare. 2012a. PicAlert!: a system for privacy-aware image classification and retrieval.. In *CIKM*, Xue wen Chen, Guy Lebanon, Haixun Wang, and Mohammed J. Zaki (Eds.). ACM, 2710–2712.

Wei Zhang, S. S. Cheung, and Minghua Chen. 2005. Hiding privacy information in video surveillance system. In *IEEE International Conference on Image Processing 2005*, Vol. 3. II–868. https://doi.org/10.1109/ICIP.2005.1530530

Yuchen Zhao, Juan Ye, and Tristan Henderson. 2014. Privacy-aware Location Privacy Preference Recommendations. In *Proceedings of the 11th International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services (MOBIQUITOUS '14)*. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), ICST, Brussels, Belgium, Belgium, 120–129. https://doi.org/10.4108/icst.mobiquitous.2014.258017

Elena Zheleva and Lise Getoor. 2009. To Join or Not to Join: The Illusion of Privacy in Social Networks with Mixed Public and Private User Profiles. In *Proceedings of the 18th International Conference on World Wide Web (WWW '09)*. 531–540.

Haoti Zhong, Anna Squicciarini, and David Miller. 2018. Toward Automated Multiparty Privacy Conflict Detection. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM '18)*. ACM, New York, NY, USA, 1811–1814. https://doi.org/10.1145/3269206.3269329

Haoti Zhong, Anna Squicciarini, David Miller, and Cornelia Caragea. 2017. A Group-Based Personalized Model for Image Privacy Classification and Labeling. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*. 3952–3958.