

# Discovering Visual Concept Structure with Sparse and Incomplete Tags

Jingya Wang<sup>a</sup>, Xiatian Zhu<sup>a</sup>, Shaogang Gong<sup>a</sup>

<sup>a</sup>*School of EECS, Queen Mary University of London, Mile End Road, London, E1 4NS, UK*

---

## Abstract

Discovering automatically the semantic structure of tagged visual data (e.g. web videos and images) is important for visual data analysis and interpretation, enabling the machine intelligence for effectively processing the fast-growing amount of multi-media data. However, this is non-trivial due to the need for jointly learning underlying correlations between heterogeneous visual and tag data. The task is made more challenging by inherently sparse and incomplete tags. In this work, we develop a method for modelling the inherent visual data concept structures based on a novel Hierarchical-Multi-Label Random Forest model capable of correlating structured visual and tag information so as to more accurately interpret the visual semantics, e.g. disclosing meaningful visual groups with similar high-level concepts, and recovering missing tags for individual visual data samples. Specifically, our model exploits hierarchically structured tags of different semantic abstractness and multiple tag statistical correlations in addition to modelling visual and tag interactions. As a result, our model is able to discover more accurate semantic correlation between textual tags and visual features, and finally providing favourable visual semantics interpretation even with highly sparse and incomplete tags. We demonstrate the advantages of our proposed approach in two fundamental applications, visual data clustering and missing tag completion, on benchmarking video (i.e. TRECVID MED 2011) and image (i.e. NUS-WIDE) datasets.

**Keywords:** Visual semantic structure; Tag hierarchy; Tag correlation; Sparse tags; Incomplete tags; Data clustering; Missing tag completion; Random forest.

---

## 1. Introduction

A critical task in visual data analysis is to automatically discover and interpret the underlying semantic concept structure of large quantities of data effectively and quickly, which allows the computing intelligence for automated organisation and management of large scale multi-media data. However, semantic structure discovery for visual data by visual feature analysis alone is inherently limited due to the semantic gap between low-level visual features and high-level semantics, particularly under the “curse” of high dimensionality, where visual features are often represented in a high-dimensional feature space [1]. On the other hand, videos and images are often attached with additional non-visual data, e.g. typically some textual sketch (Figure 1(a)). Such text information can include short tags contributed by either users or content providers, for instance, videos/images from the YouTube and Flickr websites. Often, tags may provide uncontrolled mixed levels of information but being also incomplete with respect to the visual content. This motivates (1) *multi-modality based data cluster discovery* (where visual data samples in each hidden cluster/group share the same underlying high-level concept relevant to both visual appearance and textual tags in a latent unknown space) [2, 3, 4], and (2) *instance-level tag structure completion* (where the tag set is defined as the combination of all presented tags

and missing tag revelation for each visual data sample may rely on both visual appearance and given tags) [5, 6, 7]. The former considers global data group structure, e.g. data clustering (Figure 1(b)) that serves as a critical automated data analysis strategy with important fundamental applications, such as summarising video data for automatically removing redundancy and discovering meaningful / interesting content patterns hidden in large scale data corpus without any human labelling effort [8], detecting anomalies and salient data [2], or facilitating unstructured data browsing and examination [4]. In contrast, the latter addresses local tag label structure of individual visual instances, e.g. tag completion (Figure 1(c)) that aims to automatically recover missing concepts presented in visual data. In this multi-modality data learning context, it is necessary to highlight and distinguish three fundamental notions: (1) visual content, (2) visual features, and (3) textual tags. Among them, the latter two are different representations of the former, i.e. visual content – the actual target data/objects of our problem. By visual concept structure, we particularly refer to the concept structure of “visual content” rather than “visual features”.

Exploiting readily accessible textual tags in visual content interpretation has shown to be beneficial [3, 4, 6]. Nonetheless, existing methods are restricted in a number of ways: (1) Tags are assumed with similar abstractness (or flattened tag structure). Intrinsic hierarchical tag structures are ignored in model design; (2) Tag statistical correlations and interactions between visual and tag data are not fully exploited, partly due to model complexity and design limitation. Incorporating such information into existing models effectively is not straightforward.

---

*Email addresses:* jingya.wang@qmul.ac.uk (Jingya Wang),  
xiatian.zhu@qmul.ac.uk (Xiatian Zhu), s.gong@qmul.ac.uk  
(Shaogang Gong)

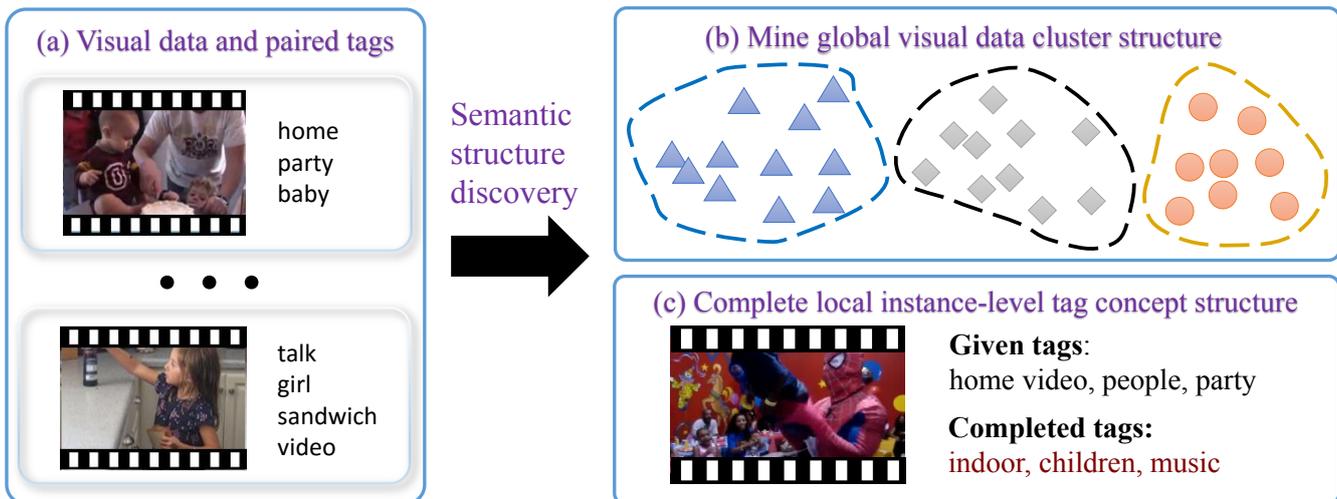


Figure 1: Problem illustration: we aim to develop an automated visual semantics discovery approach by exploiting (a) both visual and sparse tag data for (b) mining global visual data cluster structure and (c) completing local instance-level tag concept structure.

In general, joint learning of visual and text information, two different heterogeneous data modalities, in a shared representational space is non-trivial because: (1) The heteroscedasticity problem [9], that is, disparate data modalities significantly differ in representation (continuous or categorical) and distribution characteristics with different scales and covariances. In addition, the dimensionality of visual data often exceeds that of tag data by a large extent, like thousands vs. tens/hundreds. Because of this dimensionality discrepancy problem, a simple concatenation of heterogeneous feature spaces may result in an incoherent representation favourably inclined towards one dominant modality data and leading to suboptimal results. (2) Visual features can be inaccurate and unreliable, due to the inherently ambiguous and noisy visual data, and the imperfect nature of feature extraction. It is challenging to suppress the negative influence of unknown noisy visual features in data structure modelling. (3) The available text tags are often sparse and incomplete. This causes an inevitable problem that the visual (with much richer but also noisier and redundant information) and tag (being often sparse and incomplete although complementary) data are not always completely aligned and correlated.

In this work, we develop a model for robust visual semantic structure discovery and interpretation by employing both visual features and available sparse/incomplete text tags associated with the videos/images. The **contributions** of this work are as follows: **(I)** We formulate a novel approach capable of effectively extracting and fusing information from ambiguous/noisy visual features and sparse/incomplete textual tags for precisely discovering and mining the inherent visual semantic structures. This is made possible by introducing a new Hierarchical-Multi-Label Random Forest (HML-RF) model with a reformulated information gain function that allows to model the interactions between visual features and incomplete tags simultaneously. Specifically, our model is designed to minimise the uncertainty of tag distributions in an “abstract-to-specific” hierarchical fashion so as to exploit the high-order skeletal guidance

knowledge embedded in tag hierarchy structure. **(II)** We introduce a unified tag dependency based algorithm to cope with the tag sparseness and incompleteness problem. In particular, we formulate a principled way of locally integrating multiple statistical correlations (co-occurrence and mutual-exclusion) among tags during model optimisation. **(III)** We develop a data clustering method based on the proposed HML-RF model by measuring pairwise similarity between visual samples for accurately discovering the semantic global group structure of all visual data. **(IV)** We design three HML-RF tree structure driven tag prediction algorithms to recover missing tags for completing the local tag concept structure of individual visual data instances. We demonstrated the efficacy and superiority of our proposed approach on the TRECVID MED 2011 [4] (web videos) and NUS-WIDE [10] (web images) datasets through extensive comparisons with related state-of-the-art clustering, multi-view learning and tag completion methods.

## 2. Related Work

We review contemporary related studies on global structure analysis (e.g. data clustering) and local concept structure recovery (e.g. missing tag completion) using tagged visual data, tag correlation and hierarchy, and random forest models.

**Tagged visual data structure analysis:** Compared with low-level visual features, textual information provides high-level semantic meanings which can help bridge the gap between video features and human cognition. Textual tags have been widely employed along with visual features to help solve a variety of challenging computer vision problems, such as visual recognition [11] and retrieval [12], image annotation [13]. Rather than these supervised methods, we focus on structurally-constrained learning approach without the need of particular human labelling. Whilst a simple combination of visual features and textual tags may give rise to the difficult heteroscedasticity problem, Huang et al. [14] alternatively seek an optimal combina-

tion of similarity measures derived from different data modalities. The fused pairwise similarity can be then utilised for data clustering by existing graph based clustering algorithms such as spectral clustering [15]. As the interaction between visual appearance and textual tags is not modelled in the raw feature space but on the similarity graphs, the information loss in graph construction can not be recovered. Also, this model considers no inter-tag correlation.

Alternatively, multi-view learning/embedding methods are also able to jointly learn visual and text data by inferring a latent common subspace, such as multi-view metric learning [16], Restricted Boltzmann Machine and auto-encoders [17, 18], visual-semantic embedding [19], Canonical Correlation Analysis (CCA) and its variants [20, 21, 22, 23, 24]. Inspired by the huge success of deep neural networks, recently a few works have attempted to combine deep feature learning and CCA for advancing multi-view/modality data modelling [25, 26]. However, these methods usually assume a reasonably large number of tags available. Otherwise, the learned subspace may be subject to sub-optimal cross-modal correlation, e.g. in the case of significantly sparse tags. In addition, whilst incomplete tags can be considered as a special case of noisy labels, existing noise-tolerant methods [27, 28, 29] are not directly applicable. This is because they usually handle classification problems where a separate training dataset is required for model building, which however is not available in our context.

More recently, Zhou et al. [3] devised a Latent Maximum Margin Clustering (Latent MMC) model for assisting tagged video grouping. This model separates the whole task into two isolated stages: tag model learning and clustering, and thus their interaction is ignored. To tackle the above problem, Arash et al. [4] proposed a Structural MMC model where the correlations between visual features, tags and clusters are jointly modelled and optimised. The best results of clustering tagged videos are attained by Flip MMC [4] with the idea of flipping tags mainly for addressing the tag sparseness problem. In both MMC variants, tags are organised and used in a flat structure, whilst different tags may correspond to varying degrees of concept abstractness. Further, the statistical correlations between tags are neglected during optimisation. These factors may cause either degraded data modelling or knowledge loss, as shown in our experiments. Compared with these existing methods above, the proposed approach in this work is capable of jointly considering interactions between visual and tag data modalities, tag abstractness hierarchical structure and tag statistical correlations within a unified single model.

**Missing tag completion:** Text tags associated with videos and images are often sparse and incomplete, particularly those provided by web users. This may impose negative influence on tag-based applications and thus requires effective methods for tag completion. Different from conventional tag annotation [30, 31], tag completion does not require an extra completely annotated training dataset. Liu et al. [32] formulated tag completion as a non-negative data factorisation problem. Their method decomposes the global image representation into regional tag representations, on which the appearance of individ-

ual tags is characterised and visual-tag consistency is enforced. Wu et al. [5] performed tag recovery by searching for the optimal tag matrix which maximises the consistency with partially observed tags, visual similarity (e.g. visually similar samples are constrained to have common tags) and tag co-occurrence correlation. Lin et al. [7] developed a sparsity based tag matrix reconstruction method jointly considering visual-visual similarity, visual-tag association and tag-tag concurrence in completion optimisation. Similarly, Feng et al. [6] proposed another tag matrix recovery approach based on the low rank matrix theory [33]. Visual-tag consistency is also integrated into optimisation by exploring the graph Laplacian technique. However, all these methods ignore tag abstractness hierarchy structure, which may affect negatively the tag correlation and visual consistency modelling. Additionally, they depend on either global or regional visual similarity measures which can suffer from unknown noisy visual features or incomplete tags. Compared with these existing methods, we investigate an alternative strategy for tag completion, that is, to discover visual concept structure for identifying meaningful neighbourhoods and more accurate tag inference. To that end, we formulate a new Hierarchical-Multi-Label Random Forest (HML-RF) capable of jointly modelling tag and visual data, exploiting the intrinsic tag hierarchy knowledge, and the inherent strengths of a random forest for feature selection. We compare quantitatively our method with the state-of-the-art alternative tag completion models in extensive experiments and demonstrate the clear advantages of the proposed HML-RF model (Section 4.3).

**Tag hierarchy and correlations:** Hierarchy (a pyramid structure) is a natural knowledge organisation structure of our physical world, from more abstract to more specific in a top-down order [34, 35], and has been widely used in numerous studies, for example tag recommendation [36], semantic image segmentation [37], and object recognition [38]. Typically, an accurate hierarchy structure is assumed and utilised [37, 38]. But this is not always available, e.g. tag data extracted from some loosely structured meta-data source can only provide a rough hierarchy with potentially inaccurate relations, as the meta-data associated with videos in the TRECVID dataset. So are the user-provided tags from social media websites like Flickr. Such noisy hierarchy imposes more challenges but still useful if used properly. To that end, we exploit hierarchical tag structures in a more robust and coherent way for effective semantic structure modelling of sparsely tagged video/image data.

One of the most useful information encoded in hierarchy is inter-tag correlation, and *co-occurrence* should be most widely exploited, e.g. image annotation [39, 40], and object classification [38]. This positive label relation is useful since it provides a context for structuring the complexity of the real-world concepts/things. In contrast, *mutual-exclusion* is another (although less popular) relation between concepts. As opposite to co-occurrence, it is negative but complementary. Its application includes object detection [41, 42], multi-label image annotation [43], multi-task learning [44], and object recognition [38]. Unlike the above supervised settings, we investigate both correlations in a *structurally-constrained learning* manner. Also,

we do not assume their availability as in the case of [38]. Instead, we automatically mine these correlations from sparsely labelled data. Different from [43] where the tag structure is regarded as flat, we consider the co-occurrence and mutual-exclusive correlation between tags across layers of the tag hierarchy. We learn this pairwise relation, rather than assuming as prior knowledge as in [38]. Further, we relax the stringent assumption of accurate tags as made in [41, 42, 43] and the model is designed specifically to tolerate tag incompleteness and sparseness. Our goal is to exploit automatically the tag correlations and the available tag hierarchy structure effectively for inferring semantics on visual data and discovering visual concept structures.

**Random forest models:** Random forests have been shown to be effective for many computer vision tasks [45, 46, 47, 48]. Below we review several most related random forest variants. Montillo et al. [49] presented an Entangled Decision Forest for helping image segmentation by propagating knowledge across layers, e.g. dependencies between pixels and objects. Recently, Zhao et al. [50] proposed a multi-task forest for face analysis via learning different tasks at distinct layers according to the correlations between multi-tasks (e.g. head pose, facial landmarks). All these models are supervised. In contrast, our forest model performs structurally-constrained learning since we aim to discover and obtain semantic data structure using heterogeneous tags that are not target category labels but merely some semantic constraints. Furthermore, our model is unique in its capability of handling missing data, which is not considered in [50, 49]. The Constrained Clustering Forest (CC-Forest) [51, 52] is the most related to our HML-RF model, in that it is also utilised for data structure analysis e.g. measuring data affinity. The advantage of our model over CC-Forest are two-folds: (1) The capability for exploiting the tag hierarchical structure knowledge and (2) The superior effectiveness of tackling missing data, as shown in our experiments (Section 4).

### 3. Methodology

**Rational for model design:** We want to formulate a unified visual semantic structure discovery model capable of addressing the aforementioned challenges and limitations of existing methods. Specifically, to mitigate the heteroscedasticity and dimension discrepancy problems, we need to isolate different characteristics of visual and tag data, yet can still fully exploit the individual modalities as well as cross-modality interactions in a balanced manner. For handling tag sparseness and incompleteness, we propose to utilise the constraint information derived from inter-tag statistical correlations [39, 41, 38]. To that end, we wish to explore random forest [53, 54, 45] because of: (1) Its flexible training objective function for facilitating multi-modal data modelling and reformulation; (2) The decision tree’s hierarchical structures for flexible integration of abstract-to-specific structured tag topology; (3) Its inherent feature selection mechanism for handling inevitable data noise. Also, we need to resolve several shortcomings of the conventional clustering forest [54], as in its original form it is not best suited for solving our

problems in an unsupervised way. Specifically, clustering forest expects a fully concatenated representation as input during model training, it therefore does not allow a balanced utilisation of two modalities simultaneously (the dimension discrepancy problem), nor exploit interactions between visual and tag features. The existing classification forest is also not suitable as it is supervised and aims to learn a prediction function with class labelled training data (usually a single type of tag) [53]. Typical video/image tags do not offer class category labels. However, it is interesting to us that in contrast to the clustering forest, the classification forest offers a more balanced structure for using visual (as split variables) and tag (as semantic evaluation) data that is required for tackling the heteroscedasticity problem by isolating the two heterogeneous modalities during learning.

**Approach overview:** We want to reformulate the classification forest for automatically disclosing the semantic structure of videos or images with tags. To that end, we propose a novel *Hierarchical-Multi-Label Random Forest* (HML-RF). Our model goes beyond the classification forest in the following aspects: (1) Employing tags to constrain tree structure learning, rather than learning a generalised prediction function as [53, 45]; (2) Introducing a new objective function allowing *acceptance of multi-tags*, *exploitation of abstract-to-specific tag hierarchy* and *accommodation of multiple tag correlations* simultaneously. Instead of learning a classifier, HML-RF is designed to infer visual semantic concept structure for more accurately revealing both global visual data group structures and local tag structures of individual visual data samples. These structural relationships among data samples imply their underlying data group/cluster relations (obtained using a standard graph based clustering algorithm on the similarity graph estimated by our HML-RF model), as well as the specific tag concept structures of individual samples (predicted using the discovered semantic neighbourhoods encoded in the tree structures of HML-RF). An overview of the proposed visual concept structure discovery approach is depicted in Figure 3.

**Notations:** We consider two data modalities, (1) *Visual data modality* - We extract a  $d$ -dimensional visual descriptor from the  $i$ -th video/image sample denoted by  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d}) \in R^d, i = 1, \dots, n$ . All visual features are formed as  $X = \{\mathbf{x}_i\}_{i=1}^n$ . (2) *Tag data modality* - Tags associated with videos/images are extracted from the meta-data files or given by independent users. We represent  $m$  types of binary tag data ( $Z = \{1, \dots, m\}$ ) attached with the  $i$ -th video/image as  $\mathbf{y}_i = (y_{i,1}, \dots, y_{i,m}) \in [0, 1]^m$ . All tag data is defined as  $Y = \{\mathbf{y}_i\}_{i=1}^n$ . More details are provided in Section 4.1.

#### 3.1. Conventional Random Forests

Let us briefly introduce conventional random forests before detailing the proposed HML-RF model.

**Classification forests:** A classification forest [53] contains an ensemble of  $\tau$  binary decision trees. Growing a decision tree involves a recursive node splitting procedure until some stopping criterion is satisfied. The training of each split node is a

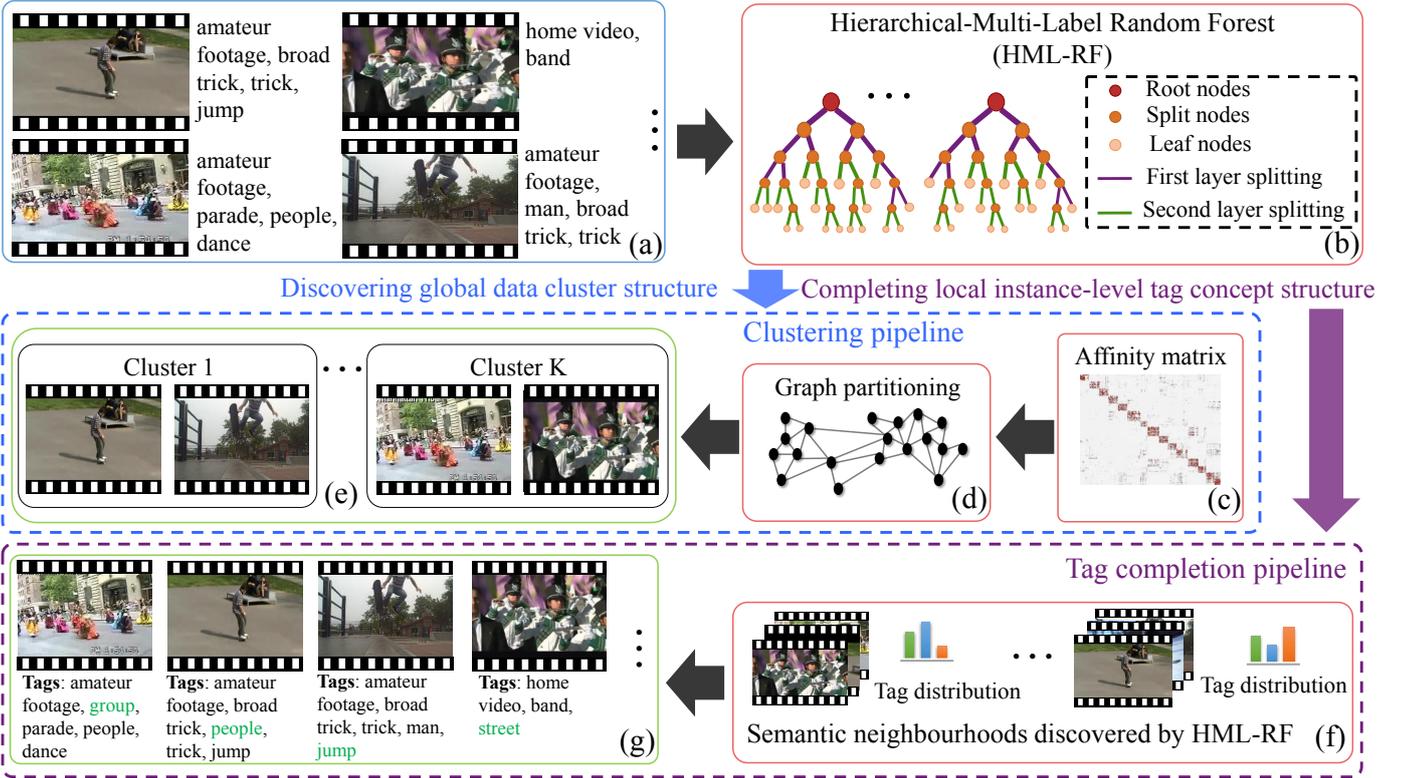


Figure 2: An overview of the proposed visual semantic structure discovery approach. (a) Example videos and associated tags; (b) The proposed HML-RF model designed to exploit inherent tag hierarchy for modelling correlations between ambiguous visual features and sparse tags, discover visual concept structures in two aspects: **Discovering global data cluster structure:** (c) Semantically constrained affinity matrix induced by HML-RF  $\rightarrow$  (d) Graph-based clustering to discover semantic groups  $\rightarrow$  (e) Resulting clusters with semantic similarity despite significant visual disparity. **Completing local instance-level tag concept structure:** (f) Semantic neighbourhood structures discovered by the proposed HML-RF model, which can then be exploited for (g) inferring missing tags to complete local concept structure at the data sample level.

process of binary split function optimisation, defined as

$$h(\mathbf{x}, \mathbf{w}) = \begin{cases} 0 & \text{if } x_f < \theta, \\ 1 & \text{otherwise.} \end{cases} \quad (1)$$

with two parameters  $\mathbf{w} = [f, \theta]$ : (i) a feature dimension  $x_f$  with  $f \in \{1, \dots, d\}$ , and (ii) a feature threshold  $\theta$ . The optimal split parameter  $\mathbf{w}^*$  is chosen via

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w} \in W} \Delta\psi_{sl}, \quad (2)$$

where the parameter search space  $W = \{\mathbf{w}_i\}_{i=1}^{\nu_{\text{try}}(|S|-1)}$  is formed by enumerating the threshold (or cut-point) on each of  $\nu_{\text{try}}$  randomly selected features (without replacement), with  $S$  denoting the sample set reaching the split node  $s$ . More specifically, the cut-points of each feature are defined as the unique midpoints of the intervals between ordered values from this feature on samples  $S$ . Thus, there is  $|S| - 1$  candidate cut-points for every chosen feature, with  $|\cdot|$  referring to the cardinality of a set. The information gain  $\Delta\psi_{sl}$  is formulated as

$$\Delta\psi_{sl} = \psi_s - \frac{|L|}{|S|}\psi_l - \frac{|R|}{|S|}\psi_r, \quad (3)$$

where  $L$  and  $R$  denote the data set routed into the left  $l$  and right  $r$  children, and  $L \cup R = S$ . The uncertainty  $\psi$  over the

label distribution can be computed as the Gini impurity [55] or entropy [45]. We used the former in our HML-RF model due to its simplicity and efficiency, i.e. the complexity of computing  $\psi_{sl}$  is  $O(1)$  as it is computed over the label distribution.

**Clustering forests:** Clustering forests aim to obtain an optimal data partitioning based on which pairwise similarity measures between samples can be inferred. In contrast to classification forests, clustering forests require no ground truth label information during the training phase. Similarly, a clustering forest consists of binary decision trees. The leaf nodes in each tree define a spatial partitioning of the training data. Interestingly, the training of a clustering forest can be performed using the classification forest optimisation approach by adopting the pseudo two-class algorithm [53, 54]. With this data augmentation strategy, the clustering problem becomes a canonical classification problem that can be solved by the classification forest training method as discussed above. The key idea behind this algorithm is to partition the augmented data space into dense and sparse regions [56]. One limitation of clustering forests is the limited ability in mining multiple modalities, as shown in Section 4.

### 3.2. Hierarchical-Multi-Label Random Forest

Our HML-RF can be considered as an extended hybrid model of classification and clustering forests. The model inputs in-

clude visual features  $x$  and tag data  $y$  of visual data samples (analogous to classification forest), and the output is semantic tree structures which can be used to predict an affinity matrix  $A$  over input samples  $X$  (similar to clustering forest). Conventional classification forests [53] typically assume single label type. In contrast, HML-RF can accept multiple types simultaneously as follows.

**Accommodating multiple tags:** A HML-RF model uses visual features as splitting variables to grow trees (HML-trees) as in Equation (1), but exploits all types of tag data *together* as tree structuring constraints in optimising  $w = [f, \theta]$ . Formally, we extend the conventional single-label based information gain function Equation (3) to multi-labels for training HML-trees:

$$\Delta\psi_{\text{ml}} = \sum_{i=1}^m \Delta\psi_{\text{sl}}^i \quad (4)$$

This summation merges all individual information gains  $\Delta\psi_{\text{sl}}^i$  from the  $i$ -th tag in an intuitive way for simultaneously enforcing knowledge of multiple tags into the HML-tree training process. Hence, the split functions are optimised in a similar way as supervised classification forests, and semantics from multiple tags are enforced simultaneously.

*Discussion:* In the context of structure discovery, e.g. tagged video/image clustering, it should be noted that our way of exploiting tags is different from conventional supervised classification forests since the tags are not target classes but semantic constraints. We call this “*structurally-constrained learning*”. Additionally, the interactions between visual features (on which split functions are defined) and tags (used to optimise split functions) are also modelled during learning by identifying the most discriminative visual features w.r.t. a collection of textual tags. Importantly, this separation of visual and tag data solves naturally the dimensionality discrepancy problem and addresses the heteroscedasticity challenge. Moreover, HML-RF benefits from the feature selection mechanism inherent to random forest for coping with noisy visual data by selecting the most discriminative localised split functions (Equation (1)) over multiple tags simultaneously.

**Incorporating tag hierarchy:** Equation (4) implies that all the tags have similar abstractness, as all of them are used in every split node (i.e. a flatten structure of tags). However, diverse tags may lie in multiple abstractness layers and how to exploit this information is critical for visual data structure modelling. The intuition is that tag hierarchy encodes approximately some relation knowledge between different underlying data structures and likely provides useful high-order skeletal guidance during the data structure inference process. The tag hierarchy structure can be roughly available from data source or automatically estimated by text analysis(see Section 4.1). To further exploit the abstractness guidance information in tag hierarchy, we introduce an adaptive hierarchical multi-label information gain function as:

$$\Delta\psi_{\text{hml}} = \sum_{k=1}^{\mu} \left( \prod_{j=1}^{k-1} (1 - \alpha_j) \alpha_k \sum_{i \in Z_k} \Delta\psi_{\text{sl}}^i \right) \quad (5)$$

where  $Z_k$  denotes the tag index set of the  $k$ -th layer in the tag hierarchy (totally  $\mu$  layers), with  $\cup_{k=1}^{\mu} Z_k = Z$ , and  $\forall_{j \neq k} Z_j \cap Z_k = \emptyset$ . Binary flag  $\alpha_k \in \{0, 1\}$  indicates the impurity of the  $k$ -th tag layer,  $k \in \{1, \dots, \mu\}$ , i.e.  $\alpha_k = 0$  when tag values are identical, i.e. pure, across all the training samples  $S$  of split node  $s$  in any tag  $i \in Z_k$ ,  $\alpha_k = 1$  otherwise. Note,  $\alpha$  is designed to be non-continuous so HML-tree per-node optimisation can focus on mining the underlying interactive information of visual-textual data at one specific semantic abstractness level. This shares a similar spirit to the “divide-and-conquer” learning strategy, e.g. reducing the local learning difficulty by considering *first* more homogeneous concepts only in training individual weak tree node models, *before* finally making the whole model to capture better semantic structure information. This is in contrast to solving the more difficult holistic optimisation problem on the entire tag set with a mixture of different abstractness levels. The target layer is  $k$  in case that  $\alpha_k = 1$  and  $\forall \alpha_j = 0, 0 < j < k$ .

*Discussion:* This layer-wise design allows the data partition optimisation to concentrate on the *most abstract* and *impure* tag layer (i.e. the target layer) so that the abstractness skeletal information in the tag hierarchy can be gradually embedded into the top-down HML-tree growing procedure for guiding the interaction modelling between visual and tag data in an abstract-to-specific fashion. This design and integration shall be natural and coherent because both tag hierarchy and HML-tree model are in the shape of pyramid and the divide-and-conquer modelling behaviour in HML-RF is intuitively suitable for the abstract-to-specific tag structure. We will show the empirical effectiveness of this layer-wise information gain design in our experiments (Section 4.2.3).

**Handling tag sparseness and incompleteness:** We further improve the HML-RF model by employing tag statistical correlations for addressing tag sparseness problem, as follows: We wish to utilise the dependences among tags to infer missing tags with a confidence measure (continuous soft tags), and exploit them along with labelled (binary hard) tags in localised split node optimisation, e.g. Equations (3) and (5).

In particular, two tag correlations are considered: *co-occurrence* - often co-occur in the same video/image samples thus positively correlated, and *mutual-exclusion* - rarely simultaneously appear so negatively correlated. They are complementary to each other, since for a particular sample, co-occurrence helps predict the *presence* degree of some missing tag based on another frequently co-occurrent tag who is labelled, whilst mutual-exclusion can estimate the *absence* degree of a tag according to its negative relation with another labelled tag. Therefore, we infer tag positive  $\{\hat{y}_{\cdot,i}^+\}$  and negative  $\{\hat{y}_{\cdot,i}^-\}$  confidence scores based upon tag co-occurrent and mutual-exclusive correlations, respectively. Note that  $\{\hat{y}_{\cdot,i}^+\}$  and  $\{\hat{y}_{\cdot,i}^-\}$  are not necessarily binary but more likely real number, e.g.  $[0, 1]$ . In our layered optimisation, we restrict the notion of missing tag to samples  $S_{\text{miss}} = \{\hat{x}\}$  where no tag in the target layer is labelled, and consider cross-layer tag correlations considering that a hierarchy is typically shaped as a pyramid, with more specific tag categories at lower layers where likely

more labelled tags are available. Suppose we compute the correlations between the tag  $i \in Z_k$  (the target tag layer) and the tag  $j \in \{Z_{k+1}, \dots, Z_\mu\}$  (subordinate tag layers).

*Co-occurrence:* We compute the co-occurrence  $\varrho_{i,j}$  as

$$\varrho_{i,j} = co_{i,j}/o_j, \quad (6)$$

where  $co_{i,j}$  denotes the co-occurrence frequency of tags  $i$  and  $j$ , that is, occurrences when both tags simultaneously appear in the same video/image across all samples; and  $o_j$  denotes the number of occurrences of tag  $j$  over all samples. Note that these statistics are collected from the available tags. The denominator  $o_j$  here is used to down-weight over-popular tags  $j$ : Those often appear across the dataset, and their existence thus gives a weak positive cue of supporting the simultaneous presence of tag  $i$ . For example, tag ‘people’ may appear in most videos and so brings a limited positive correlation to others. In spirit, this design shares the principle of Term Frequency Inverse Document Frequency [57, 58], which considers the inverse influence of total term occurrence times across the entire dataset as well. Once  $\varrho_{i,j}$  is obtained, for a potentially missing tag  $i \in Z_k$  of  $\hat{x} \in S_{\text{miss}}$ , we estimate its positive score  $\hat{y}_{\cdot,i}^+$  via:

$$\hat{y}_{\cdot,i}^+ = \sum_{j \in \{Z_{k+1}, \dots, Z_\mu\}} \varrho_{i,j} y_{\cdot,j} \quad (7)$$

where  $y_{\cdot,j}$  refers to the  $j$ -th tag value of  $\hat{x}$ . With Equation (7), we accumulate the positive support from all labelled subordinate tags to estimate the presence confidence of tag  $i$ .

---

#### Algorithm 1: Split function optimisation in a HML-tree

---

**Input:** At a split node  $s$  of a HML-tree  $t$ :

- Visual data  $X_s$  of training samples  $S$  arriving at  $s$ ;
- Corresponding labelled tag data  $Y_s$ ;
- Soft tag estimation using tag correlations:
  - \* Positive scores  $\{\hat{y}_{\cdot,i}^+\}$  estimated with Equations (6) and (7);
  - \* Negative scores  $\{\hat{y}_{\cdot,i}^-\}$  estimated with Equations (8) and (9);

**Output:**

- The best feature cut-point  $w^*$ ;
- The associated child node partition  $\{L^*, R^*\}$ ;

**1 Optimisation:**

- 2 Initialise  $L^* = R^* = \emptyset$ ,  $\Delta\psi_{\text{hml}}^* = 0$ ,  $w^* = [-1, -\infty]$ ;
- 3 **for**  $k \leftarrow 1$  **to**  $\nu_{\text{try}}$  **do**
- 4     Select a visual feature  $x_k \in \{1, \dots, d\}$  randomly without replacement;
- 5     **for** each possible cut-point of  $x_k$  **do**
- 6         Split  $S$  into a candidate partition  $\{L, R\}$ ;
- 7         Compute  $\Delta\psi_{\text{hml}}$  with Equations (3) and (5);
- 8         **if**  $\Delta\psi_{\text{hml}} > \Delta\psi_{\text{hml}}^*$  **then**
- 9             Update  $w^*$  with  $x_k$  and current threshold;
- 10            Update  $\Delta\psi_{\text{hml}}^* = \Delta\psi_{\text{hml}}$ ,  $L^* = L$ , and  $R^* = R$ .
- 11         **end**
- 12     **end**
- 13 **end**

---

*Mutual-exclusion:* We calculate this negative correlation as

$$\epsilon_{i,j} = \max(0, r_{i,j}^- - r_i^-)/(1 - r_i^-), \quad (8)$$

where  $r_i^-$  refers to the negative sample percentage on tag  $i$  across all samples, and  $r_{i,j}^-$  the negative sample percentage

on tag  $i$  over samples with positive tag  $j$ . The denominator  $(1 - r_i^-)$  is the normalisation factor. Hence,  $\epsilon_{i,j}$  measures statistically the relative increase in negative sample percentage on tag  $i$  given positive tag  $j$ . This definition reflects statistical exclusive degree of tag  $j$  against tag  $i$  intuitively. The cases of  $\epsilon < 0$  are not considered since they are already measured in the co-occurrence. Similarly, we predict the negative score  $\hat{y}_{\cdot,i}^-$  for  $\hat{x}$  on tag  $i$  with:

$$\hat{y}_{\cdot,i}^- = \sum_{j \in \{Z_{k+1}, \dots, Z_\mu\}} \epsilon_{i,j} y_{\cdot,j}. \quad (9)$$

Finally, we normalise both  $\hat{y}_{\cdot,i}^+$  and  $\hat{y}_{\cdot,i}^-$ ,  $i \in Z_p$ , into the unit range  $[0, 1]$ . Algorithm 1 summarises the split function optimisation procedure in a HML-tree.

### 3.3. Discovering Global Data Cluster Structure

Our HML-RF model is designed to discover visual semantic structures, e.g. global group structure over data samples. Inspired by clustering forests [53, 54, 45], this can be achieved by first estimating pairwise proximity between samples and then applying graph based clustering methods to obtain data groups (Figure 3(c,d,e)).

#### Inducing affinity graph from the trained HML-RF model:

Specifically, the  $t$ -th ( $t \in \{1, \dots, \tau\}$ ) tree within the HML-RF model partitions the training samples at its leaves. Each leaf node forms a *neighbourhood*, which contains a subset of data samples that share visual and semantic commonalities. All samples in a neighbourhood are *neighbours* to each other. These neighbours are considered similar both visually and semantically due to the proposed split function design (Equation (5)). More importantly, tag correlations and tag hierarchy structure knowledge are also taken into account in quantifying the semantic concept relationships. With these neighbourhoods, we consider an affinity model without any parameter to tune. Specifically, we assign pairwise similarity ‘1’ for sample pair  $(x_i, x_j)$  if they fall into the same HML-tree leaf node (i.e. being neighbours), and ‘0’ otherwise. This results in a tree-level affinity matrix  $A^t$ . A smooth affinity matrix  $A$  can be obtained through averaging all the tree-level affinity matrices:

$$A = \frac{1}{\tau} \sum_{t=1}^{\tau} A^t \quad (10)$$

with  $\tau$  the tree number of HML-RF. Equation (10) is adopted as the ensemble model of HML-RF due to its advantage of suppressing the noisy tree predictions, although other alternatives as the product of tree-level predictions are possible [45]. Intuitively, the multi-modality learning strategies of HML-RF enable its data similarity measure to be more meaningful. This can benefit significantly video/image clustering using a graph-based clustering method, as described next. **Forming global clusters:** Once the affinity matrix  $A$  is obtained, one can apply any off-the-shelf graph-based clustering model to acquire the final clustering result, e.g. spectral clustering [15]. Specifically, we firstly construct a sparse  $\kappa$ -NN graph, (Figure 3(d)), whose

edge weights are defined by  $\mathbf{A}$  (Figure 3(c)). Subsequently, we symmetrically normalise  $\mathbf{A}$  to obtain  $\mathbf{S} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ , where  $\mathbf{D}$  denotes a diagonal degree matrix with elements  $D_{i,i} = \sum_{j=1}^n \mathbf{A}_{i,j}$  ( $n$  denotes the video/image sample number). Given  $\mathbf{S}$ , we perform spectral clustering to discover the latent clusters of videos/images (Figure 3(e)). Each sample  $\mathbf{x}_i$  is then assigned to a cluster index  $c_i \in C$ , where  $C = \{1, \dots, p\}$  contains a total of  $p$  cluster indices.

### 3.4. Completing Local Instance-Level Concept Structure

In addition to inferring the global group structure, the learned semantic structure by the HML-RF model can also be exploited for reasoning the local concept structures of individual samples which are often partial and incomplete due to sparsely labelled tags. This task is known as *tag completion* [5]. Intuitively, the potential benefit of HML-RF for tag completion is due to semantic neighbourhoods over data samples formed during the model training phase (Section 3.2). More specifically, as data splits in HML-RF consider both correlations between visual features and tags, and dependencies between tags in abstractness hierarchy and statistics, visually similar neighbour samples (e.g. sharing the same leaves) may enjoy common semantic context and/or tags, and thus helpful and indicative in recovering missing tags. Formally, we aim to predict the existence probability  $p(\mathbf{x}_*, j)$  of a missing tag  $j \in Z$  in a sample  $\mathbf{x}_*$ . Given estimated  $p(\mathbf{x}_*, j)$ , those with top probabilities are considered as missing tags. To that end, we derive three tree-structure driven missing tag completion algorithms as below.

**(I) Completion by local neighbourhoods:** We estimate  $p(\mathbf{x}_*, j)$  by local neighbourhoods formed in HML-RF. Specifically, we first identify the neighbourhood  $N^t$  of  $\mathbf{x}_*$  in each HML-tree  $t \in \{1, 2, \dots, \tau\}$  by retrieving the leaf node that  $\mathbf{x}_*$  falls into. Second, for each  $N_{\mathbf{x}_*}^t$ , we compute the distribution pdf( $t, j$ ) of tag  $j$  over  $\mathbf{x}_*$ 's neighbours. As these neighbours are similar to  $\mathbf{x}_*$ , we use pdf( $t, j$ ) as a tree-level prediction. However, some neighbourhoods are unreliable due to the inherent visual ambiguity and tag sparseness, we thus ignore them and consider only confident ones with pdf( $t, j$ ) = 0 (called negative neighbourhood) or pdf( $t, j$ ) = 1 (called positive neighbourhood). Finally, we calculate  $p(\mathbf{x}_*, j)$  as

$$p(\mathbf{x}_*, j) = \frac{|P_j^+|}{|P_j^+| + |P_j^-|} \quad (11)$$

where  $|P_j^+|$  and  $|P_j^-|$  are the sets of positive and negative neighbourhoods, respectively. As such, the negative impact of unreliable neighbourhoods can be well suppressed. We denote this Local Neighbourhoods based method as “**HML-RF(LN)**”.

**(II) Completion by global structure:** Similar to local neighbourhoods of HML-RF, the data clusters (obtained with the method as described in Section 3.3) can be considered as global neighbourhoods. Therefore, we may alternatively exploit them for missing tag prediction. In particular, we assume that  $\mathbf{x}_*$  is assigned with cluster  $c$ . We utilise the cluster-level data distribution for missing tag estimation as:

$$p(\mathbf{x}_*, j) = \frac{|X_c^+|}{|X_c^+| - 1} \quad (12)$$

where  $X_c$  are data samples in cluster  $c$ , and  $X_c^+ \subset X_c$  are samples with labelled positive tag  $j$ . The intuition is that visual samples from the same cluster (thus of same high-level semantics/concept) are likely to share similar tags. Note, this is also a tree-structure based inference method in that these clusters are induced from tree-structure driven similarity measures (Section 3.3). We denote this Global Cluster based prediction algorithm as “**HML-RF(GC)**”.

**(III) Completion by affinity measure:** Similar to k-nearest neighbour classification [59, 60], we perform tag completion using affinity measures. Specifically, we utilise the tag information of  $\kappa$  nearest neighbours  $N_\kappa$  by adaptive weighting:

$$p(\mathbf{x}_*, j) = \frac{1}{|\kappa|} \sum_{i \in N_\kappa} y_{i,j} \mathbf{A}_{i,*} \quad (13)$$

where  $y_{i,j}$  denotes the tag  $j$  value of the  $i$ -th nearest neighbour  $\mathbf{x}_i$ ,  $\mathbf{A}_{i,*}$  is the pairwise similarity between  $\mathbf{x}_i$  and  $\mathbf{x}_*$  estimated by Equation (10), or the weight. Different from HML-RF(LN) that models the individual neighbourhoods within tree leaves, this method considers weighted pairwise relationship across all HML-trees, i.e. how many times two samples fall into the same leaf nodes. Conceptually, this can be considered as a hybrid model of HML-RF(LN) and HML-RF(GC) due to the inherent relation with both local neighbourhoods (i.e. tree leaves) and global clusters (the same similarity estimation). We denote this HML-RF Affinity Measure based tag recovery algorithm as “**HML-RF(AM)**”.

## 4. Experiments

### 4.1. Datasets and Experimental Settings

**Datasets:** We utilised two web-data benchmarks, the TRECVID MED 2011 video dataset [61] and the NUS-WIDE image dataset [10], for evaluating the performance of our proposed HML-RF model. Figure 3 shows a number of samples from the two datasets.

*TRECVID MED 2011:* It contains 2379 web videos from 15 clusters which we aim to discover in global structure analysis as in [3, 4]. This dataset is challenging for clustering using only visual features, in that videos with the same high-level concepts can present significant variety/dynamics in visual appearance. This necessitates the assistance of other data modalities, e.g. tags automatically extracted from textual judgement files associated with video samples [4]. Specifically, a total of 114 tags were obtained and used in our evaluation. On average, around 4 tags (3.5% of all tags) were extracted per video, thus very sparse and incomplete with the need for recovering many unknown missing tags. The tag hierarchy was established according to the structure presented in the meta-data files with two levels of tag abstractness. For example, tag “party” is more structurally abstract than tags “people/food/park” in the context of TRECVID videos where a number of semantic events (e.g. with respect to wedding ceremony and birthday celebration) may be meaningfully related with tag “party” whilst tags “people/food/park” should be very general and common to many

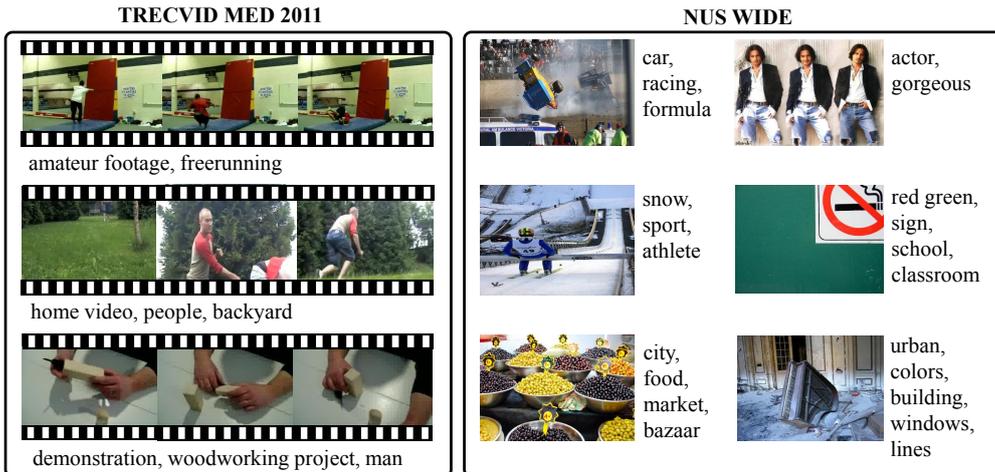


Figure 3: Examples from TRECVID MED 2011 [61] and NUS-WIDE [10].

different events and thus structurally specific. For video clustering, we aim to discover the underlying event category groups of web videos, given the ground-truth annotation available. This is similar to that of [62, 4]. For evaluating the performance of missing tag completion, we manually completed a subset of 200 video samples on 51 randomly selected tags as ground truth [6].

**NUS-WIDE:** We further evaluated the HML-RF model on a tagged web image dataset, NUS-WIDE [10]. We randomly sampled 30 clusters, each of which contains over 500 images and a total of 17523 images were selected for the evaluation of both global image clustering and local tag concept completion. This dataset contains 1000 different tags. Every image is labelled with 4.8 tags (i.e. 0.48% of all tags) on average.

For NUS-WIDE, we need to establish the tag hierarchy since tags are given in a flat structure. Inspired by [63, 22], we estimate the tag abstractness degree by mining and employing tag-image data statistics information. To be more precise, we first apply term frequency inverse document frequency (tf-idf) weighting to the binary tag vector  $\mathbf{y}_i = [y_{i,1}, \dots, y_{i,m}]$  of each image  $i$  ( $m$  denotes the tag type number), and get a new tag representation  $\tilde{\mathbf{y}}_i = [\tilde{y}_{i,1}, \dots, \tilde{y}_{i,m}]$ . This allows  $\tilde{\mathbf{y}}_i$  encoding the importance of each tag against the corresponding image by taking into account the tag-image statistic relation among the entire dataset. Then, we perform K-means over these tf-idf weighted tag vectors  $\{\tilde{\mathbf{y}}_i\}$  of all images to obtain  $E$  topic clusters. In each cluster  $e$  where  $\{\tilde{\mathbf{y}}_i\}$  fall into, we compute the abstractness or representativeness score for tag  $j$  as  $\sigma_j^e = \sum \tilde{y}_{i,j}^e$  and select the tags with top- $\eta$  highest  $\sigma_j^e$  scores into the current hierarchy layer. By performing this selection on all clusters, we form the current layer with selected most abstract tags whilst the remaining tags drop into lower layers. Similarly, we build one or multiple lower hierarchy layers on the remaining tags with the same steps above. Actually, we can consider this tag hierarchy formation as a process of revealing underlying topics in a layer-wise fashion. We select more tags per cluster for lower layers considering the potentially pyramid hierarchy shape, e.g. choosing top  $\eta = 3 \times i$  tags from every cluster for the  $i$ -th hierarchy layer. On tagged NUS-WIDE images,

tag “race” is considered more structurally abstract than tags “sky/street/house/men” by our proposed method above. This is reasonable because there exist some underlying groups (e.g. regarding Formula-1 and raft competition) that are semantically relevant with tag “race” whilst tags “sky/street/house/men” describe concrete objects that may be possibly shared by much more different data structures and hence structurally specific. Our proposed HML-RF model is formulated particularly to accommodate such abstractness skeletal knowledge in rough tag hierarchy for discovering and interpreting sparsely and/or incompletely tagged visual data, beyond conventional multi-modality correlation learning methods that often attempt to straightly correlate visual features and textual tags whilst totally ignoring tag hierarchy information. In the following experiments, we start with a two-layer tag hierarchy, then evaluate the effect of tag layer number on the model performance.

For image clustering, our aim is to reveal the category groups of the dominant scene or event presented in these web images, given the ground-truth available in group metadata [64, 65]. To evaluate the performance of different tag completion methods, we divided the full tag labels into two parts: observed part (60%) with the remaining (40%) as ground truth [6]. The observed tags were randomly chosen.

**Visual features:** For TRECVID MED 2011, we used HOG3D features [66] as visual representation of videos. In particular, we first generated a codebook of 1000 words using K-means [2]. With this codebook, we created a 1000-D histogram feature vector for each video. Finally, the approximated Histogram Intersection Kernel via feature extension [67] was adopted to further enhance the expressive capability of visual features. For NUS-WIDE, we exploited a VGG-16 convolutional neural network (CNN) [68] pre-trained on the ImageNet Large-Scale Visual Recognition Challenge 2012 dataset [69] to extract image features. This allows the image description benefiting from auxiliary rich object image annotation. Specifically, we used the output (4096-D feature vector) from the first Fully-Connected CNN layer as image feature representation.

**Implementation details:** The default parameter settings are as follows. The forest size  $\tau$  was fixed to 1000 for all random forest models. The depth of each tree was automatically determined by setting the sample number in the leaf node,  $\phi$ , which we set to 3. We set  $\nu_{\text{try}} = \sqrt{d}$  with  $d$  the data feature dimension (Equation (2)) and  $\kappa = 20$  (Equation 13). For fair comparison, we used the exactly same number of clusters, visual features and tag data in all compared methods. For any random forest model, we repeated 10 folds and reported the average results. In addition to the default settings above, we also evaluated the influence of two important HML-RF parameters, e.g.  $\tau$  and  $\phi$  (Section 4.2.3).

#### 4.2. Evaluation on Discovering Global Data Cluster Structure

**Input data modes:** For comparison, we tested four modes of input data: (1) ViFeat: videos are represented by HOG3D visual features; (2) BiTag: binary tag vectors are used instead of visual features; (3) DetScore [4]: tag classifiers (e.g. SVM) are trained for individual tags using the available tags with visual features and their detection scores are then used as model input<sup>1</sup>; (4) ViFeat&BiTag: both visual and tag data are utilised. More specifically, the two modalities may be combined into one single feature vector (called ViFeat&BiTag-cmb), or modelled separately in some balanced way (called ViFeat&BiTag-blm), depending on the design nature of specific methods.

**Baseline models:** We extensively compared our HML-RF model against the following related state-of-the-art methods: (1) K-means [2]: The most popular clustering algorithm. (2) Spectral Clustering (SpClust) [15]: A popular and robust clustering mechanism based on the eigen-vector structures of affinity matrix. In ViFeat&BiTag mode, the averaging over separate normalised affinity matrices of visual and tag data (SpClust-blm) was also evaluated, in addition to the combined single feature (SpClust-cmb). (3) Affinity Propagation (AffProp) [70]: An exemplar based clustering algorithm whose input is also affinity matrix. This method is shown insensitive to exemplar initialisation as all data samples are simultaneously considered as potential cluster centres. (4) Clustering Random Forest (ClustRF) [53, 54]: A feature selection driven data similarity computing model. It was used to generate the data affinity matrix, followed by SpClust for obtaining the final clusters. (5) Constrained-Clustering Forest (CC-Forest) [51]: A state-of-the-art multi-modality data based clustering forest characterised by joint learning of heterogeneous data. Its output is affinity matrix induced from all data modalities. Similarly, the clusters are generated by SpClust. (6) Affinity Aggregation for Spectral Clustering (AASC) [14]: A state-of-the-art multi-modal spectral clustering method that searches for an optimal weighted combination of multiple affinity matrices, each from a single data modality. (7) CCA+SpClust [20]: The popular Canonical Correlation Analysis (CCA) model that maps two views (e.g. visual and tag features) to a common latent space with the objective of maximising the correlation between the two. In this

common space, we computed pairwise similarity between samples and applied the spectral clustering algorithm to obtain clusters. (8) 3VCCA+SpClust [22]: A contemporary three-view CCA algorithm extended from the conventional CCA by additionally considering the third view about high-level semantics. Specifically, we utilised the first layer of abstract tags as the data of third view. Similarly, we used spectral clustering on the similarity measures in the induced common space for data clustering. (9) Maximum Margin Clustering (MMC) [71]: A widely used clustering model based on maximising the margin between clusters. (10) Latent Maximum Margin Clustering (L-MMC) [3]: An extended MMC model that allows to accommodate latent variables, e.g. tag labels, during maximum cluster margin learning. (11) Structural MMC (S-MMC) [4]: A variant of MMC model assuming structured tags are labelled on data samples. (12) Flip MMC (F-MMC) [4]: The state-of-the-art tag based video clustering method capable of handling the missing tag problem, beyond S-MMC. (13) Deep Canonical Correlation Analysis (DCCA) [25]: a deep neural network (DNN) based extension of CCA [20] where a separate DNN is used for extracting features of each data modality, followed by canonical correlation maximisation between across-modal features. (14) Deep Canonically Correlated Autoencoders (DCCAE) [26]: a state-of-the-art deep multi-view learning method that combines the reconstruction errors of split autoencoder [18] and the correlation maximisation of DCCA [25] in model formulation.

**Evaluation metrics:** We adopted five metrics to evaluate the clustering accuracy: (1) *Purity* [3], which calculates the averaged accuracy of the dominating class in each cluster; (2) *Normalised Mutual Information* (NMI) [72], which considers the mutual dependence between the predicted and ground-truth partitions; (3) *Rand Index* (RI) [73], which measures the ratio of agreement between two partitions, i.e. true positives within clusters and true negatives between clusters; (4) *Adjusted Rand Index* (ARI) [74], an adjusted form of RI that additionally considers disagreement, and equals 0 when the RI equals its expected value; (5) *Balanced F1 score* (F1) [75], which uniformly measures both precision and recall. All metrics lie in the range of  $[0, 1]$  except ARI in  $[-1, 1]$ . For each metric, higher values indicate better performance. Whilst there may exist some inconsistency between different metrics due to their property discrepancy [76], using all them allows to various aspects of performance measure.

##### 4.2.1. Clustering Evaluation on TRECVID MED 2011

We evaluated the effectiveness of distinct models for tag-based video clustering, using the *full* tag data along with visual features. The results are reported in Table 1. With visual features alone, all clustering methods produce poor results, e.g. the best NMI is 0.20, achieved by SpClust. Whereas binary tag representations provide much more information about the underlying video data structure than visual feature modality, e.g. all models can double their scores or even more in most metrics. Interestingly, using the detection scores can lead to even better results than the original binary tags. The plausible reason is that missing tags can be partially recovered after using

<sup>1</sup> We only compared the reported results in [4] since we cannot reproduce the exact evaluation setting due to the lack of experimental details.

Table 1: Comparing clustering methods on TRECVID MED 2011 [61].

Input mode	Method	Purity	NMI	RI	F1	ARI
ViFeat	K-means[2]	0.26	0.19	0.88	0.14	0.08
	SpClust[15]	0.25	0.20	0.88	0.15	0.07
	ClustRF[53]	0.23	0.17	0.87	0.14	0.08
	AffProp[70]	0.23	0.16	0.87	0.14	0.07
	MMC[71]	0.25	0.19	0.88	0.14	0.09
BiTag	K-means[2]	0.51	0.52	0.86	0.30	0.23
	SpClust[15]	0.71	0.73	0.93	0.56	0.60
	ClustRF[53]	0.77	0.81	0.94	0.64	0.60
	AffProp[70]	0.50	0.44	0.87	0.28	0.21
	MMC[71]	0.76	0.72	0.95	0.64	0.60
DetScore	K-means[2]	0.63	0.60	0.93	0.50	-
	SpClust[15]	0.82	0.76	0.96	0.69	-
	MMC[71]	0.83	0.78	0.96	0.73	-
	L-MMC[3]	0.86	0.82	0.97	0.79	-
ViFeat&BiTag-cmb	K-means[2]	0.51	0.49	0.90	0.34	0.24
	SpClust-cmb[15]	0.76	0.74	0.94	0.62	0.66
	ClustRF[53]	0.23	0.17	0.87	0.15	0.08
	AffProp[70]	0.51	0.46	0.86	0.29	0.21
ViFeat&BiTag-blm	SpClust-blm[15]	0.75	0.72	0.95	0.62	0.59
	CCA+SpClust[20]	0.85	0.81	0.97	0.77	0.75
	3VCCA+SpClust[22]	0.86	0.86	0.97	0.78	0.77
	CC-Forest[51]	0.41	0.33	0.89	0.41	0.19
	AASC[14]	0.30	0.15	0.87	0.13	0.06
	MMC[71]	0.79	0.72	0.95	0.66	0.66
	DCCA[25]	0.84	0.80	0.96	0.74	0.72
	DCCAE[26]	0.84	0.80	0.97	0.75	0.73
	S-MMC[4]	0.87	0.84	0.97	0.79	-
	F-MMC[4]	0.90	0.88	<b>0.98</b>	0.84	-
	HML-RF(Ours)	<b>0.94</b>	<b>0.90</b>	<b>0.98</b>	<b>0.88</b>	<b>0.87</b>

the detection scores. When using both data modalities, we observed superior results than either single modality with many methods like SpClust, AffProp, MMC. This confirms the overall benefits from jointly learning visual and tag data because of their complementary effect. Also, it is shown that separate and balanced use of visual and tag features (ViFeat&BiTag-blm) is more likely to surpass methods using concatenated visual and tag vectors (ViFeat&BiTag-cmb). A possible reason is that visual and tag features are heterogeneous to each other, a direct combination leads to an unnatural and inconsistent data representation thus likely increases the modelling difficulty and deteriorates the model performance.

For the performance of individual methods, the proposed HML-RF model evidently provides the best results by a significant margin over the second best Flip MMC in most metrics, except RI which is a less-sensitive measure due to its practical narrower range [76]. This is resulted from the joint exploitation of interactions between visual and tag data, tag hierarchical structure, and tag correlations with a unified HML-RF model (Algorithm 1), different from MMC and its variants wherein tags are exploited in a flat organisation and no tag dependences are considered. K-means hardly benefits from visual and tag combination, due to its single distance function based grouping mechanism therefore is very restricted in jointly exploiting multi-modal data.

Among all affinity based models, ClustRF is surprisingly dominated by visual data when using visual features & tag as input. This may be because that visual features with large vari-

ances may be mistakenly considered as optimum due to larger information gain induced on them. CC-Forest suffers less by separately exploiting the two modalities, but still inferior than HML-RF due to ignoring the intrinsic tag structure and the tag sparseness challenge. AASC yields much poorer clustering results than HML-RF, suggesting that the construction of individual affinity matrices can lose significant information, such as the interactions between the visual and tag data, as well as statistical tag correlations.

The methods of AffProp and SpClust-cmb also suffer from the heteroscedasticity problem in that the input affinity matrix is constructed from the heterogeneous concatenation of visual and tag data and thus ineffective to exploit the knowledge embedded across modalities and tag statistical relationships. However, separating visual and tag features does not bring benefit to SpClust (SpClust-blm). This may be due to tag sparseness and the lack of correlation modelling between visual and tag data. Whilst through correlating and optimising cross-modal latent common space, correlation analysis models (e.g. CCA, DCCA, DCCAE and 3VCCA) overcome somewhat the heterogeneous data learning challenge but remain suboptimal and inferior due to over-sparse tags and the ignorance of tag hierarchy and inter-tag correlations.

Table 2: Comparing clustering methods on NUS-WIDE [10].

Input mode	Method	Purity	NMI	RI	F1	ARI
ViFeat	K-means[2]	0.28	0.26	0.94	0.13	0.11
	SpClust[15]	0.27	0.24	0.94	0.14	0.11
	ClustRF[53]	0.27	0.24	0.94	0.14	0.11
	AffProp[70]	0.25	0.22	0.91	0.13	0.09
	MMC[71]	0.24	0.20	0.94	0.12	0.09
BiTag	K-means[2]	0.46	0.64	0.77	0.20	0.15
	SpClust[15]	0.51	0.59	0.72	0.12	0.06
	ClustRF[53]	0.57	0.60	0.90	0.15	0.33
	AffProp[70]	0.50	0.59	0.76	0.15	0.16
	MMC[71]	0.54	0.61	0.94	0.24	0.40
DetScore	K-means[2]	0.51	0.65	0.79	0.22	0.17
	SpClust[15]	0.55	0.61	0.75	0.16	0.09
	ClustRF[53]	0.60	0.62	0.92	0.17	0.35
	AffProp[70]	0.54	0.60	0.78	0.17	0.18
	MMC[71]	0.59	0.64	0.95	0.25	0.41
ViFeat&BiTag-cmb	K-means[2]	0.29	0.26	0.94	0.14	0.11
	SpClust-cmb[15]	0.28	0.24	0.94	0.14	0.11
	ClustRF[53]	0.28	0.24	0.93	0.13	0.09
	AffProp[70]	0.26	0.22	0.91	0.13	0.10
ViFeat&BiTag-blm	SpClust-blm[15]	0.58	0.56	0.87	0.19	0.14
	CCA+SpClust[20]	0.48	0.41	0.95	0.24	0.28
	3VCCA+SpClust[22]	0.52	0.45	<b>0.96</b>	0.25	0.32
	CC-Forest[51]	0.26	0.23	0.91	0.12	0.07
	AASC[14]	0.28	0.24	0.94	0.13	0.10
	MMC[71]	0.24	0.20	0.94	0.12	0.09
	DCCA[25]	0.61	0.62	0.89	0.30	0.27
	DCCAE[26]	0.62	0.63	0.89	0.30	0.27
HML-RF(Ours)	<b>0.67</b>	<b>0.67</b>	<b>0.96</b>	<b>0.32</b>	<b>0.45</b>	

#### 4.2.2. Clustering Evaluation on NUS-WIDE

We further evaluated the proposed HML-RF model and its competitors on tagged image dataset NUS-WIDE [10]. In this experiment, we utilised a two-layer tag hierarchy in HML-RF. The clustering results are reported in Table 2. It is evident



Figure 4: Example clusters discovered by the HML-RF model on NUS-WIDE [10]. Tags are shown at the right of corresponding images. The inconsistent samples are indicated with red bounding box.

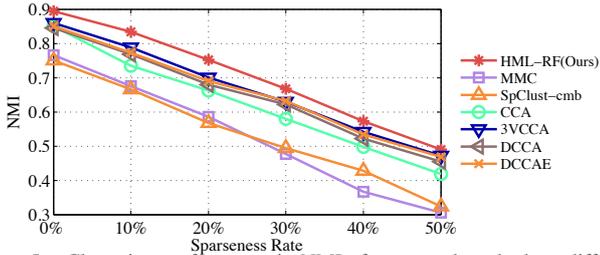


Figure 5: Clustering performance in NMI of compared methods at different tag sparseness rates on TRECVID MED 2011 [61].

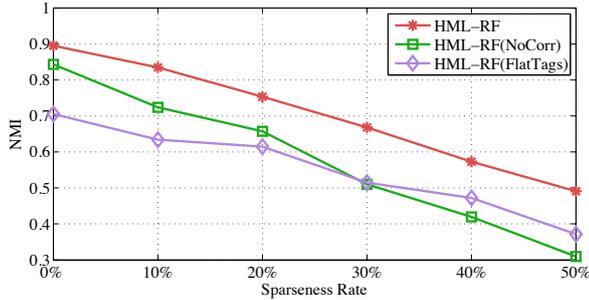


Figure 6: Evaluating the effectiveness of specific HML-RF components on TRECVID MED 2011 [61].

that our HML-RF model surpasses all baseline methods, consistent with the findings in clustering TRECVID videos. Specifically, methods based on SpClust obtain generally more accurate clusters. Interestingly, simple combination of affinity matrices (SpClust-bltn) is shown superior than latent common subspace learning (CCA and 3VCCA). This is opposite from the observations on the TRECVID videos above. A possible explanation may be due to the additional difficulty for joint subspace learning caused by the greater tag sparseness on NUS-WIDE images, e.g. missing tags making the learned projection inaccurate and suboptimal. Deep leaning based DCCA and DCCAE methods also suffer from the same problem although their stronger modelling capability can improve considerably the quality of learned subspaces. By incorporating tag hierarchy knowledge and employing automatically mined tag correlations, our HML-RF model mitigates more effectively such tag sparsity and incomplete cross-modal data alignment challenges. This again suggests the capability and effectiveness of our method in exploiting sparse tags for discovering global visual data concept structure. Example of image clusters discovered by our HML-RF are shown in Figure 4.

#### 4.2.3. Further Analysis

We further conducted a series of in-depth evaluations and analysis: (1) model robustness against tag sparseness; (2) HML-RF model component effect; (3) HML-RF model parameter sensitivity; and (4) tag hierarchy structure effect.

**Model robustness against tag sparseness:** We conducted a scalability evaluation against tag sparseness and incompleteness. This is significant since we may have access to merely a small size of tags in many practical settings. To simulate these scenarios, we randomly removed varying ratios (10% ~ 50%) of tag data on the TRECVID MED 2011 dataset. We utilised

Table 3: Comparing relative drop in NMI of top-7 clustering models, given different tag sparseness rates on TRECVID MED 2011 [61]. Smaller is better.

Sparseness rate (%)	10	20	30	40	50
SpClust-cmb[15]	0.11	0.24	0.34	0.43	0.57
MMC[71]	0.12	0.24	0.38	0.52	0.60
CCA+SpClust[20]	0.14	0.22	0.32	0.41	0.51
3VCCA+SpClust[22]	0.08	0.19	0.27	0.37	<b>0.45</b>
DCCA[20]	0.09	0.20	0.26	0.38	0.46
DCCAE[22]	0.09	0.19	0.26	0.37	<b>0.45</b>
HML-RF(Ours)	<b>0.07</b>	<b>0.16</b>	<b>0.25</b>	<b>0.36</b>	<b>0.45</b>

both visual and tag data as model input since most methods can benefit from using both<sup>2</sup>. The most common metric NMI [2] was used in this experiment.

The results by top-7 clustering methods are compared in Figure 5. Given less amount of tag data, as expected we observe a clear performance drop trend across all these models. However, the relative drops in the performance of HML-RF model due to tag incompleteness are the smallest among all compared methods at 10% ~ 40% sparseness rate (less is more sparse). This performance degradation is comparable among three best models (HML-RF, 3VCCA and DCCAE) at 50% sparseness rate, as shown in Table 3. This demonstrates the robustness and benefits of the proposed HML-RF model with respect to tag sparseness and incompleteness, and making it more practically useful when fewer tags are available. This also demonstrates that a joint exploitation of visual features, tags hierarchy as well as tag correlations can bring about significant benefits to visual semantic structure interpretation and global video clustering with sparse/incomplete tags. For qualitative visualisation, an example of clusters formed by our HML-RF under the most sparse case is given in Figure 7.

**HML-RF model component effect:** We explicitly examined two components of the proposed HML-RF for casting light on model formulation: (1) the effect of exploiting tag abstractness hierarchy structure; and (2) the influence of tag statistical correlations. To that end, we build two stripped-down variants of HML-RF: (I) HML-RF(FlatTags): A HML-RF without exploiting tag hierarchy and tag correlations (Equation (4)); (II) HML-RF(NoCorr): A HML-RF without tag correlation (Equation (5)). Contrasting the performance between HML-RF(FlatTags) and HML-RF(NoCorr) allows for measuring the former, whilst that between HML-RF(NoCorr) and HML-RF for the later. We repeated the same experiments as above with the two variants.

It is evident from Figure 6 that both components make significant differences but their relative contribution varies under different tag sparseness cases. Particularly, given the full tags, tag abstractness hierarchy plays a significant role, e.g. boosting NMI from 0.71 to 0.84; but when more sparse tag data is utilised, the performance gain decreases and even drops at > 30% sparseness rates. However, combining with tag correlations can effectively increase the clustering accuracy. This indicates that the tag hierarchy component works under certain tag

<sup>2</sup> Structural MMC and Flip MMC models [4] were not included in this evaluation due to the difficulties in reproducing their models from a lack of sufficient implementation details.

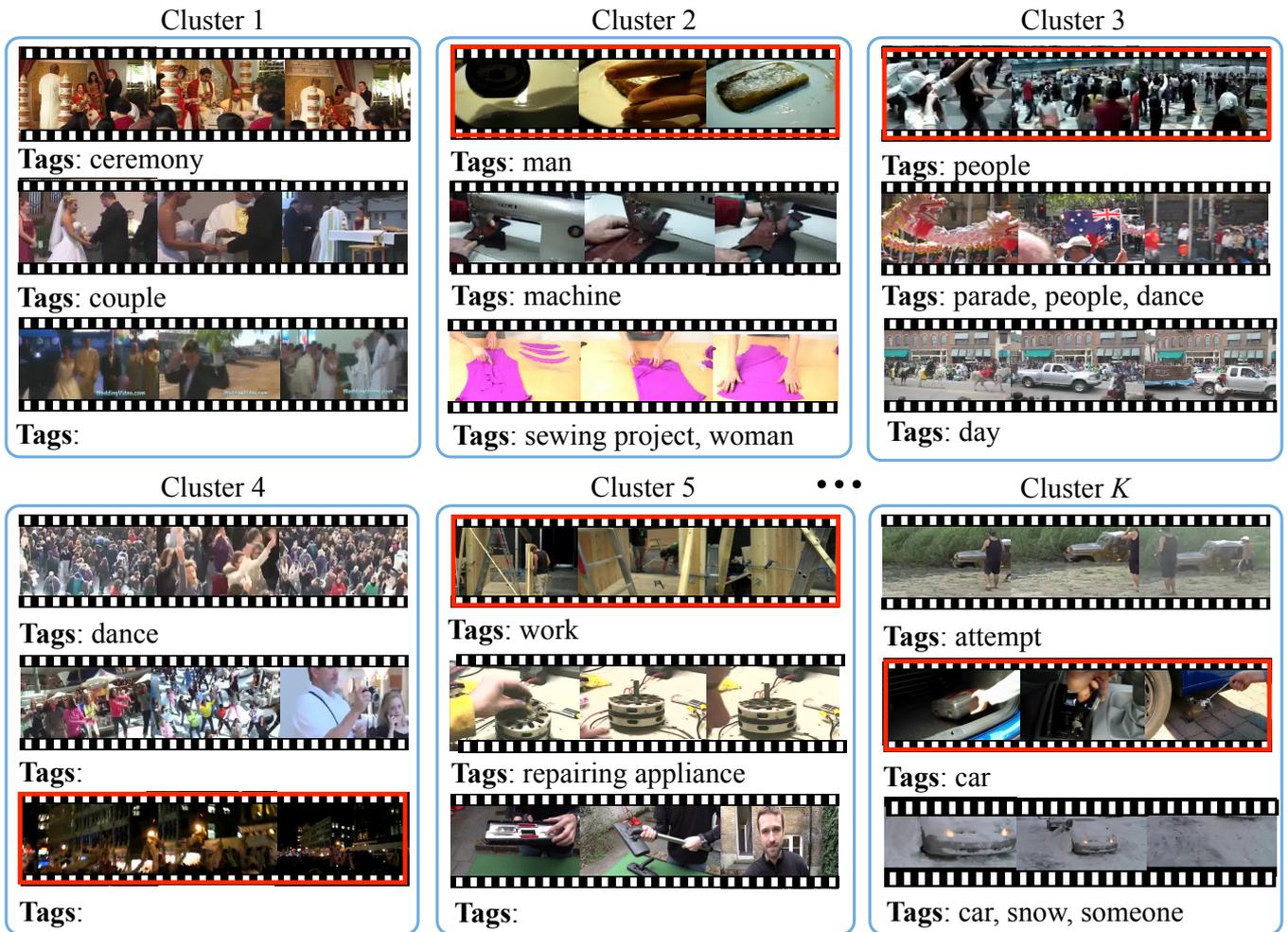


Figure 7: Example clusters formed by the HML-RF model given 50% tag sparseness rate on TRECVID MED 2011 [61]. Tags are shown underneath the corresponding video. Not that some videos have no tag data. Inconsistent samples are indicated with red bounding box.

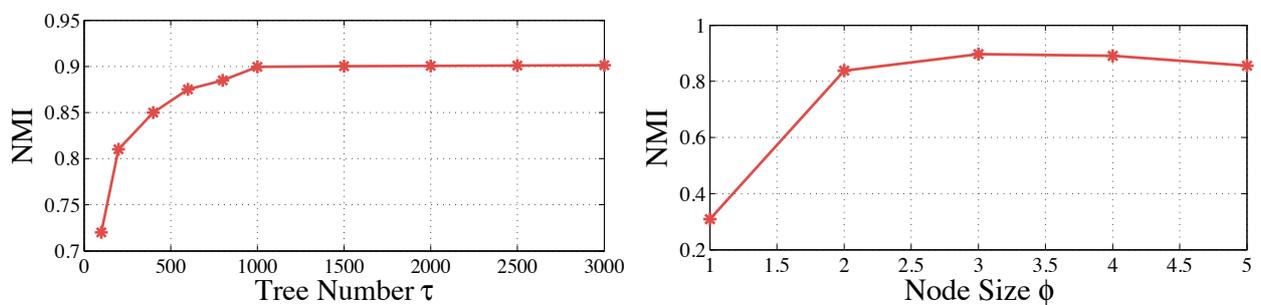


Figure 8: Clustering performance in NMI of HML-RF over different forest sizes ( $\tau$ ) and node size ( $\phi$ ) on TRECVID MED 2011 [61].

Table 4: Evaluating the effect of tag hierarchy layer number in clustering performance by our HML-RF model on NUS-WIDE [10].

Tag layer number	Purity	NMI	RI	F1	ARI
2	0.67	0.67	<b>0.96</b>	0.32	0.45
3	0.68	0.68	<b>0.96</b>	0.34	0.47
4	0.70	0.69	<b>0.96</b>	<b>0.35</b>	<b>0.48</b>
5	<b>0.71</b>	<b>0.70</b>	<b>0.96</b>	<b>0.35</b>	<b>0.48</b>
6	<b>0.71</b>	<b>0.70</b>	<b>0.96</b>	<b>0.35</b>	<b>0.48</b>
7	<b>0.71</b>	<b>0.70</b>	<b>0.96</b>	<b>0.35</b>	<b>0.48</b>

densities and coordinates well with tag correlations particularly in sparse tag cases. On the other hand, an opposite phenomenon takes place with tag correlations, i.e. it brings large benefit (from 0.31 to 0.49) in the most sparse case. These observations suggest that the two components are complementary and both are important constituents of the unified HML-RF model.

**HML-RF model parameter sensitivity:** We evaluated two key parameters in HML-RF: Tree number  $\tau$  and leaf node size  $\phi$ . The results are given in Figure 8. It is evident that when more trees are trained and utilised, the clustering accuracy increases monotonically and starts to converge from  $\tau = 1000$ . This is consistent with the findings in [45, 77]. When  $\phi = 1$ , weaker clustering results are obtained. This makes sense because HML-trees are overly grown, e.g. they enforce very similar data samples to be separated and thus make the pairwise affinity estimation inaccurate (Section 3.3). Setting small values to  $\phi$  significantly improves the clustering accuracy, and is shown to be insensitive w.r.t. specific numbers.

**Tag hierarchy structure effect:** Apart from two-layer tag hierarchy, we further evaluated the effect of tag layer number on the clustering performance of our HML-RF model on the NUS-WIDE [10] dataset. Specifically, we evaluated different tag hierarchies ranging from 3 to 7 layers, and the results are shown in Table 4. We made these observations: (1) The layer number of tag hierarchy can affect the results of data structure discovery by our HML-RF model; (2) The NUS-WIDE tags may lie in multiple abstractness layers, which leads to better discovered cluster structure than that by two layers; (3) The performance starts to get saturated from five layers and appending further more layers has little effect on data structure discovery, probably due to that over specific tags have little influence on data structure. These findings imply the effectiveness and robustness of HML-RF in accommodating tag hierarchies of various structures and qualities.

**Tag abstractness effect:** We further evaluated the benefit of tag abstractness by comparing (i) the 2-layers tag hierarchy structure with (ii) a 1-layer structure of the most specific tags in the proposed HML-RF model. Table 5 shows a significant performance advantage from exploiting a hierarchical tag abstractness structure for data clustering on both the TRECVID MED 2011 and the NUS-WIDE datasets. This demonstrates more clearly the effectiveness of HML-RF in mining and exploiting semantic information from multiple levels of tag abstractness for global data structure analysis.

### 4.3. Evaluation on Completing Local Instance-Level Concept Structure

**Baseline methods:** We compared our missing tag completion method (all three algorithms) for completing local instance-level semantic concept against the following three contemporary approaches: (1) Linear Sparse Reconstructions (LSR) [7]: A state-of-the-art image-specific and tag-specific Linear Sparse Reconstruction scheme for tag completion. (2) Tag Completion by Matrix Recovery (TCMR) [6]: A recent tag matrix recovery based completion algorithm that captures both underlying tag dependency and visual consistency. (3) A group of cluster based completion methods: Specifically, we used the same algorithm as HML-RF(GC) for missing tag recovery (Section 3.4). The clusters were obtained by the compared methods in Section 3.3. For HML-RF, we utilised the clustering results by the five-layer hierarchy. Similarly, we name these completion methods in form of ‘‘ClusteringMethodName(GC)’’, e.g. MMC(GC).

**Evaluation metrics:** We utilised three performance measures: (1) AP@N, which measures Average Precision of N recovered tags. (2) AR@N, which calculates Average Recall of N recovered tags, i.e. the percentage of correctly recovered tags over all ground truth missing tags. (3) Coverage@N, which denotes the percentage of samples with at least one correctly recovered tag when N tags are completed.

#### 4.3.1. Missing Tag Completion Evaluation on TRECVID

The tag completion results on TRECVID MED 2011 are given in Tables 6 and 7. It is evident that the proposed completion algorithms outperform all compared methods. In particular, it is observed that global clusters provide strong cues for missing tag recovery, e.g. DCCAE is superior than or similar to the state-of-the-art completion methods TCMR and LSR at AP@1. This suggests the intrinsic connection between global and local semantic structures, and validates our motivation for bridging the two visual data structure analysis tasks (Section 3.4). By more accurate global group structure revelation, HML-RF(GC) enables even better missing tag completion, e.g. obtaining higher average precision and recall than other clustering methods. Moreover, HML-RF(GC) produces better tag recovery than our local neighbourhood based completion method HML-RF(LN), particularly in cases of completing multiple tags. This further indicates the positive restricting effect of global data structures over inferring local instance-level semantic concept structures. However, HML-RF(LN) provides best AR@1, which should be due to its strict rule on selecting neighbourhoods. While TCMR considers both tag correlation as well as visual consistency, it is still inferior to the proposed HML-RF owing potentially to (1) the incapability of exploiting the tag abstract-to-specific hierarchy knowledge; and (2) the assumptions on low rank matrix recovery may be not fully satisfied given real-world visual data. These observations and analysis demonstrate the superiority of our HML-RF in instance-level tag completion, owing to its favourable capability in jointly learning heterogeneous visual and tag data and thus more accurate semantic visual structure disclosure.

Table 5: Evaluating the effect of tag abstractness in the HML-RF model on data clustering.

Dataset	Tag Structure	Purity	NMI	RI	F1	ARI
TRECVID MED 2011 [61]	1-Layer Most Specific Tags	0.39	0.33	0.88	0.24	0.18
	2-Layers Hierarchy Tags	<b>0.94</b>	<b>0.90</b>	<b>0.98</b>	<b>0.88</b>	<b>0.87</b>
NUS-WIDE [10]	1-Layer Most Specific Tags	0.25	0.24	0.92	0.11	0.07
	2-Layers Hierarchy Tags	<b>0.67</b>	<b>0.67</b>	<b>0.96</b>	<b>0.32</b>	<b>0.45</b>

Table 6: Comparing Precision & Recall between tag completion methods on TRECVID MED 2011 [61].

Metric	AP@N					AR@N					
	Recovered tag # N	1	2	3	4	5	1	2	3	4	5
LSR [7]	0.31	0.25	0.22	0.20	0.17	0.15	0.24	0.32	0.38	0.40	
TCMR [6]	0.35	0.27	0.24	0.22	0.20	0.17	0.26	0.35	0.43	0.48	
AASC(GC)	0.23	0.17	0.14	0.13	0.11	0.12	0.18	0.21	0.25	0.27	
SpClust-blnc(GC)	0.31	0.27	0.25	0.23	0.21	0.15	0.25	0.36	0.44	0.48	
CC-Forest(GC)	0.28	0.24	0.18	0.15	0.14	0.15	0.23	0.26	0.27	0.31	
CCA+SpClust(GC)	0.34	0.29	0.26	0.26	0.23	0.16	0.26	0.34	0.47	0.52	
3VCCA+SpClust(GC)	0.35	0.29	0.26	0.26	0.23	0.17	0.27	0.34	0.47	0.52	
MMC(GC)	0.32	0.25	0.23	0.24	0.21	0.15	0.24	0.36	0.45	0.49	
DCCA(GC)	0.35	0.29	0.26	0.26	0.23	0.17	0.27	0.34	0.47	0.52	
DCCAE(GC)	0.36	0.29	<b>0.27</b>	0.26	0.24	0.17	0.27	0.35	0.47	0.53	
HML-RF(GC)	0.36	<b>0.31</b>	<b>0.27</b>	<b>0.27</b>	<b>0.25</b>	0.17	<b>0.29</b>	<b>0.37</b>	<b>0.49</b>	<b>0.56</b>	
HML-RF(LN)	0.37	0.29	0.25	0.23	0.20	<b>0.19</b>	0.28	0.34	0.44	0.49	
HML-RF(AM)	<b>0.38</b>	0.30	0.26	0.24	0.22	0.18	0.27	0.36	0.44	0.50	

Table 7: Comparing Coverage@N between different tag completion methods.

Dataset	TRECVID MED 2011 [61]					NUS-WIDE [10]					
	Recovered tag # N	1	2	3	4	5	1	2	3	4	5
LSR [7]	0.31	0.43	0.52	0.59	0.61	0.30	0.35	0.38	0.40	0.42	
TCMR [6]	0.35	0.46	0.57	0.66	0.71	0.25	0.33	0.39	0.43	0.46	
AASC(GC)	0.23	0.33	0.38	0.43	0.46	0.09	0.14	0.17	0.22	0.22	
SpClust-blnc(GC)	0.31	0.44	0.55	0.63	0.65	0.15	0.21	0.25	0.29	0.33	
CC-Forest(GC)	0.28	0.43	0.47	0.48	0.52	0.08	0.13	0.17	0.21	0.21	
CCA+SpClust(GC)	0.34	0.45	0.56	0.65	0.70	0.15	0.22	0.27	0.32	0.36	
3VCCA+SpClust(GC)	0.35	0.46	0.56	0.65	0.70	0.16	0.23	0.28	0.32	0.36	
MMC(GC)	0.32	0.42	0.55	0.66	0.70	0.10	0.15	0.18	0.24	0.23	
DCCA(GC)	0.35	0.46	0.56	0.66	0.70	0.18	0.21	0.27	0.29	0.33	
DCCAE(GC)	0.36	0.47	0.57	0.66	0.71	0.18	0.23	0.27	0.29	0.33	
HML-RF(GC)	0.36	<b>0.49</b>	<b>0.59</b>	<b>0.68</b>	<b>0.75</b>	0.20	0.26	0.30	0.32	0.35	
HML-RF(LN)	0.37	<b>0.49</b>	0.56	0.65	0.68	0.29	0.35	0.39	0.41	0.42	
HML-RF(AM)	<b>0.38</b>	0.47	0.58	0.65	0.70	<b>0.34</b>	<b>0.41</b>	<b>0.45</b>	<b>0.48</b>	<b>0.50</b>	

Table 8: Comparing Precision & Recall between different tag completion methods on NUS-WIDE [10].

Metric	AP@N					AR@N					
	Recovered tag # N	1	2	3	4	5	1	2	3	4	5
LSR [7]	0.30	0.22	0.18	0.15	0.13	0.15	0.21	0.24	0.27	0.28	
TCMR [6]	0.25	0.19	0.16	0.15	0.13	0.13	0.19	0.23	0.26	0.29	
AASC(GC)	0.09	0.09	0.09	0.07	0.07	0.05	0.07	0.10	0.12	0.15	
SpClust-blnc(GC)	0.15	0.12	0.10	0.08	0.09	0.09	0.13	0.15	0.20	0.20	
CC-Forest(GC)	0.08	0.09	0.09	0.07	0.07	0.04	0.07	0.09	0.12	0.15	
CCA+SpClust(GC)	0.15	0.13	0.12	0.11	0.09	0.09	0.13	0.16	0.20	0.21	
3VCCA+SpClust(GC)	0.16	0.14	0.13	0.11	0.09	0.09	0.14	0.17	0.20	0.23	
MMC(GC)	0.10	0.09	0.09	0.07	0.07	0.06	0.07	0.11	0.12	0.17	
DCCA(GC)	0.18	0.12	0.11	0.09	0.09	0.10	0.13	0.15	0.18	0.19	
DCCAE(GC)	0.18	0.13	0.11	0.09	0.09	0.10	0.13	0.15	0.18	0.19	
HML-RF(GC)	0.20	0.15	0.13	0.10	0.09	0.11	0.14	0.16	0.18	0.19	
HML-RF(LN)	0.29	0.20	0.17	0.14	0.11	0.15	0.20	0.23	0.25	0.26	
HML-RF(AM)	<b>0.34</b>	<b>0.24</b>	<b>0.20</b>	<b>0.17</b>	<b>0.15</b>	<b>0.18</b>	<b>0.24</b>	<b>0.28</b>	<b>0.30</b>	<b>0.32</b>	

### TRECVID MED 2011



**Obs:** home video, group  
**GT:** people, man, street  
**Pr:** dance, music, **street**



**Obs:** home video, parade  
**GT:** children, school  
**Pr:** street, people, **school**



**Obs:** amateur footage, sewing-project, woman  
**GT:** people, indoors, work, hand, girl, table  
**Pr:** machine, **indoors, hand**



**Obs:** home video, people, party  
**GT:** indoors, children  
**Pr:** **indoors, children** music



**Obs:** home video, flash mob  
**GT:** music, dance, group  
**Pr:** **dance, group, music**



**Obs:** home video, parkour  
**GT:** jump, man  
**Pr:** **man, jump, boy**



**Obs:** amateur footage  
**GT:** people, street  
**Pr:** **street, music, vehicle**



**Obs:** demonstration, man  
**GT:** machine, hand  
**Pr:** **machine, wood, hand**



**Obs:** amateur footage, sandwich, truck, man  
**GT:** indoors, kitchen, hand, talk, food  
**Pr:** **indoors, children, kitchen**

### NUS WIDE



**Obs:** travel, turkey  
**GT:** bazaar, istanbul  
**Pr:** **bazaar, market, istanbul**



**Obs:** yellow, colour, wales, host  
**GT:** flowers, spring, bright, day, contrast  
**Pr:** green, **flowers, macro**



**Obs:** texas, action, blur, actor, austin  
**GT:** dance  
**Pr:** portrait, film, athlete



**Obs:** police, free  
**GT:** china, protest, riot, protesters  
**Pr:** **protesters, protest, demonstration**



**Obs:** travel, Europe, maps  
**GT:**  
**Pr:** paper, art, earth



**Obs:** people, india, market  
**GT:** portrait, closeup, bazaar  
**Pr:** **bazaar, portrait, actor**



**Obs:** hangar  
**GT:**  
**Pr:** formula, car, airport



**Obs:** blue, winter, film, cold  
**GT:** ice, cold, outdoors, frozen, crystals  
**Pr:** **crystals, ice, frost**



**Obs:** nets  
**GT:** basketball  
**Pr:** **basketball, athlete, protesters**

Figure 9: Examples of tag completion by our HML-RF(AM) method. Correctly recovered tags are highlighted in green colour. (Obs: Observed tags; GT: Ground Truth for missing tags; Pr: Predicted tags).

#### 4.3.2. Missing Tag Completion Evaluation on NUS-WIDE

Tables 8 and 7 show the comparative results for tag completion on the NUS-WIDE image dataset [10], where the available tags are more sparse (0.48%) as compared to the TRECVID MED 2011 video dataset (3.5%). Overall, our methods HML-RF(AM) outperforms all other baselines, including the state-of-the-art models LSR and TCMR, and contemporary deep-based multi-modal correlation learning methods DCCA and DCCAE. We found that our HML-RF(GC) model does not perform as strongly as on TRECVID MED 2011. This shall be due to less accurate global group structures discovered (see Table 2). By imposing stringent neighbourhood selection, HML-RF(LN) produces considerably better tag recovery accuracy than HML-RF(GC). This validates the proposed pure neighbourhood based completion strategy in handling sparse and incomplete tags where a large number of missing tags can negatively bias tag recovery (Section 3.4). HML-RF(AM) achieves the best results due to the combined benefits from both local and global neighbourhood structures. These evaluations and observations further validate the capability and efficacy of the proposed model in jointly learning heterogeneous visual and tag modalities and semantically interpreting the instance-level concept structure of ambiguous visual content in both video and image data. For qualitative evaluation, we show in Figure 9 the top-3 recovered tags per sample by our HML-RF(AM) method.

## 5. Conclusion

In this work, we presented an visual concept structure discovery framework by formulating a novel Hierarchical-Multi-Label Random Forest (HML-RF) model for jointly exploiting heterogeneous visual and tag data modalities, with the aim of creating an intelligent visual machine for automatically organising and managing large scale visual databases. The proposed new forest model, which is defined by a new information gain function, enables naturally incorporating tag abstractness hierarchy and effectively exploiting multiple tag statistical correlations, beyond modelling the intrinsic interactions between visual and tag modalities. With the learned HML-RF, we further derive a generic clustering pipeline for global group structure discovery and three tag completion algorithms for local instance-level tag concept structure recovery. Extensive comparative evaluations have demonstrated the advantages and superiority of the proposed approach over a wide range of existing state-of-the-arts clustering, multi-view embedding and tag completion models, particularly in cases where only sparse tags are accessible. Further, a detailed model component examination is provided for casting insights on our modelling principles and model robustness. In addition to the above two applications, our HML-RF model can potentially benefit other related problems, such as retrieval and manifold ranking.

## Acknowledgements

This work was partially supported by the China Scholarship Council, Vision Semantics Limited, and Royal Society Newton Advanced Fellowship Programme (NA150459). The corresponding author is Xiatian Zhu.

## References

### References

- [1] K. Beyer, J. Goldstein, R. Ramakrishnan, U. Shaft, When is nearest neighbor meaningful?, in: Database TheoryICDT99, 1999, pp. 217–235.
- [2] A. K. Jain, Data clustering: 50 years beyond k-means, *Pattern Recognition Letters* 31 (8) (2010) 651–666.
- [3] G.-T. Zhou, T. Lan, A. Vahdat, G. Mori, Latent margin clustering, in: Advances in Neural Information Processing Systems, Lake Tahoe, Nevada, USA, 2013, pp. 28–36.
- [4] A. Vahdat, G.-T. Zhou, G. Mori, Discovering video clusters from visual features and noisy tags, in: European Conference on Computer Vision, Zurich, Switzerland, 2014, pp. 526–539.
- [5] L. Wu, R. Jin, A. K. Jain, Tag completion for image retrieval, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (3) (2013) 716–727.
- [6] Z. Feng, S. Feng, R. Jin, A. K. Jain, Image tag completion by noisy matrix recovery, in: European Conference on Computer Vision, Zurich, Switzerland, 2014, pp. 424–438.
- [7] Z. Lin, G. Ding, M. Hu, J. Wang, X. Ye, Image tag completion via image-specific and tag-specific linear sparse reconstructions, in: IEEE Conference on Computer Vision and Pattern Recognition, Portland, Oregon, United States, 2013, pp. 1618–1625.
- [8] B. T. Truong, S. Venkatesh, Video abstraction: A systematic review and classification, *ACM Transactions on Multimedia Computing, Communications, and Applications* 3 (1) (2007) 3.
- [9] R. Duin, M. Loog, Linear dimensionality reduction via a heteroscedastic extension of lda: the chernoff criterion, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 26 (6) (2004) 732–739.
- [10] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, Y. Zheng, Nus-wide: a real-world web image database from national university of singapore, in: ACM international Conference on Image and Video Retrieval, Santorini, Greece., 2009, p. 48.
- [11] A. Vahdat, G. Mori, Handling uncertain tags in visual recognition, in: IEEE International Conference on Computer Vision, Sydney, Australia, 2013, pp. 737–744.
- [12] P. Natarajan, S. Wu, S. Vitaladevuni, X. Zhuang, S. Tsakalidis, U. Park, R. Prasad, P. Natarajan, Multimodal feature fusion for robust event detection in web videos, in: IEEE Conference on Computer Vision and Pattern Recognition, Providence, Rhode Island, 2012, pp. 1298–1305.
- [13] A. Makadia, V. Pavlovic, S. Kumar, A new baseline for image annotation, in: European Conference on Computer Vision, Berlin, Heidelberg, 2008, pp. 316–329.
- [14] H.-C. Huang, Y.-Y. Chuang, C.-S. Chen, Affinity aggregation for spectral clustering, in: IEEE Conference on Computer Vision and Pattern Recognition, Marseille, France, 2012, pp. 773–780.
- [15] A. Y. Ng, M. I. Jordan, Y. Weiss, et al., On spectral clustering: Analysis and an algorithm, in: Advances in Neural Information Processing Systems, Vol. 2, Vancouver, British, 2002, pp. 849–856.
- [16] N. Quadrianto, C. H. Lampert, Learning multi-view neighborhood preserving projections, in: International Conference on Machine Learning, Bellevue, Washington, United States, 2011, pp. 425–432.
- [17] N. Srivastava, R. R. Salakhutdinov, Multimodal learning with deep boltzmann machines, in: Advances in Neural Information Processing Systems, Lake Tahoe, Nevada, United States, 2012, pp. 2222–2230.
- [18] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, A. Y. Ng, Multimodal deep learning, in: International Conference on Machine Learning, 2011, pp. 689–696.
- [19] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al., Devise: A deep visual-semantic embedding model, in: Advances in Neural Information Processing Systems, 2013, pp. 2121–2129.
- [20] D. R. Hardoon, S. Szedmak, J. Shawe-Taylor, Canonical correlation analysis: An overview with application to learning methods, *Neural computation* 16 (12) (2004) 2639–2664.
- [21] S. J. Hwang, K. Grauman, Learning the relative importance of objects from tagged images for retrieval and cross-modal search, *International Journal of Computer Vision* 100 (2) (2012) 134–153.
- [22] Y. Gong, Q. Ke, M. Isard, S. Lazebnik, A multi-view embedding space for modeling internet images, tags, and their semantics, *International Journal of Computer Vision* 106 (2) (2014) 210–233.
- [23] P. Rai, H. Daume, Multi-label prediction via sparse infinite cca, in: Advances in Neural Information Processing Systems, Vancouver, British Columbia, Canada, 2009, pp. 1518–1526.
- [24] A. Sharma, A. Kumar, H. Daume III, D. W. Jacobs, Generalized multi-view analysis: A discriminative latent space, in: IEEE Conference on Computer Vision and Pattern Recognition, Providence, Rhode Island, United States, 2012, pp. 2160–2167.
- [25] G. Andrew, R. Arora, J. A. Bilmes, K. Livescu, Deep canonical correlation analysis., in: International Conference on Machine Learning, 2013, pp. 1247–1255.
- [26] W. Wang, R. Arora, K. Livescu, J. Bilmes, On deep multi-view representation learning, in: International Conference on Machine Learning, 2015, pp. 1083–1092.
- [27] N. Natarajan, I. S. Dhillon, P. K. Ravikumar, A. Tewari, Learning with noisy labels, in: Advances in Neural Information Processing Systems, 2013, pp. 1196–1204.
- [28] S. Sukhbaatar, J. Bruna, M. Paluri, L. Bourdev, R. Fergus, Training convolutional networks with noisy labels, *Workshop of International Conference in Learning Representations*.
- [29] B. Frénaý, M. Verleysen, Classification in the presence of label noise: a survey, *IEEE Transactions on Neural Networks and Learning Systems* 25 (5) (2014) 845–869.
- [30] R. S. Cabral, F. Torre, J. P. Costeira, A. Bernardino, Matrix completion for multi-label image classification, in: Advances in Neural Information Processing Systems, Granada, Spain, 2011, pp. 190–198.
- [31] Y. Mu, J. Dong, X. Yuan, S. Yan, Accelerated low-rank visual recovery by random projection, in: IEEE Conference on Computer Vision and Pattern Recognition, Colorado Springs, USA, 2011, pp. 2609–2616.
- [32] X. Liu, S. Yan, T.-S. Chua, H. Jin, Image label completion by pursuing contextual decomposability, *ACM Transactions on Multimedia Computing, Communications, and Applications* 8 (2) (2012) 21.
- [33] E. J. Candès, B. Recht, Exact matrix completion via convex optimization, *Foundations of Computational Mathematics* 9 (6) (2009) 717–772.
- [34] C. Fellbaum, *WordNet*, Wiley Online Library, 1998.
- [35] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: IEEE Conference on Computer Vision and Pattern Recognition, Miami, Florida, United States, 2009, pp. 248–255.
- [36] Personalized recommendation in social tagging systems using hierarchical clustering, in: Proceedings of the 2008 ACM conference on Recommender systems, New York, NY, USA, 2008, pp. 259–266.
- [37] S. Zheng, M.-M. Cheng, J. Warrell, P. Sturgess, V. Vineet, C. Rother, P. H. Torr, Dense semantic image segmentation with objects and attributes, in: IEEE Conference on Computer Vision and Pattern Recognition, Columbus, Ohio, 2014, pp. 3214–3221.
- [38] J. Deng, N. Ding, Y. Jia, A. Frome, K. Murphy, S. Bengio, Y. Li, H. Neven, H. Adam, Large-scale object classification using label relation graphs, in: European Conference on Computer Vision, Zurich, Switzerland, 2014, pp. 48–64.
- [39] T. Griffiths, Z. Ghahramani, Infinite latent feature models and the indian buffet process, in: Advances in Neural Information Processing Systems, 2005, pp. 475–482.
- [40] X. Chen, Y. Mu, S. Yan, T.-S. Chua, Efficient large-scale image annotation by probabilistic collaborative multi-label propagation, in: ACM International Conference on Multimedia, New York, NY, USA, 2010, pp. 35–44.
- [41] M. J. Choi, J. J. Lim, A. Torralba, A. S. Willsky, Exploiting hierarchical context on a large database of object categories, in: IEEE Conference on Computer Vision and Pattern Recognition, San Francisco, California, USA, 2010, pp. 129–136.
- [42] C. Desai, D. Ramanan, C. C. Fowlkes, Discriminative models for multi-class object layout, *International Journal of Computer Vision* 95 (1)

- (2011) 1–12.
- [43] X. Chen, X.-T. Yuan, Q. Chen, S. Yan, T.-S. Chua, Multi-label visual classification with label exclusive context, in: *IEEE International Conference on Computer Vision*, Barcelona, Spain, 2011, pp. 834–841.
- [44] Y. Zhou, R. Jin, S. Hoi, Exclusive lasso for multi-task feature selection, in: *International Conference on Artificial Intelligence and Statistics*, Chia Laguna Resort, Sardinia, Italy, 2010, pp. 988–995.
- [45] A. Criminisi, J. Shotton, E. Konukoglu, Decision forests: A unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning, *Foundations and Trends® in Computer Graphics and Vision* 7 (2–3) (2012) 81–227.
- [46] X. Zhu, C. C. Loy, S. Gong, Constrained clustering with imperfect oracles, *IEEE Transactions on Neural Networks and Learning Systems* 27 (6) (2016) 1345–1357.
- [47] X. Zhu, C. C. Loy, S. Gong, Constrained clustering: Effective constraint propagation with imperfect oracles, in: *IEEE International Conference on Data Mining*, 2013, pp. 1307–1312.
- [48] X. Zhu, C. Change Loy, S. Gong, Constructing robust affinity graphs for spectral clustering, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1450–1457.
- [49] A. Montillo, J. Shotton, J. Winn, J. E. Iglesias, D. Metaxas, A. Criminisi, Entangled decision forests and their application for semantic segmentation of ct images, in: *Information Processing in Medical Imaging*, Kloster Irsee, Germany, 2011, pp. 184–196.
- [50] X. Zhao, T.-K. Kim, W. Luo, Unified face analysis by iterative multi-output random forests, in: *IEEE Conference on Computer Vision and Pattern Recognition*, Columbus, Ohio, 2014, pp. 1765–1772.
- [51] X. Zhu, C. C. Loy, S. Gong, Video synopsis by heterogeneous multi-source correlation, in: *IEEE International Conference on Computer Vision*, Sydney, Australia, 2013, pp. 81–88.
- [52] X. Zhu, C. C. Loy, S. Gong, Learning from multiple sources for video summarisation, *International Journal of Computer Vision* 117 (3) (2016) 247–268.
- [53] L. Breiman, Random forests, *Machine Learning* 45 (1) (2001) 5–32.
- [54] T. Shi, S. Horvath, Unsupervised learning with random forest predictors, *Journal of Computational and Graphical Statistics* 15 (1) (2006) 118–138.
- [55] L. Breiman, J. Friedman, C. Stone, R. Olshen, *Classification and regression trees*, Chapman & Hall/CRC, 1984.
- [56] B. Liu, Y. Xia, P. S. Yu, Clustering through decision tree construction, in: *Ninth international conference on Information and knowledge management*, McLean, VA, USA, 2000, pp. 20–29.
- [57] A. Berger, R. Caruana, D. Cohn, D. Freitag, V. Mittal, Bridging the lexical chasm: statistical approaches to answer-finding, in: *ACM SIGIR Conference on Research and Development in Information Retrieval*, Athens, Greece, 2000.
- [58] J. Sivic, A. Zisserman, Video Google: A text retrieval approach to object matching in videos, in: *IEEE International Conference on Computer Vision*, Nice, France, 2003, pp. 1470–1477.
- [59] T. M. Cover, P. E. Hart, Nearest neighbor pattern classification, *IEEE Transactions on Information Theory* 13 (1) (1967) 21–27.
- [60] K. Q. Weinberger, L. K. Saul, Distance metric learning for large margin nearest neighbor classification, *The Journal of Machine Learning Research* 10 (2009) 207–244.
- [61] P. Over, G. M. Awad, J. Fiscus, B. Antonishek, M. Michel, A. F. Smeaton, W. Kraaij, G. Quénot, Trecvid 2010—an overview of the goals, tasks, data, evaluation mechanisms, and metrics (2011) 1–52.
- [62] B. Zhao, F. Wang, C. Zhang, Efficient multiclass maximum margin clustering, in: *International Conference on Machine Learning*, ACM, 2008, pp. 1248–1255.
- [63] X. Wei, W. B. Croft, Lda-based document models for ad-hoc retrieval, in: *ACM SIGIR Conference on Research and Development in Information Retrieval*, ACM, 2006, pp. 178–185.
- [64] J. Johnson, L. Ballan, L. Fei-Fei, Love thy neighbors: Image annotation by exploiting image metadata, in: *IEEE International Conference on Computer Vision*, 2015, pp. 4624–4632.
- [65] H. Hu, G.-T. Zhou, Z. Deng, Z. Liao, G. Mori, Learning structured inference neural networks with label relations, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2960–2968.
- [66] A. Klaser, M. Marszałek, C. Schmid, A spatio-temporal descriptor based on 3d-gradients, in: *British Machine Vision Conference*, Leeds, UK, 2008, pp. 275–1.
- [67] A. Vedaldi, A. Zisserman, Efficient additive kernels via explicit feature maps 34 (3) (2012) 480–492.
- [68] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: *International Conference on Learning Representation*, 2015.
- [69] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, *International Journal of Computer Vision* 115 (3) (2015) 211–252.
- [70] B. J. Frey, D. Dueck, Clustering by passing messages between data points, *Science* 315 (5814) (2007) 972–976.
- [71] L. Xu, J. Neufeld, B. Larson, D. Schuurmans, Maximum margin clustering, in: *Advances in Neural Information Processing Systems*, 2004, pp. 1537–1544.
- [72] N. X. Vinh, J. Epps, J. Bailey, Information theoretic measures for clusterings comparison: is a correction for chance necessary?, in: *International Conference on Machine Learning*, Montreal, Canada, 2009, pp. 1073–1080.
- [73] W. M. Rand, Objective criteria for the evaluation of clustering methods, *American Statistical association* 66 (336) (1971) 846–850.
- [74] D. Steinley, Properties of the hubert-arable adjusted rand index., *Psychological methods* 9 (3) (2004) 386.
- [75] N. Jardine, C. J. van Rijsbergen, The use of hierarchic clustering in information retrieval, *Information storage and retrieval* (1971) 217–240.
- [76] N. X. Vinh, J. Epps, J. Bailey, Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance, *The Journal of Machine Learning Research* 11 (October) (2010) 2837–2854.
- [77] J. Shotton, M. Johnson, R. Cipolla, Semantic texton forests for image categorization and segmentation, in: *IEEE Conference on Computer Vision and Pattern Recognition*, Anchorage, Alaska, United States, 2008, pp. 1–8.