# Min-Entropy Latent Model for Weakly Supervised Object Detection

Fang Wan, *Student Member, IEEE*, Pengxu Wei, *Student Member, IEEE*, Zhenjun Han, *Member, IEEE*, Jianbin Jiao, *Member, IEEE*, and Qixiang Ye, *Senior Member, IEEE*

**Abstract**—Weakly supervised object detection is a challenging task when provided with image category supervision but required to learn, at the same time, object locations and object detectors. The inconsistency between the weak supervision and learning objectives introduces significant randomness to object locations and ambiguity to detectors. In this paper, a min-entropy latent model (MELM) is proposed for weakly supervised object detection. Min-entropy serves as a model to learn object locations and a metric to measure the randomness of object localization during learning. It aims to principally reduce the variance of learned instances and alleviate the ambiguity of detectors. MELM is decomposed into three components including proposal clique partition, object clique discovery, and object localization. MELM is optimized with a recurrent learning algorithm, which leverages continuation optimization to solve the challenging non-convexity problem. Experiments demonstrate that MELM significantly improves the performance of weakly supervised object detection, weakly supervised object localization, and image classification, against the state-of-the-art approaches.

**Index Terms**—Weakly Supervised Learning, Object Detection, Min-Entropy Latent Model, Recurrent Learning.

✦

## 1 INTRODUCTION

SUPERVISED object detection has made great progress in recent years [1]–[6], as concluded in the object detection survey [7]. This can be attributed to the availability of large datasets with precise object annotations and deep neural networks capable of absorbing the annotation information, especially. Nevertheless, annotating a bounding-box for each object in large datasets is laborious, expensive, or even impractical. It is also not consistent with cognitive learning, which requires solely the presence or absence of a class of objects in a scene, instead of bounding-boxes that indicate the precise locations of all objects.

Weakly supervised learning (WSL) refers to methods that rely on training data with incomplete annotations to learn recognition models. Weakly supervised object detection (WSOD) requires solely the image-level annotations indicating the presence or absence of a class of objects in images to learn detectors [8]–[29]. It can leverage rich Web images with tags to learn object-level models.

To tackle the WSOD problem, existing approaches often resort to latent variable learning or multi-instance learning (MIL) by using redundant object proposals as inputs. The learning objective is designed to choose a true instance from redundant object proposals of each image to minimize the image classification loss. Due to the unavailability of object-level annotations, WSOD approaches require to collect instances from redundant proposals, as well as learning detectors that compromise the appearance of various objects. It typically requires solving a non-convex model and thus is challenged by the local minimum problem.

In the learning procedure of weakly supervised deep detection networks (WSDDN) [22], a representative WSOD approach, the problem has been observed, *i.e.*, the collected instances switch among different object parts with great randomness, Fig. 1. Various object parts were capable of minimizing image classification loss, but experienced difficulty in optimizing object detectors due to their appearance ambiguity. Recent approaches have used image segmentation [28], [30], context information [24], and instance classifier refinement [27] to empirically regularize the learning procedure. However, the issues about principally reducing localization randomness and alleviating the local minimum remain unresolved.

In this paper, we propose a clique-based min-entropy latent model (MELM) [1] to collect instances with minimum randomness, motivated by a classical thermodynamic principle: *Minimizing entropy results in minimum randomness of a system.* Min-entropy is used as a model to learn object locations and a metric to measure the randomness of localization during learning. MELM is concluded as three components: (1) Instance (object and object part) collection with a clique partition module; (2) Object clique discovery with a global min-entropy model; (3) Object localization with a local min-entropy model, Fig. 2. A clique is defined as a set of object proposals which are spatially related (*i.e.*, overlapping with each other) and class related (*i.e.*, having similar object class scores), Fig. 3. The introduction of proposal cliques can facilitate reducing the redundancy of region proposals and optimizing min-entropy models.

With the clique partition module and min-entropy models, we can collect instances with minimum randomness, activate true object extent, and suppress object parts, Fig. 1. MELM is deployed as a clique partition module and network branches concerning object clique discovery and object localization on top of a deep convolutional neural

*F. Wan, Z. Han, J. Jiao, and Q. Ye are with the School of Electronic, Electrical and Communication Engineering, University of Chinese Academy of Sciences (UCAS), Beijing, China, 100049. Emails: wanfang13@mails.ucas.ac.cn, hanzhj@ucas.ac.cn, jiaojb@ucas.ac.cn, qxye@ucas.ac.cn. Pengxu Wei is with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China. Email: weipengxu11@mails.ucas.ac.cn.*

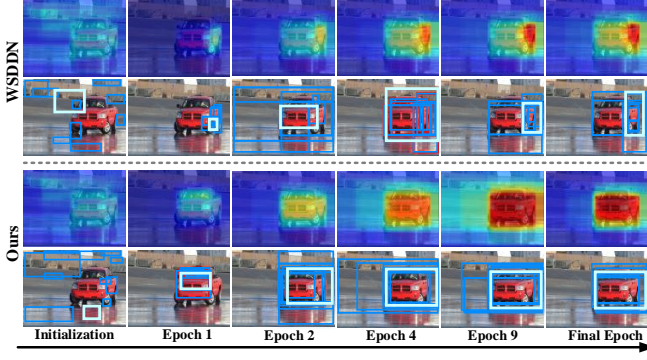1. Source code is available at https://github.com/WinFrand/MELM.

Fig. 1. Evolution of object locations during learning. Blue boxes denote proposals of high object probability and white ones detected objects. It can be seen that our approach reduces localization randomness and learns object extent. (Best viewed in color.)

network (CNN). Based on the global and local min-entropy models, we adopt a recurrent strategy to train detectors and pursue true object extent using solely image-level supervision. This is based on the priori that in deep networks the image classification task and object detection task are highly correlated, which allows MELM to recurrently transfer the weak supervision, *i.e.*, image category annotations, to object locations. By accumulating multiple iterations, MELM discovers multiple objects, if such exist, from a single image.

MELM is first proposed in our CVPR paper [31] and is promoted both theoretically and experimentally in this full version. The contributions of this paper include: (1) A min-entropy latent model that is integrated with deep networks to effectively collect instances and principally minimize the localization randomness during weakly supervised learning. (2) A clique partition module that facilitates instance collection, object extent activation, and object part suppression. (3) A recurrent learning algorithm that formulates image classification and object detection as a predictor and a corrector, respectively, and leverages continuation optimization to solve the challenging non-convexity problem. (4) State-of-the-art performance of weakly supervised detection, localization, and image classification.

The remainder of this paper can be concluded as follows. Related works are described in Section 2 and the proposed method is presented in Section 3. Experimental results are given in Section 4. We conclude this paper in Section 5.

## 2 RELATED WORK

WSOD was often solved with a pipelined approach, *i.e.*, an image was first decomposed into object proposals, with which clustering [14]–[16], latent variable learning [12]–[15], [17] or multiple instance learning [8], [10], [11], [21], [32] was used to perform proposal selection and classifier estimation. With the rise of deep learning, pipelined approaches have been evolving into multiple instance learning (MIL) networks [22]–[25], [27]–[29], [33]–[38].

**Clustering.** Various clustering methods were based on a hypothesis that a class of object instances shape a single compact cluster while the negative instances form multiple diffuse clusters. With such a hypothesis, Wang *et al.* [15],

[16] calculated clusters of object proposals using probabilistic latent Semantic Analysis (pLSA) on positive samples, and employed a voting strategy on these clusters to determine positive sub-categories. Bilen and Song [13], [14] leveraged clustering to initialize latent variables, *i.e.*, object regions, part configurations and sub-categories, and learn object detectors based on the initialization. Clustering is a simple but effective method. The disadvantage lies in that a true positive cluster could incorporate significant noise if the objects are surrounded by clutter backgrounds.

**Latent Variable Learning.** Latent SVM [26] learned object locations and detectors using an Expectation-Maximization-like algorithm. Probabilistic Latent Semantic Analysis [15], [16] learned object locations in a latent space.

Various latent variable methods were required to solve the non-convexity problem. They often got stuck in a poor local minimum during learning, *e.g.*, falsely localizing object parts or backgrounds. To pursue a stronger minimum, object symmetry and class mutual exclusion information [12], Nesterov's smoothing [17], and convex clustering [14] were introduced to the optimization function. These approaches can be regarded as regularization which enforces the appearance similarity among objects.

**Multiple Instance Learning (MIL).** A major approach for tackling WSOD is to formulate it as an MIL problem [8], which treats each training image as a "bag" and iteratively selects high-scored instances from each bag when learning detectors. However, MIL remains puzzled by random poor solutions. The multi-fold MIL [10], [11] used division of a training set and cross validation to reduce the randomness and thereby prevented training from prematurely locking onto erroneous solutions. Hoffman *et al.* [21] trained detectors with weak annotations while transferring representations from extra object classes using full supervision (bounding-box annotation) and joint optimization. To reduce the randomness of positive instances, bag splitting was used during the optimization procedure of MILinear [25].

MIL has been updated to MIL networks [22], [27], where the convolutional filters behave as detectors to activate regions of interest on the deep feature maps [39]–[41]. The beam search [42] was used to localize objects by leveraging spatial distributions and informative patterns captured in the convolutional layers. To alleviate the non-convexity problem, Li *et al.* [23] adopted progressive optimization as regularized loss functions. Tang *et al.* [27] proposed to refine instance classifiers online by propagating instance labels to spatially overlapped instances. Diba *et al.* [28] proposed weakly supervised cascaded convolutional networks (WCCN). It learned to produce a class activation map and then selected the best object locations on the map by minimizing the segmentation loss.

MIL networks [27]–[29] report state-of-the-art performance, but are misled by the inconsistency between data annotations and learning objectives. With image-level annotations, they are capable of learning effective representations for image classification. Without object-level annotation, however, their localization ability is limited. The convolutional filters learned with image-level supervision incorporate redundant patterns, *e.g.*, object parts and backgrounds, which cause localization randomness and model ambiguity.

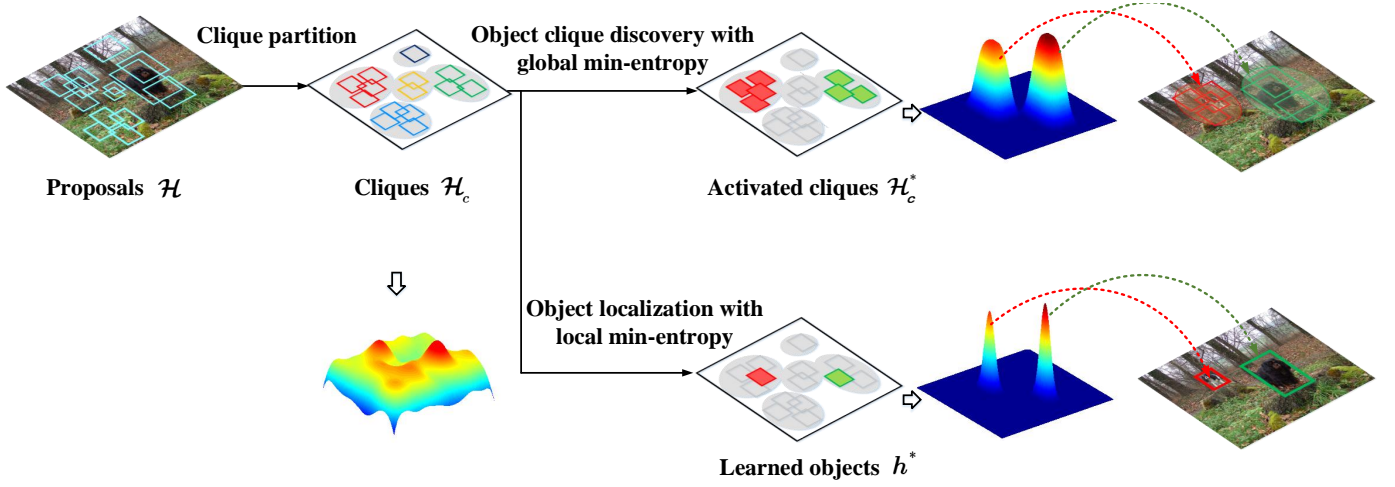Recent methods leveraged online instance classifier re-

Fig. 2. Illustration of the min-entropy latent model (MELM). A clique partition module is proposed to collect objects/parts from redundant proposals; Based on the cliques, a global min-entropy model is defined for object clique discovery; Within discovered cliques, a local min-entropy model is proposed to suppress object parts and select true objects. The three components are iteratively performed.

finement (OICR) [27], [43] and proposal clusters [29], [43] to improve localization. The iterative generation of the proposal clusters [43] with OICR prevented the network from concentrating on parts of objects. In this paper, we propose to solve the localization randomness problem by introducing proposal cliques and min-entropy latent models. Our defined proposal cliques facilitate reducing the redundancy of proposals and optimizing min-entropy models. Using the clique-based min-entropy models, we can learn instances with minimum randomness, activate object extent, and suppress object parts, Fig. 1.

To translate the image labels to object locations, the MIL network approaches [27], [43] defined multiple network branches: the first one for the basic MIL network and the others for instance classifier refinement. We inherit the multi-branch architecture but add recurrent learning to facilitate the object score feedback [44]. With recurrent learning, the network branches can directly benefit from each other.

## 3 METHODOLOGY

### 3.1 Overview

In weakly supervised learning, the inconsistency between the supervision (image-level annotation) and the objective (object-level classifier) introduces significant randomness to object localization and ambiguity to detectors. We aim at reducing this randomness to facilitate the collection of instances. To this end, we analyze two factors that cause such randomness: proposal redundancy and location uncertainty. 1) It is known that the objective functions of WSOD models are typically non-convex [8] and have many local minima. The redundant proposals deteriorate them by introducing more local minima and larger searching space. 2) As the object locations are uncertain, the learned instances may switch among object parts, *i.e.*, local minima.

To reduce the proposal redundancy, we firstly partition the redundant object proposals into cliques and collect instances which are spatially related (*i.e.*, overlapping with each other) and class related (*i.e.*, having similar object class

scores). To minimize localization randomness, we design a global min-entropy model that reflects class and spatial distributions of object cliques. By optimizing the global min-entropy model, discriminative cliques containing objects and object parts are discovered, Fig. 2, and the cliques which lack discriminative information are suppressed. The discovered cliques are used to activate true object extent.

To localize objects in the discovered cliques, a local min-entropy latent model is defined. By optimizing the local min-entropy model pseudo-objects are estimated and their spatial neighbors are estimated as hard negatives. Such pseudo-objects and hard negatives estimated under the min-entropy principle have minimized randomness during learning, and further improve the performance of object localization, Fig. 2. MELM is deployed as a clique partition module and two network branches concerning object clique discovery and object localization, Fig. 4. During learning, it leverages a clique partition module to smooth the objective function and a continuation optimization method to solve the challenging non-convexity problem.

### 3.2 Min-Entropy Latent Model

Let $x \in \mathcal{X}$ denote an image and $y \in \mathcal{Y}$ denote labels indicating if $x$ contains an object or not, where $\mathcal{Y} = \{1, 0\}$. $y = 1$ indicates that there is at least one object of positive class in the image (positive image) while $y = 0$ indicates an image without the object of positive class (negative image). $h$ denoting an object proposal (location) is a latent variable and $\mathcal{H}$ denoting object proposals in an image is the solution space. $\mathcal{H}_c$ denoting proposal clique is a subset of $\mathcal{H}$. $\theta$ denotes the network parameters. The min-entropy latent model (MELM) with object locations $h^*$ and network parameters $\theta^*$ to be learned, is defined as

$$
\begin{aligned}
\{h^*, \theta^*\} &= \underset{h, \theta}{\arg\min} \, E_{(\mathcal{X}, \mathcal{Y})} (h, \theta) \\
&= \underset{h, \theta}{\arg\min} \, E_{(\mathcal{X}, \mathcal{Y})} (\mathcal{H}_c, \theta) + \lambda E_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c)} (h, \theta) \\
&\Leftrightarrow \underset{h, \theta}{\arg\min} \, L_{(\mathcal{X}, \mathcal{Y})} (\mathcal{H}_c, \theta) + \lambda L_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c)} (h, \theta),
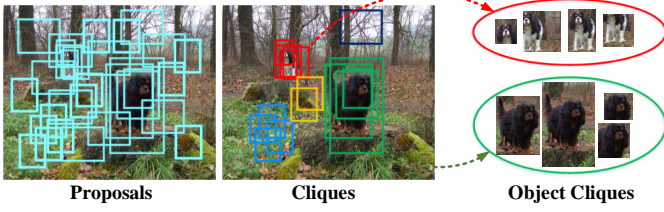\end{aligned}
$$

(1)

**Fig. 3.** The proposals of high scores are selected and dynamically partitioned into same cliques if they are spatially related (*i.e.*, overlapping with each other) and class related (*i.e.*, having similar object class scores). Clique partition targets at collecting object/object parts and activating true object extent.

where $E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta)$ and $E_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h, \theta)$ are the global and local entropy models which respectively serve for object clique discovery and object localization, Fig. 4. $\lambda$ is a regularization weight. $L_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta)$ and $L_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h, \theta)$ are loss functions based on $E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta)$ and $E_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h, \theta)$, respectively.

Given image-level annotations, *i.e.*, the presence or absence of a class of objects in images, the learning objective of MELM is to find a solution that disentangles object instances from noisy object proposals with minimum image classification loss and localization randomness. To this end, MELM is decomposed into three components including clique partition, object clique discovery, and object localization.

### 3.2.1 Clique partition

Noting that the localization randomness usually occurs among high-scored proposals, we empirically select a set of high-scored (top-200) proposals $\tilde{\mathcal{H}}$ to construct the cliques, where $\tilde{\mathcal{H}} \subseteq \mathcal{H}$.

The proposal cliques are the minimum sufficient cover to $\tilde{\mathcal{H}}$ which satisfy the following formulations, as

$$\begin{cases} \bigcup\limits_{c=1}^{C} \mathcal{H}_c = \tilde{\mathcal{H}} \\ \forall c \neq c', \ \mathcal{H}_c \cup \mathcal{H}_{c'} = \emptyset, \end{cases} \quad (2)$$

where $c, c' \in \{1, ..., C\}$ and $C$ is the number of proposal cliques. To partition cliques, the proposals are sorted by their object scores and the following two steps are iteratively performed: 1) Construct a clique using the proposal of highest object score but not belonging to any clique. 2) Find the proposals that overlap with any proposal in the clique larger than a threshold $\tau$ and merge them into the clique.

### 3.2.2 Object clique discovery with global min-entropy

During the learning procedure, it is required that the cliques evolve with minimum randomness. At the same time, it is required to discover discriminative cliques containing objects and object parts. The network parameters fine-tuned with such cliques can activate true object extent. To this end, a global min-entropy model is defined as

$$\begin{aligned} \mathcal{H}_c^* &= \arg\min_{\mathcal{H}_c} E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta) \\ &= \arg\min_{\mathcal{H}_c} -\log \sum_c p(y, \mathcal{H}_c; \theta), \end{aligned} \quad (3)$$

where $p(y, \mathcal{H}_c; \theta)$ is the class probability of a clique $\mathcal{H}_c$ defined on the object score $s(y, h; \theta)$, as

$$p(y, \mathcal{H}_c; \theta) = \frac{\exp\left(1/|\mathcal{H}_c| \sum\limits_{h \in \mathcal{H}_c} s(y, h; \theta)\right)}{\sum\limits_c \sum\limits_y \exp\left(1/|\mathcal{H}_c| \sum\limits_{h \in \mathcal{H}_c} s(y, h; \theta)\right)}, \quad (4)$$

where $|\cdot|$ calculates proposal number in a clique. $s(\cdot)$ denotes the last FC layer in the object clique discovery branch that outputs object scores for proposals.

To ensure that the discovered cliques can best discriminate the positive images from negative ones, we further introduce a classification-related weight $w_{\mathcal{H}_c}$. Based on the prior that the object class probabilities of proposals are correlated with their image class probabilities, the global min-entropy is then defined as

$$E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta) = -\log \sum_c w_{\mathcal{H}_c} p(y, \mathcal{H}_c; \theta), \quad (5)$$

where $w_{\mathcal{H}_c}$, defined as

$$w_{\mathcal{H}_c} = \frac{p(y, \mathcal{H}_c; \theta)}{\sum\limits_y p(y, \mathcal{H}_c; \theta)}, \quad (6)$$

is the classification-related weight of clique $\mathcal{H}_c$. Eq. (5) belongs to the Aczél and Daróczy (AD) entropy [45], [46] family and is derivable. Eq. (6) shows that when $y = 1$, $w_{\mathcal{H}_c} \in [0, 1]$ is positively correlated to object score of the positive class in a clique, but negatively correlated to scores of all other classes.

With above definitions, we implement an object clique discovery branch on top of the network, Fig. 4, and define a loss function to learn network parameters, as

$$\begin{aligned} L_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta) = {}& y E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c, \theta) \\ &- (1 - y) \sum_h \log(1 - p(y, h; \theta)). \end{aligned} \quad (7)$$

For positive images, $y = 1$, the second term is zero and only the global min-entropy term is optimized. For negative images, $y = 0$, the first term is zero and the second term (image classification loss) is optimized.

### 3.2.3 Object localization with local min-entropy

The cliques discovered by the global min-entropy model constitute good initialization for object localization, but nonetheless incorporate random false positives, e.g., object parts and/or partial objects with backgrounds. This is caused by the learning objective of object clique discovery, which selects proposals to discriminate positive images from negative ones but does not consider how to precisely localize objects.

A local min-entropy latent model is then defined to localize objects based on the discovered cliques, as

$$h^* = \arg\min_{h \in \mathcal{H}_c^*} E_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c^*)}(h, \theta), \quad (8)$$

where

$$E_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h, \theta) = -\sum_{h \in \Omega_{h^*}} w_h p(y, h; \theta) \log p(y, h; \theta) \quad (9)$$
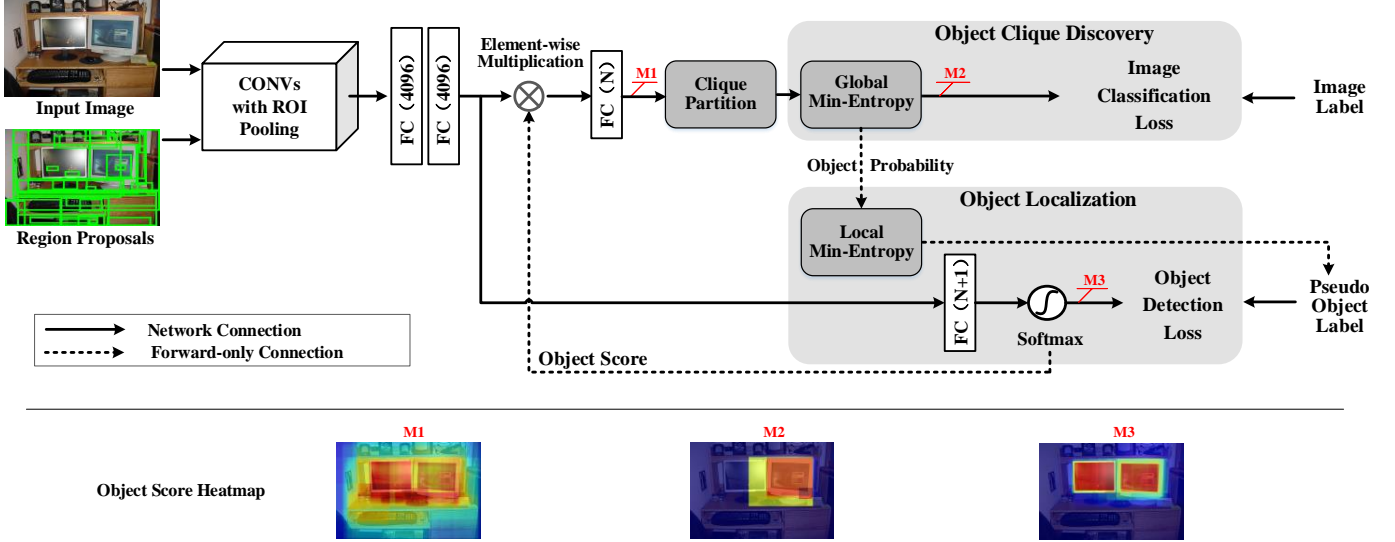
Fig. 4. MELM is deployed as a clique partition module and two network branches for object clique discovery and object localization. These two network branches are unified with feature learning and optimized with a recurrent learning algorithm. "M1", "M2" and "M3" are heatmaps about proposal scores without min-entropy, with global min-entropy, and with local min-entropy, respectively. $N$ is the number of object categories.

also belongs to the AD entropy [45], [46] family and is also derivable. Different from Eq. (5) which considers the sum of the proposal probabilities globally to predict the image labels, Eq. (9) is designed to locally discriminate each proposal to be positive or negative. $w_h$ is defined as

$$w_h = \frac{\sum\limits_{h \in \Omega_{h^*}} g(h, h^*) p(y, h; \theta)}{p(y, h; \theta) \sum\limits_{h \in \Omega_{h^*}} g(h, h^*)}, \tag{10}$$

where $\Omega_{h^*}$ denotes neighborhoods of $h^*$ in the clique. $g(h, h^*) = e^{-a(1-O(h,h^*))^2}$ is a Gaussian kernel function with parameter $a$. $O(h, h^*)$ is the IoU of two proposals. The Gaussian kernel function returns a high value when $O(h, h^*)$ is large, and a low value when $O(h, h^*)$ is small. With Eq. (10), we define a "soft" proposal labeling strategy for object localization, which is validated to be less sensitive to noises [47] compared to the hard thresholding approach defined in [31].

Accordingly, the loss function of the object localization branch is defined as

$$L_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c)}(h, \theta) = E_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c^*)}(h, \theta). \tag{11}$$

According to the definition of $w_h$, the proposals close to $h^*$ tend to be true objects, and those far from $h^*$, i.e., $O(h; h*) < 0.5$, are hard negatives. Optimizing the loss function produces sparse object proposals of high object probability $p(y, h; \theta)$ and suppresses object parts in clique $\mathcal{H}_c^*$. During the learning procedure, the localization capability of detectors is progressively improved.

### 3.3 Model Implementation

MELM is implemented with an integrated deep network, with a clique partition module and two network branches added on top of the FC layers, Fig. 4. The first network branch, designated as the *object clique discovery* branch, has a global min-entropy layer, which defines the distribution

---

**Algorithm 1** Recurrent Learning

**Input:** Image $x \in \mathcal{X}$, image label $y \in \mathcal{Y}$, and object proposals $h \in \mathcal{H}$

**Output:** Network parameters $\theta$
1: Initialize object score $s(h) = s(y, h; \theta) = 1$ for all $h$
2: **for** $i = 1$ **to** $MaxIter$ **do**
3:    $\phi_h \leftarrow$ Compute deep features for all $h$ through forward score
4:    $\phi_h \leftarrow \phi_h \cdot s(h)$ Aggregate features by object score
5:    **Clique partition:**
6:      $\mathcal{H}_c \leftarrow$ Clique partition using Eq. (2)
7:    **Object clique discovery:**
8:      $\mathcal{H}_c^* \leftarrow$ Optimize $E_{(\mathcal{X}, \mathcal{Y})}(\mathcal{H}_c, \theta)$ using Eq. (5)
9:      $L_{(\mathcal{X}, \mathcal{Y})}(\mathcal{H}_c, \theta) \leftarrow$ Compute using Eq. (7)
10:   **Object localization:**
11:     $h^* \leftarrow$ Optimize $E_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c^*)}(h, \theta)$ using Eq. (8)
12:     $L_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c)}(h, \theta) \leftarrow$ Compute using Eq. (11)
13:   **Network parameter update:**
14:     $\theta \leftarrow$ Back-propagate by miniminzing Eq. (7) and Eq. (11)
15:   $s(h) \leftarrow$ Update object score using parameters $\theta$
16: **end for**

---

of object probability and targets at finding candidate object cliques by optimizing the global entropy and the image classification loss. The second branch, designated as the *object localization* branch, has a local min-entropy layer and a soft-max layer. The local min-entropy layer classifies the object candidates in a clique into pseudo objects[2] and hard negatives by optimizing the local entropy and pseudo object detection loss.

In the learning phase, object proposals are firstly generated for each image. An ROI-pooling layer atop the convolutional layer (CONV5) is used for efficient feature extraction for these proposals. The MELMs are optimized with a recurrently learning algorithm, which uses forward propagation to select sparse proposals as object instances, and back-

---

2. Pseudo objects are the instantaneously learned objects.

propagation to optimize the network parameters with the gradient defined in Appendix. The object probability of each proposal is recurrently aggregated by multiplying by the object probability learned in the preceding iteration. In the detection phase, the learned object detectors, *i.e.*, the parameters for the soft-max and FC layers, are used to classify proposals and localize objects.

### 3.4 Model Learning

The objective of model learning is transferring the image category supervision to object locations with min-entropy constraints, *i.e.*, minimum localization randomness.

**Recurrent Learning.** A recurrent learning algorithm is implemented to transfer the image-level (weak) supervision using an integrated forward- and back-propagation procedure, Fig. 5(a). In a feed-forward procedure, the min-entropy latent models discover object cliques and localize objects which are used as pseudo-objects for detector learning. With the learned detectors the object localization branch assigns all proposals new object probability, which is used to aggregate the object scores with an element-wise multiply operation in the next learning iteration. In the back-propagation procedure, the object clique discovery and object localization branches are jointly optimized with an SGD algorithm, which propagates gradients generated with image classification loss and pseudo-object detection loss. With forward- and back-propagation procedures, the network parameters are updated and the image classifiers and object detectors are mutually enforced. The recurrent learning algorithm is described in Alg. 1.

**Accumulated Recurrent Learning.** Fig. 5(b) shows the proposed accumulated recurrent learning (ARL). In ARL, we add multiple object localization branches, which may localize objects different from those discovered by previous branches. We thus accumulates objects from all previous branches. Doing so not only endows this approach the capability to localize multiple objects in a single image but also improves the robustness about object appearance diversity by learning various objects with multiple detectors.

### 3.5 Model analysis

With the clique partition module and recurrent learning, MELM implements the idea of continuation optimization [48] to alleviate the non-convexity problem.

In continuation optimization, a complex non-convex objective function is denoted as $E(\theta)$, where $\theta$ denotes the model parameters. Optimizing $E(\theta)$ is to find the solution

$$\theta^* = \arg\min_{\theta} E(\theta). \tag{12}$$

While directly optimizing Eq. (12) causes local minimum solutions, a smoothed function $E(\theta, \lambda)$ is introduced to approximate $E(\theta)$ and facilitate the optimization, as

$$E(\theta, \lambda) = E(\theta) - \lambda \mathcal{E}(\theta), \tag{13}$$

where $\lambda \in [0, 1]$ controls the smoothness of the approximate function $E(\theta, \lambda)$ and $\mathcal{E}(\theta)$ is a correction function. The traditional continuation method traces an implicitly defined curve from a starting point $(\theta^0, 1)$ to a solution point $(\theta^*, 0)$, where $\theta^0$ is the solution of $E(\theta, \lambda)$ when $\lambda=1$.



Fig. 5. Flowchart of (a) the recurrent learning algorithm and (b) unfolded accumulated recurrent learning algorithm. The solid lines denote network connections and dotted lines denote forward-only connections.

During the procedure, if $E(\theta, \lambda)$ is smooth and its solution is close to $E(\theta)$, we need only to fill the gap between them. This is done by defining a consequence of predictions and corrections to iteratively approximate the original objective function and approach the globally optimal solution $\theta^*$.

The objective function of MELM, defined in Eq. (1), is to find the solution $\{h^*, \theta^*\}$,

$$\{h^*, \theta^*\} = \arg\min_{h, \theta} E_{(\mathcal{X}, \mathcal{Y})}(h, \theta). \tag{14}$$

For the complexity and non-convexity of $E_{(\mathcal{X}, \mathcal{Y})}(h, \theta)$, we propose to optimize an approximate function,

$$E_{(\mathcal{X}, \mathcal{Y})}(\mathcal{H}_c, \theta) = E_{(\mathcal{X}, \mathcal{Y})}(h, \theta) - \lambda E_{(\mathcal{X}, \mathcal{Y}, \mathcal{H}_c)}(h, \theta), \tag{15}$$

which corresponds to Eq. (1). $E_{(\mathcal{X}, \mathcal{Y})}(\mathcal{H}_c, \theta)$ is defined by the clique partition module and is smoother than $E_{(\mathcal{X}, \mathcal{Y})}(h, \theta)$. This is achieved by reducing the solution space from thousands of proposals to tens of cliques in each image and averaging the class probability of all proposals in each clique, as defined by Eq. (4).

With the approximate function defined, we explore recurrent predictions and corrections to optimize the model. The gap between $E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c,\theta)$ and $E_{(\mathcal{X},\mathcal{Y})}(h,\theta)$ is that the former is defined to discover object cliques but the latter to localize objects. As the solution of $E_{(\mathcal{X},\mathcal{Y})}(h,\theta)$ (object) is included in the solution of $E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c,\theta)$ (clique), the gap can be simply filled by designing a correction model $E_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h,\theta)$ to localize the object in the clique. With recurrent learning, the original objective function is thus progressively approximated.

Accordingly, the weakly supervised learning problem is decomposed into an object clique discovery problem (prediction) and object localization problem (correction). The non-convex optimization problem is turned into a proximate problem, which is easier to be optimized [49], [50].

## 3.6 Object Detection

By optimizing the min-entropy latent models, we obtain object detectors, which are applied to detect objects from test images. The detection procedure involves feature extraction and object localization Fig. 4. With redundant object proposals extracted by the Selective Search [51] or the EdgeBox method [52], a test image is fed to the feature extraction module, and then a ROI-pooling layer is used to extract features for each proposal. The detector outputs object scores for each proposal and a Non-Maximum Suppression (NMS) procedure is used to remove the overlapped proposals.

## 4 EXPERIMENTS

The PASCAL VOC 2007, 2010, 2012 datasets [53], the ILSVRC 2013 dataset [54], and the MSCOCO 2014 dataset [55] are used to evaluate the proposed approach. In what follows, the datasets and experimental settings are first described. The evaluation of the model and comparison with the state-of-the-art approaches are then presented.

## 4.1 Experimental Settings

**Datasets.** The VOC datasets have 20 object categories. The VOC 2007 datasets contains 9963 images which are divided into three subsets: 5011 for $train$ and $val$, and 4952 for $test$. The VOC 2010 dataset contains 19740 images of which 10103 for $train$ and $val$, and 9637 for $test$. The VOC 2012 dataset contains 22531 images which are divided into three subsets: 11540 for $train$ and $val$, and 10991 for $test$. The ILSVRC 2013 detection dataset is more challenging for object detection as it has 200 object categories, containing 464278 images where 424126 image for $train$ and $val$, and 40152 images for $test$. For comparison with the previous works, we split the $val$ set of ILSVRC 2013 detection dataset into $val1$ and $val2$ as in [1], which was used for training and test, respectively. Although it has more training images, the number of images for each object category is much less than that in the VOC datasets. The MSCOCO 2014 dataset contains 80 object categories, with challenging aspects including multiple objects, multiple classes, and small objects. On the PASCAL VOC and ILSVRC 2013 datasets the mean average precision (mAP) is used for evaluation. On the MSCOCO 2014 dataset the mAP under multiple IoUs is used.

**CNN Models.** MELM is implemented with two popular CNN models pre-trained on the ImageNet ILSVRC 2012 dataset. The first CNN model VGG-CNN-F (VGGF for short) [56] has a similar architechture as the AlexNet [57] which has 5 convolutional layers and 3 fully connected layers. The second CNN model is VGG16 [58], which has 13 convolutional layers and 3 fully connected layers. For these two CNN models, we replaced the spatial pooling layer after the last convolution layer with the ROI-pooling layer as [2]. The FC8 layer in the two CNN models was removed and the MELM model was added.

**Object Proposals.** The Selective Search [51] or Edge-Boxes method [52] was used to extract about 2000 object proposals for each image. As the conventional object detection task, we used the fast setting when generating proposals by Selective Search. We also removed the proposals whose width or height are less than 20 pixels.

**Learning settings.** Following [22], [24], [27], [28], the input images were re-sized into 5 scales {480, 576, 688, 864, 1200} with respect to the larger side, height or width. The scale of a training image was randomly selected and each image was randomly flipped. In this way, each test image was augmented into 10 images. For recurrent learning, we employed the SGD algorithm with momentum 0.9, weight decay 5e-4, and batch size 1. The model iterated 20 epochs where the learning rate was 5e-3 for the first 15 epochs and 5e-4 for the last 5 epochs. The output scores of each proposal from the 10 augmented images were averaged.

## 4.2 Model Effect and Analysis

### 4.2.1 Clique Affect

Fig. 6 shows that in discovered cliques discriminative objects and object parts were collected and the proposals which lack discriminative information were suppressed. With the proposals about objects and object parts, the global min-entropy model could activate object extent during the back-propagation procedure. It can also be seen that the true object in a clique can be precisely localized after the recurrent learning procedure.

Fig. 7 shows the object cliques from different learning epochs. It can be seen that in the early training stage (Epoch 2), the object clique collected the object extent, $i.e$, object and object parts. This ensured the object extent activation by the object clique discovery branch. The object localization branch further suppressed the object parts in the object clique (Epoch 4). MELM finally activated the true object extent, suppressed the object part and detected objects accurately (Epoch 20).

### 4.2.2 Randomness Analysis

Fig. 8a shows the evolution of global and local entropy, suggesting that our approach optimizes the min-entropy objective during learning. Fig. 8b provides the gradient evolution of the FC layers. In the early learning epochs, the gradient of the global min-entropy module was slightly larger than that of the local min-entropy module, suggesting that the network focused on optimizing the image classifiers. As learning proceeded, the gradient of the global min-entropy module decreased such that the local min-entropy

|  |  |  |  |  |  |
|---|---|---|---|---|---|
| **Cliques** | **Object Cliques** | **Localized Objects** | **Cliques** | **Object Cliques** | **Localized Objects** |

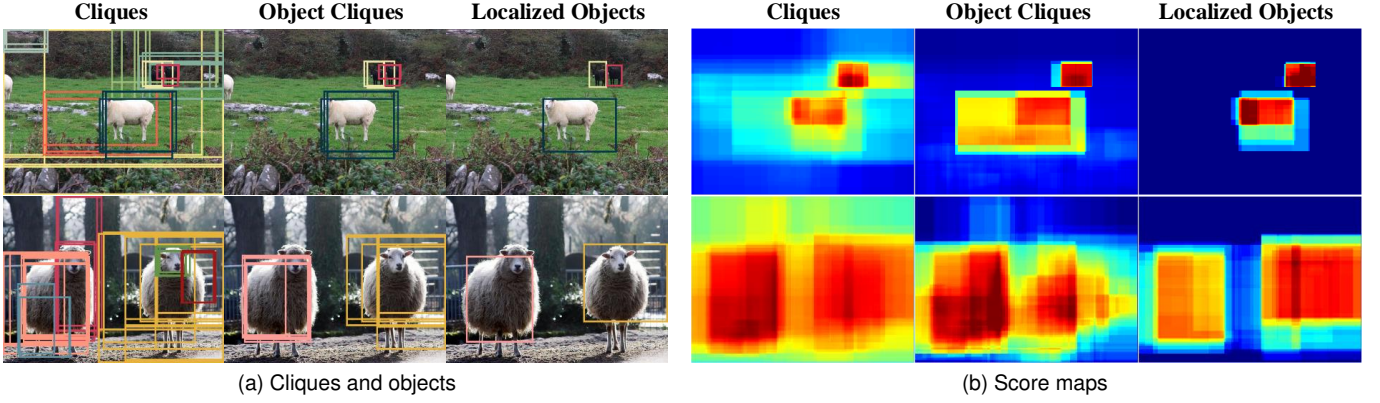(a) Cliques and objects                    (b) Score maps

Fig. 6. Visualization of the clique partition, object clique discovery, and object localization results. (a) Bounding boxes of different colors denote proposals from different cliques. (b) Score maps of cliques and objects. (Best viewed in color)
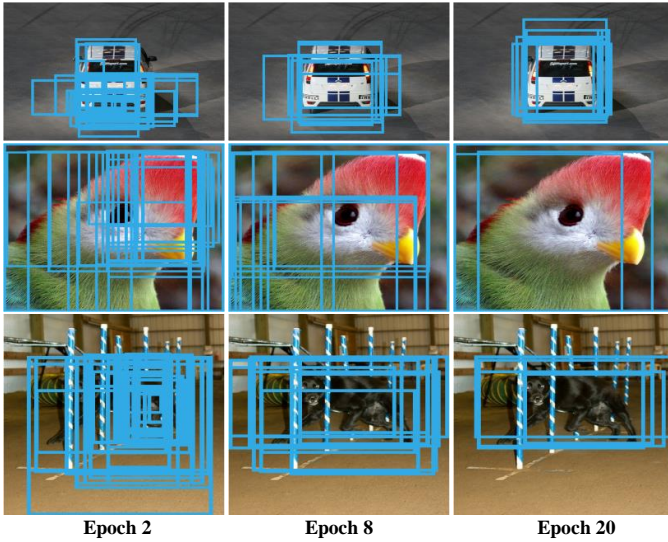


**Epoch 2**          **Epoch 8**          **Epoch 20**

Fig. 7. Evolution of cliques. (Best viewed in color).



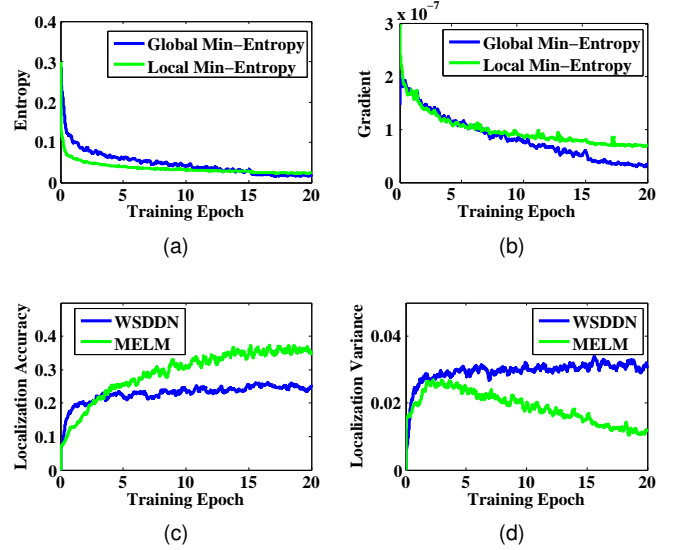(a)          (b)

(c)          (d)

Fig. 8. Gradient, entropy, and localization on the PASCAL VOC 2007 $trainval$ set. (a) The evolution of entropy. (b) The evolution of gradient. (c) Localization accuracy. (d) Localization variance.

module dominated the training of the network, indicating that the object detectors were being optimized.

To evaluate the effect of min-entropy, the randomness of object locations was evaluated with localization accuracy and localization variance. Localization accuracy was calculated by weighted averaging the overlaps between the ground-truth object boxes and the learned object boxes, by using $p(y, h; \theta)$ as the weight. Localization variance was defined as the weighted variance of the overlaps by using $p(y, h; \theta)$ as the weight. Fig. 8c and Fig. 8d show that the proposed MELM had significantly greater localization accuracy and lower localization variance than WSDDN. This strongly indicates that our approach effectively reduces localization randomness during weakly supervised learning.

Such an effect was further illustrated in Fig. 9, where we compared WSDDN with MELM by the localization accuracy and localization variance during the learning. As shown in Fig. 9, MELM significantly reduced the localization randomness and achieved higher localization accuracy than WSDDN. Take the "bicycle" in Fig. 9 for example. In the

early training epochs, both WSDDN and MELM failed to localize the objects. In the following training epochs MELM reduced the randomness and achieved high localization accuracy. In contrast, WSDDN switched among object parts and failed to localize the true objects.

### 4.2.3 Ablation Experiments

**Baseline.** The baseline approach was derived by simplifying Eq. (7) to solely model the global entropy $E_{(\mathcal{X}, \mathcal{Y})}(\mathcal{H}_c, \theta)$. This is similar to WSDDN without the spatial regulariser [22] where the single learning objective is to minimize the image classification loss. This baseline, referred to as "MELM-base" in Table 1, achieved 31.5% mAP using the VGGF network.

**Clique Effect.** By dividing the object proposals into cliques, the "MELM-base" approach was promoted to "MELM-base+Clique". Table 1 shows that the introduction of proposal cliques improved the detection performance by

Fig. 9. Comparison of the learned object locations by WSDDN [22] and MELM. The yellow boxes in the first column denote ground-truth objects. The white boxes denote the learned object locations and the blue boxes denote the high-scored proposals. It can be seen that for WSDDN the learned object locations evolved with large randomness, *i.e.*, switch among the proposals around the objects. In contrast the object locations learned by MELM are consistent and have small randomness, which is quantified by the localization variance curves in the last column. (Best viewed in color)

TABLE 1
Detection mean average precision (%) on the PASCAL VOC 2007 test set. Ablation experimental results of MELM.

| CNN | Method | mAP |
|---|---|---|
| VGGF | MELM-base | 31.5 |
| | MELM-base+Clique | 33.9 |
| | MELM-D | 33.6 |
| | MELM-L | 36.0 |
| | MELM-D+RL | 34.1 |
| | MELM-L+RL | **38.4** |
| VGG16 | MELM-base+Clique | 29.5 |
| | MELM-D | 32.6 |
| | MELM-L | 40.1 |
| | MELM-D+RL | 34.5 |
| | MELM-L+RL | 42.6 |
| | MELM-D+ARL | 37.4 |
| | MELM-L1+ARL | 46.4 |
| | MELM-L2+ARL | **47.3** |

TABLE 2
Detection mean average precision (%) on the PASCAL VOC 2007 *val* set. Performance with different clique sizes (controlled by $\tau$) of MELM.

| $\tau$ | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 1 |
|---|---|---|---|---|---|---|
| mAP | 32.6 | 34.3 | 34.4 | 35.3 | 33.5 | 34.4 |

without using the recurrent learning. Table 1 shows that with VGGF we achieved 33.6% and 36.0% mAP for object clique discovery and object localization branches, which improved the baseline "MELM-base" by 2.1% and 5.5%. For VGG16, "MELM-L" significantly improved the "MELM-base+Clique" from 29.5% to 40.1%, with a 10.6% margin at most. This fully demonstrated that the min-entropy models and their implementation with object clique discovery and object localization branches were pillars of our approach.

**Recurrent learning.** In Table 1, the recurrent learning algorithms "MELM-D+RL" and "MELM-L+RL", respectively achieved 34.5% and 42.6% mAP, improving the "MELM-L" (without recurrent learning) by 0.5% and 2.4%. When using VGG16, "MELM-D+RL" and "MELM-L+RL" respectively achieved 34.5% and 42.6% mAP, improving the "MELM-L" by 1.9% and 2.5%. These improvements showed that with recurrent learning, Fig. 4, the object clique discovery and object localization branches benefited from each other and thus were mutually enforced.

**Accumulated recurrent learning.** The models with accumulated recurrent learning were denoted by "MELM-D+ARL", "MELM-L1+ARL", and "MELM-L2+ARL" in Table 1. In the learning procedure, the high scored proposals were accumulated into the next branch. When using two object localization branches, "MELM-L2-ARL" significantly

2.4% (from 31.5% to 33.9%). That occurred because using partitioned cliques reduced the solution space of the latent variable learning, thus readily reducing the redundancy of object proposals and facilitating a better solution. We also conducted experiments with different $\tau$ values, which controls the clique size as defined in Sec. 3.2.1, and summarized the results in Table. 2. Accordingly, we empirically set $\tau$ to be 0.7 in other experiments.

**Min-entropy models.** We denoted the min-entropy models by "MELM-D" and "MELM-L" in Table 1, which respectively corresponded to object clique discovery and object localization. We trained the models by simply cascading the object clique discovery and object localization branches,

TABLE 3
Detection mean average precision (%) on the PASCAL VOC 2007 test set. Comparison of MELM to the state-of-the-arts.

| CNN | Method | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv | mAP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| VGGF/ AlexNet | MILinear [25] | 41.3 | 39.7 | 22.1 | 9.5 | 3.9 | 41.0 | 45.0 | 19.1 | 1.0 | 34.0 | 16.0 | 21.3 | 32.5 | 43.4 | **21.9** | 19.7 | 21.5 | 22.3 | 36.0 | 18.0 | 25.4 |
| | Multi-fold MIL [11] | 39.3 | 43.0 | 28.8 | 20.4 | 8.0 | 45.5 | 47.9 | 22.1 | 8.4 | 33.5 | 23.6 | 29.2 | 38.5 | 47.9 | 20.3 | 20.0 | 35.8 | 30.8 | 41.0 | 20.1 | 30.2 |
| | PDA [23] | 49.7 | 33.6 | 30.8 | 19.9 | 13.0 | 40.5 | 54.3 | 37.4 | **14.8** | 39.8 | 9.4 | 28.8 | 38.1 | 49.8 | 14.5 | **24.0** | 27.1 | 12.1 | 42.3 | 39.7 | 31.0 |
| | LCL+Context [16] | 48.9 | 42.3 | 26.1 | 11.3 | 11.9 | 41.3 | 40.9 | 34.7 | 10.8 | 34.7 | 18.8 | 34.4 | 35.4 | 52.7 | 19.1 | 17.4 | 35.9 | 33.3 | 34.8 | 46.5 | 31.6 |
| | WSDDN [22] | 42.9 | 56.0 | 32.0 | 17.6 | 10.2 | 61.8 | 50.2 | 29.0 | 3.8 | 36.2 | 18.5 | 31.1 | 45.8 | 54.5 | 10.2 | 15.4 | 36.3 | 45.2 | 50.1 | 43.8 | 34.5 |
| | ContextNet [24] | **57.1** | 52.0 | 31.5 | 7.6 | 11.5 | 55.0 | 53.1 | 34.1 | 1.7 | 33.1 | **49.2** | **42.0** | 47.3 | 56.6 | 15.3 | 12.8 | 24.8 | **48.9** | 44.4 | 47.8 | 36.3 |
| | WCCN [28] | 43.9 | **57.6** | **34.9** | 21.3 | 14.7 | 64.7 | 52.8 | 34.2 | 6.5 | 41.2 | 20.5 | 33.8 | 47.6 | 56.8 | 12.7 | 18.8 | **39.6** | 46.9 | **52.9** | 45.1 | 37.3 |
| | OICR [27] | 53.1 | 57.1 | 32.4 | 12.3 | 15.8 | 58.2 | 56.7 | **39.6** | 0.9 | 44.8 | 39.9 | 31.0 | **54.0** | **62.4** | 4.5 | 20.6 | 39.2 | 38.1 | 48.9 | 48.6 | 37.9 |
| | MELM | 56.4 | 54.7 | 30.9 | 21.1 | **17.3** | 52.8 | **60.0** | 36.1 | 3.9 | **47.8** | 35.5 | 28.9 | 30.9 | 61.0 | 5.8 | 22.8 | 38.8 | 39.6 | 42.1 | **54.8** | **38.4** |
| VGG16 | WSDDN [22] | 39.4 | 50.1 | 31.5 | 16.3 | 12.6 | 64.5 | 42.8 | 42.6 | 10.1 | 35.7 | 24.9 | 38.2 | 34.4 | 55.6 | 9.4 | 14.7 | 30.2 | 40.7 | 54.7 | 46.9 | 34.8 |
| | PDA [23] | 54.5 | 47.4 | **41.3** | 20.8 | **17.7** | 51.9 | 63.5 | 46.1 | 21.8 | 57.1 | 22.1 | 34.4 | 50.5 | 61.8 | 16.2 | 29.9 | 40.7 | 15.9 | 55.3 | 40.2 | 39.5 |
| | OICR [27] | **58.0** | 62.4 | 31.1 | 19.4 | 13.0 | 65.1 | 62.2 | 28.4 | 24.8 | 44.7 | 30.6 | 25.3 | 37.8 | 65.5 | 15.7 | 24.1 | 41.7 | 46.9 | **64.3** | **62.6** | 41.2 |
| | Self-Taught [29] | 52.2 | 47.1 | 35.0 | 26.7 | 15.4 | 61.3 | 66.0 | **54.3** | 3.0 | 53.6 | 24.7 | 43.6 | 48.4 | 65.8 | 6.6 | 18.8 | 51.9 | 43.6 | 53.6 | 62.4 | 41.7 |
| | WCCN [28] | 49.5 | 60.6 | 38.6 | **29.2** | 16.2 | **70.8** | 56.9 | 42.5 | 10.9 | 44.1 | 29.9 | 42.2 | 47.9 | 64.1 | 13.8 | 23.5 | 45.9 | 54.1 | 60.8 | 54.5 | 42.8 |
| | TS²C [59] | 59.3 | 57.5 | 43.7 | 27.3 | 13.5 | 63.9 | 61.7 | 59.9 | 24.1 | 46.9 | 36.7 | 45.6 | 39.9 | 62.6 | 10.3 | 23.6 | 41.7 | 52.4 | 58.7 | 56.6 | 44.3 |
| | WeakRPN [60] | 57.9 | 70.5 | 37.8 | 5.7 | 21.0 | 66.1 | 69.2 | 59.4 | 3.4 | 57.1 | 57.3 | 35.2 | 64.2 | 68.6 | 32.8 | 28.6 | 50.8 | 49.5 | 41.1 | 30.0 | 45.3 |
| | MELM | 55.6 | **66.9** | 34.2 | 29.1 | 16.4 | 68.8 | **68.1** | 43.0 | **25.0** | **65.6** | **45.3** | **53.2** | 49.6 | **68.6** | 2.0 | 25.4 | **52.5** | **56.8** | 62.1 | 57.1 | **47.3** |
| Ens. | OICR-Ens. [27] | 58.5 | 63.0 | 35.1 | 16.9 | 17.4 | 63.2 | 60.8 | 34.4 | 8.2 | 49.7 | 41.0 | 31.3 | 51.9 | 64.8 | **13.6** | 23.1 | 41.6 | 48.4 | 58.9 | 58.7 | 42.0 |
| | MELM-Ens. | **60.3** | 65.0 | 39.5 | 29.0 | 17.5 | 66.1 | 66.4 | 44.8 | 18.6 | 59.0 | 48.4 | 53.2 | 53.0 | 67.2 | 11.0 | **26.5** | 50.0 | 55.7 | 63.1 | 62.4 | 47.8 |

TABLE 4
Correct localization rate (%) on the PASCAL VOC 2007 *trainval* set. Comparison of MELM to the state-of-the-arts.

| CNN | Method | mAP |
|---|---|---|
| VGGF/ AlexNet | MILinear [25] | 43.9 |
| | LCL+Context [16] | 48.5 |
| | PDA [23] | 49.8 |
| | WCCN [28] | 52.6 |
| | Multi-fold MIL [11] | 54.2 |
| | WSDDN [22] | 54.2 |
| | ContextNet [24] | 55.1 |
| | MELM | **58.4** |
| VGG16 | PDA [23] | 52.4 |
| | WSDDN [22] | 53.5 |
| | WCCN [28] | 56.7 |
| | MELM | **61.4** |

TABLE 5
Detection mean average precision (%) on the PASCAL VOC 2010, 2012, and the ILSVRC 2013 datasets. Comparison of MELM to the state-of-the-arts.

| Dataset | CNN | Method | Dataset Splitting | mAP |
|---|---|---|---|---|
| PASCAL VOC 2010 | VGGF/ AlexNet | PDA [23] | train/val | 21.4 |
| | | WCCN [28] | trainval/test | 28.8 |
| | | MELM | train/val | **35.6** |
| | | MELM | trainval/test | **36.3** |
| | VGG16 | PDA [23] | train/val | 30.7 |
| | | WCCN [28] | trainval/test | 39.5 |
| | | MELM | train/val | **37.1** |
| | | MELM | trainval/test | **39.9** |
| PASCAL VOC 2012 | VGGF/ AlexNet | PDA [23] | train/val | 22.4 |
| | | MILinear [25] | train/val | 23.8 |
| | | WCCN [28] | trainval/test | 28.4 |
| | | ContextNet [24] | trainval/test | 35.3 |
| | | OICR-VGGM [27] | trainval/test | 34.6 |
| | | MELM | train/val | **36.2** |
| | | MELM | trainval/test | **36.4** |
| | VGG16 | PDA [23] | train/val | 29.1 |
| | | Self-Taught [29] | train/val | 39.0 |
| | | WCCN [28] | trainval/test | 37.9 |
| | | OICR [27] | trainval/test | 37.9 |
| | | Self-Taught [29] | trainval/test | 38.3 |
| | | TS²C [59] | trainval/test | 40.0 |
| | | MELM | train/val | **40.2** |
| | | MELM | trainval/test | **42.4** |
| ILSVRC 2013 | VGGF/ AlexNet | MILinear [25] | - | 9.6 |
| | | PDA [23] | val1/val2 | 7.7 |
| | | WCCN [28] | - | 9.8 |
| | | MELM | val1/val2 | **13.4** |

improved the mAP of "MELM-L-RL" from 42.6% to 46.4% (+3.8%). It further improved the mAP from 46.4% to 47.3% (+0.9%) when using three branches, but did not significantly improve when using four.

## 4.3 Performance and Comparison

### 4.3.1 PASCAL VOC datasets

**Weakly Supervised Object Detection.** Table 3 compared the detection performance of MELM with the state-of-the-art approaches on the PASCAL VOC 2007 dataset. It can be seen that MELM respectively achieved 38.4% and 47.3% with the VGGF and VGG16 models. With the popular VGG16 model, MELM respectively outperformed the OICR [27], Self-Taught [29], WCCN [28], WeakRPN [60], and TS²C [59] by 6.1% (47.3% vs. 41.2%), 5.6% (47.3% vs. 41.7%), 4.5% (47.3% vs. 42.8%), 3.0% (47.3% vs. 44.3%) and 2.0% (47.3% vs. 45.3%), which were significant margins in terms of the challenging WSOD task. MELM using multiple networks (MELM-Ens.) outperformed OICR-Ens. (47.8% mAP vs.

42.0% mAP). To further improve the detection performance, we re-trained a Fast-RCNN detector using learned pseudo objects and a ResNet-101 network, and achieved 49.0% mAP.

Table 5 compared the detection performance of MELM with the state-of the-art approaches on the VOC 2010 and

TABLE 6
Detection and localization performance (%) on MSCOCO 2014.
Comparison of MELM to the state-of-the-arts.

| Image Classification | | | | | | | |
|---|---|---|---|---|---|---|---|
| Method | mAP | F1-C | P-C | R-C | F1-O | P-O | R-O |
| CAM [61] | 54.4 | - | - | - | - | - | - |
| SPN [37] | 56.0 | - | - | - | - | - | - |
| ResNet-101 [62] | 75.2 | 69.5 | 80.8 | 63.4 | 74.4 | 82.2 | 68.0 |
| MELM-VGG16 | 79.1 | 72.0 | 79.3 | 68.6 | 76.8 | 82.5 | 71.9 |

| Pointing Localization (with class prediction) | | | | | |
|---|---|---|---|---|---|
| Method | WeakSup [34] | Pronet [63] | DFM [42] | SPN [37] | MELM |
| mAP | 41.2 | 43.5 | 49.2 | 55.3 | 65.1 |

| Object Detection | | | |
|---|---|---|---|
| Method | CNN | mAP@.5 | mAP@[.5,.95] |
| WSDDN [22] | VGGF | 10.1 | 3.1 |
| MELM | VGGF | 11.9 | 4.1 |
| | VGG16 | 18.8 | 7.8 |

TABLE 7
Image classification mAP (%) on the PASCAL VOC 2007 *test* set.
Comparison of MELM to the state-of-the-arts.

| CNN | Method | mAP |
|---|---|---|
| VGGF/ AlexNet | MILinear [25] | 72.0 |
| | AlexNet [57] | 82.4 |
| | WSDDN [22] | 85.3 |
| | WCCN [28] | **87.8** |
| | MELM | **87.8** |
| VGG16 | VGG16 [58] | 89.3 |
| | WSDDN [22] | 89.7 |
| | WCCN [28] | 90.9 |
| | MELM | **93.1** |

VOC 2012 datasets. It can be seen that MELM usually outperformed the state-of-the-art approaches. On the VOC 2010 dataset, MELM with VGGF significantly outperformed WCCN [28] by 7.5% (36.3% vs. 28.8%) with a VGGF model, and was comparable to it with a VGG16 model. On the VOC2012 dataset, with a VGGF model, MELM respectively outperformed WCCN [28] and OICR [27] by 8.0% ( 36.4% vs. 28.4%) and 1.8% (36.4% vs. 34.6%). With a VGG16 model, MELM respectively outperformed WCCN [28], Self-Taught [29], OICR [27], and TS$^2$C [59] by 4.5% (42.4% vs. 37.9%), 4.1% (42.4% vs. 38.3%), 4.5% (42.4% vs. 37.9%) and 2.4% (42.4% vs. 40.0%).

Specifically, the detection performance for "bicycle" (+4.5%), "cow" (+8.5%), "dining-table" (+14.7%), "dog" (+9.6%) significantly improved, which shows the general effectiveness of MELM

Despite of the average good performance, our approach failed on the "person" class, as shown in the last image of Fig. 10(a). "Person" is one of the most challenging class as people often involve great appearance variance from clothes, poses, and occlusions. Furthermore, the definition for ??person?? is not consistent. A "person" could be defined as a pedestrian, a head-and-shoulder, or just a human face. Given such ambiguous definition, what the algorithm can do is to localize the most discriminative part of a "person", e.g., the face. We also note that although the performance of "person" decreased, the average performance for all class significantly increased.

For the object classes with large appearance variance, we observed that the algorithm correctly classified the object regions but often failed to precisely localize them, i.e., the IoU between the learned bounding boxes and the groundtruth is smaller than 0.5. When using the "pointing localization" metric [37], the "person" class achieved 97.1% localization accuracy, which shows potential to practical applications.

Fig. 10 shows some of the detection examples. It can be seen that MELM precisely localize objects from clutter background and correctly localized multiple object regions in a single image.

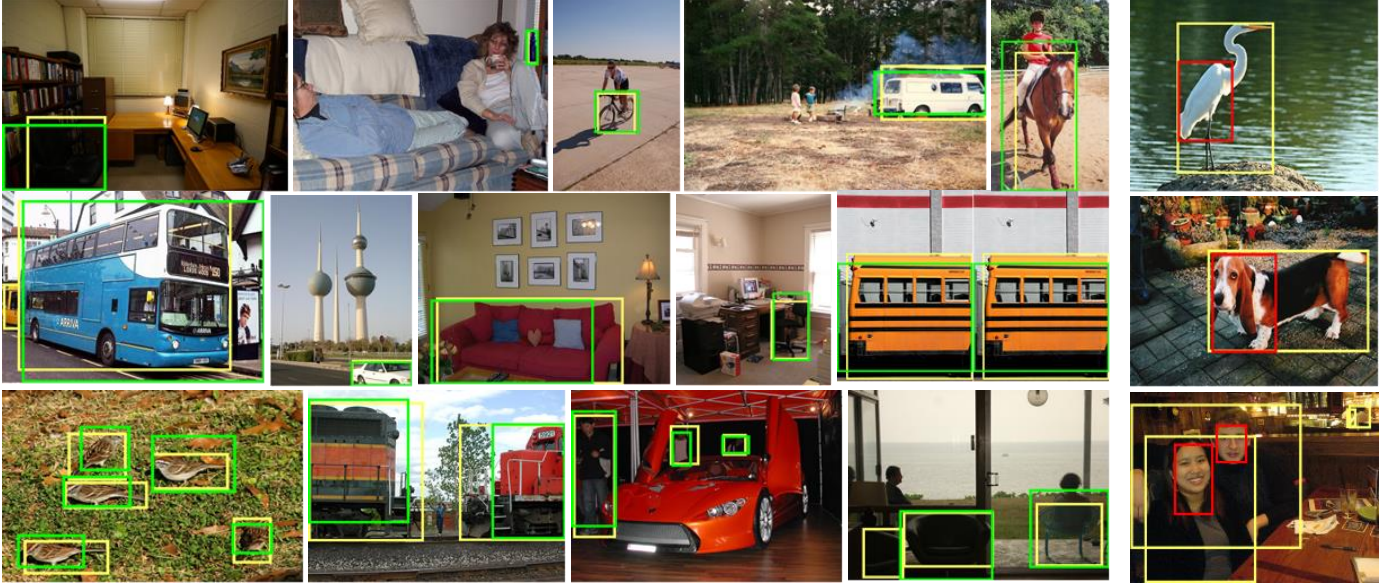**Weakly Supervised Object Localization.** The Correct Localization (CorLoc) metric [18] was employed to evaluate the localization accuracy. CorLoc is the percentage of images for which the region of highest object score has at least 0.5 interaction-over-union (IoU) with the ground-truth object region. This experiment was done on the *trainval* set because the region selection exclusively worked in the training process.

It can be seen in Table 4 that with VGGF model, the mean CorLoc of MELM respectively outperformed the state-of-the-art WSDDN [22] and WCCN [28] by 4.2% (58.4% vs. 54.2%) and 5.8% (58.4% vs. 52.6%). With the VGG16 model, it respectively outperformed the state-of-the-art WSDDN [22] and WCCN [28] by 7.9% (61.4% vs. 53.5%) and 4.7% (61.4% vs. 56.7%). Noticeably, on the "bus", "car", "chair", and "table" classes, MELM outperformed the compared state-of-the-art methods up to 7~15%. This shows that the clique-based min-entropy strategy is more effective than the image segmentation strategy used in WCCN.
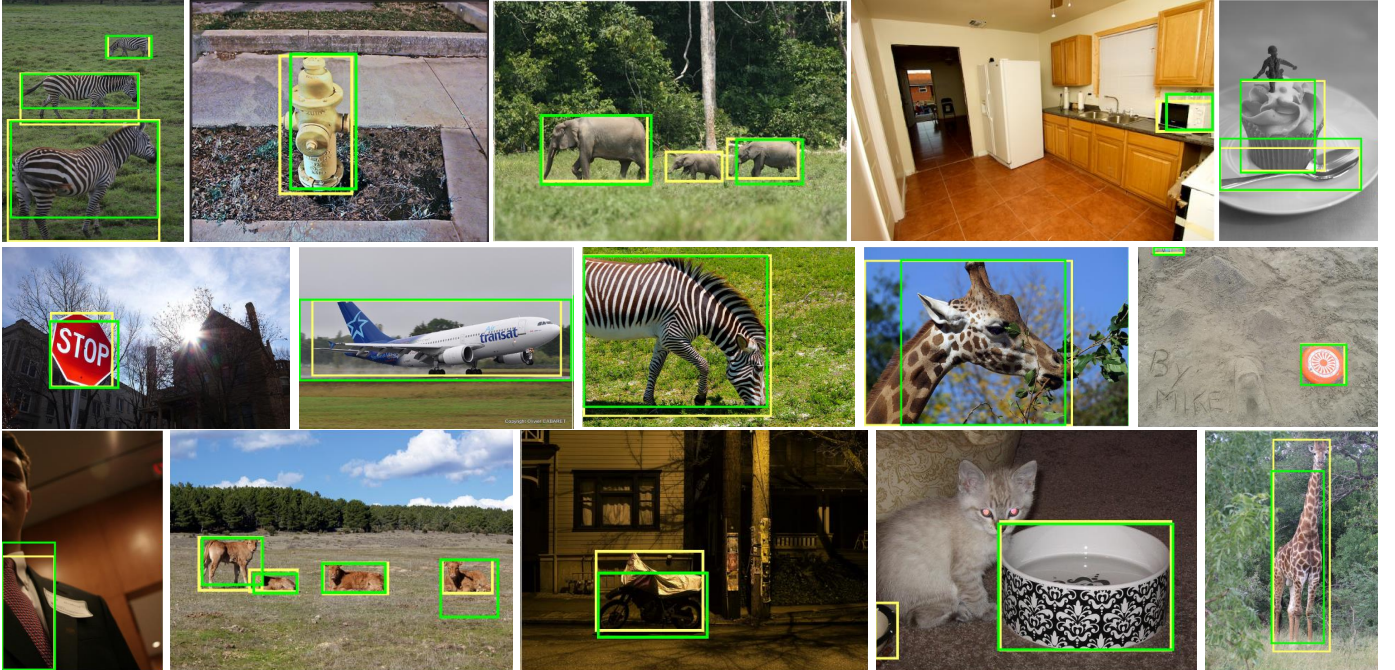
**Image Classification.** The object clique discovery and object localization components highlighted informative regions and suppressed disturbing backgrounds, which also benefited image classification. As shown in Tab. 7, with the VGGF model, MELM achieved 87.8% mAP. With the VGG16 model, MELM achieved 93.1% mAP, which respectively outperformed WSDDN [22] and WCCN [28] up to 3.4% (93.1% vs. 89.7%) and 2.2% (93.1% vs. 90.9%). It is noteworthy that MELM outperformed the VGG16 network, specifically trained for image classification, by 3.8% mAP (93.1% vs. 89.3%).

### 4.3.2 Large-scale datasets

On the ILSVRC2013 dataset with 200 object classes, Table 5, MELM with VGGF outperformed the WCCN approach by 3.6% (13.4% vs. 9.8%). On the MS COCO 2014 dataset, we evaluated the image classification, pointing localization, and object detection performance and compared it with the state-of-the-arts. The evaluation metrics for image classification included macro/micro precision (P-C and P-O), macro/micro recall (R-C and R-O), macro/micro F1-measure (F1-C and F1-O) [64]. It can be seen in Table. 6 that for image classification MELM outperformed SPN [37] by 23.1% (79.1% vs. 56%). For pointing localization, MELM outperformed SPN by 9.8% (65.1% vs. 55.3%). For object detection, MELM outperformed WSDDN. With these experiments, we set new baselines for weakly supervised object detection on large-scale datasets.

(a) PASCAL VOC 2012



(b) MSCOCO 2014

Fig. 10. Object detection examples on the PASCAL VOC 2012 and MS COCO 2014 datasets. Yellow bounding boxes denote ground-truth annotations, green boxes correct detection results and red boxes false detection results. (Best viewed in color).

## 5 CONCLUSION

In this paper, we proposed an effective deep min-entropy latent model (MELM) for weakly supervised object detection (WSOD). MELM was deployed as three components of clique partition, object clique discovery, and object localization, and was unified with the deep learning framework in an integrated manner. By partitioning and discovering cliques, MELM provided a new way to learn latent object regions from redundant object proposals. With the min-entropy principle, it can principally reduce the variance of positive instances and alleviate the ambiguity of detectors. With the recurrent learning algorithm, MELM improved the performance of weakly supervised detection, weakly supervised localization, and image classification, in striking contrast with state-of-the-art approaches. The underlying reality is that min-entropy results in minimum randomness of an information system and the recurrent learning takes advantages of continuation optimization, which provides fresh insights for weakly supervised learning problems.

## APPENDIX

For succinct representation, we denote $E_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c,\theta)$, $E_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h,\theta)$, $L_{(\mathcal{X},\mathcal{Y})}(\mathcal{H}_c,\theta)$, and $L_{(\mathcal{X},\mathcal{Y},\mathcal{H}_c)}(h,\theta)$ as $E(\mathcal{H}_c,\theta)$, $E(h,\theta)$, $L(\mathcal{H}_c,\theta)$, and $L(h,\theta)$, respectively.

**Derivation for object clique discovery.** Given the object score $s(y, h; \theta)$ as the input of the entropy models, its gradient can be computed as

$$\frac{\partial L(\mathcal{H}_c, \theta)}{\partial s(y, h; \theta)} = \sum_{y', h'} \frac{\partial L(\mathcal{H}_c, \theta)}{\partial p(y', h'; \theta)} \frac{\partial p(y', h'; \theta)}{\partial s(y, h; \theta)}$$

$$= \sum_{y', h'} \left( y' \frac{\partial E(\mathcal{H}_c, \theta)}{\partial p(y', h'; \theta)} + \frac{y' - 1}{1 - p(y', h'; \theta)} \right) \frac{\partial p(y', h'; \theta)}{\partial s(y, h; \theta)}, \tag{16}$$

where the partial derivation of $E(\mathcal{H}_c, \theta)$ with respect to $p(y, h; \theta)$ is computed as

$$\frac{\partial E(\mathcal{H}_c, \theta)}{\partial p(y', h'; \theta)} = \frac{-1}{\sum_c w_{\mathcal{H}_c} \sum_{h \in \mathcal{H}_c} p(y, h; \theta)} \cdot$$

$$\left( \frac{1}{|\mathcal{H}_c'|} \left( \sum_{h \in \mathcal{H}_c'} p(y', h; \theta) \right) \left( \sum_{y \neq y'} p(y, h; \theta) \right) \right. \tag{17}$$

$$\left. / \left( \sum_y p(y, h; \theta) \right)^2 + w_{\mathcal{H}_c'} \right),$$

where $\mathcal{H}_{\rfloor}^{'}$ is the clique including $h'$. The partial derivation of $p(y, h; \theta)$ with respect to $s(y, h; \theta)$ is computed as

$$\frac{\partial p(y', h'; \theta)}{\partial s(y, h; \theta)} = \begin{cases} -s(y', h'; \theta) s(y, h; \theta), h \neq h' \text{ or } y \neq y', \\ s(y', h'; \theta) - s(y, h; \theta)^2, otherwise. \end{cases} \tag{18}$$

**Derivation for object localization.** In Eq. (11), the term $w_h p(y, h; \theta)$ is used as a pseudo label for $h$, which does not back-propagate gradients. Therefore, the derivation for object localization can be simply computed as

$$\frac{\partial L(h, \theta)}{\partial s(y, h; \theta)} = \sum_{y', h'} \frac{\partial L(h, \theta)}{\partial p(y', h'; \theta)} \frac{\partial p(y', h'; \theta)}{\partial s(y, h; \theta)}$$

$$= \sum_{y', h' \in \{\mathcal{H}_1^*, \mathcal{H}_2^*, \ldots\}} w_{h'} \frac{\partial p(y', h'; \theta)}{\partial s(y, h; \theta)}. \tag{19}$$

The partial derivation of $L(h, \theta)$ with respect to $s(y, h; \theta)$ is calculated with Eq. (18) and Eq. (19).

## ACKNOWLEDGMENTS

## REFERENCES

[1] G. Ross, D. Jeff, D. Trevor, and M. Jagannath, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2014, pp. 580–587.

[2] R. Girshick, "Fast r-cnn," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 1440–1448.

[3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Adv. in Neural Inf. Process. Syst. (NIPS)*, 2015, pp. 91–99.

[4] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2016, pp. 21–37.

[5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 779–788.

[6] L. Liu, J. Chen, P. Fieguth, G. Zhao, R. Chellappa, and M. Pietikäinen, "From bow to cnn: Two decades of texture representation for texture classification," *Int. J. Comput. Vis*, pp. 1–36, 2018.

[7] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikinen, "Deep learning for generic object detection: A survey," *arXiv preprint arXiv:1809.02165*, 2018.

[8] A. Stuart, T. Ioannis, and H. Thomas, "Support vector machines for multiple-instance learning," in *Adv. in Neural Inf. Process. Syst. (NIPS)*, 2002, pp. 561–568.

[9] P. Megha and L. Svetlana, "Scene recognition and weakly supervised object localization with deformable part-based models," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2011, pp. 1307–1314.

[10] C. R. Gokberk, V. Jakob, and S. Cordelia, "Multi-fold mil training for weakly supervised object lcalization," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. Workshop*, 2014, pp. 2409–2416.

[11] ——, "Weakly supervised object localization with multi-fold multiple instance learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 1, pp. 189–203, 2016.

[12] B. Hakan, P. Marco, and T. Tuytelaars, "Weakly supervised object detection with posterior regularization," in *Brit. Mach. Vis. Conf. (BMVC)*, 2014, pp. 1997–2005.

[13] S. H. Oh, L. Y. Jae, J. Stefanie, and D. Trevor, "Weakly supervised discovery of visual pattern configurations," in *Adv. in Neural Inf. Process. Syst. (NIPS)*, 2014, pp. 1637–1645.

[14] B. Hakan, P. Marco, and T. Tinne, "Weakly supervised object detection with convex clustering," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 1081–1089.

[15] W. Chong, R. Weiqiang, H. Kaiqi, and T. Tieniu, "Weakly supervised object localization with latent category learning," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2014, pp. 431–445.

[16] W. Chong, H. Kaiqi, R. Weiqiang, Z. Junge, and M. Steve, "Large-scale weakly supervised object localization via latent category learning," *IEEE Trans. Image Process.*, vol. 24, no. 4, pp. 1371–1385, 2015.

[17] S. H. Oh, G. Ross, J. Stefanie, M. Julien, H. Zaid, and D. Trevor, "On learning to localize objects with minimal supervision," in *Proc. 31st Int. Conf. Mach. Learn. (ICML)*, 2014, pp. 1611–1619.

[18] D. Thomas, A. Bogdan, and F. Vittorio, "Weakly supervised localization and learning with generic knowledge," *Int. J. Comput. Vis*, vol. 100, no. 3, pp. 275–293, 2012.

[19] Z. Yimeng and C. Tsuhan, "Weakly supervised object recognition and localization with invariant high order features," in *Brit. Mach. Vis. Conf. (BMVC)*, 2010, pp. 1–11.

[20] S. Parthipan and X. Tao, "Weakly supervised object detector learning with model drift detection," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2011, pp. 343–350.

[21] H. Judy, P. Deepak, D. Trevor, and S. Kate, "Detector discovery in the wild: Joint multiple instance and representation learning," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, p. 797?23.

[22] B. Hakan and V. Andrea, "Weakly supervised deep detection networks," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 2846–2854.

[23] L. Dong, H. J. Bin, L. Yali, W. Shengjin, and Y. M. Hsuan, "Weakly supervised object localization with progressive domain adaptation," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 3512–3520.

[24] K. Vadim, O. Maxime, C. Minsu, and L. Ivan, "Contextlocnet: Context-aware deep network models for weakly supervised localization," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2016, pp. 350–365.

[25] R. Weiqiang, H. Kaiqi, T. Dacheng, and T. Tieniu, "Weakly supervised large scale object localization with multiple instance learning and bag splitting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 405–416, 2016.

[26] Q. Ye, T. Zhang, Q. Qiu, B. Zhang, J. Chen, and G. Sapiro, "Self-learning scene-specific pedestrian detectors using a progressive latent model," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 2057–2066.

[27] P. Tang, X. Wang, X. Bai, and W. Liu, "Multiple instance detection network with online instance classifier refinement," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 3059–3067.

[28] A. Diba, V. Sharma, A. Pazandeh, H. Pirsiavash, and L. Van Gool, "Weakly supervised cascaded convolutional networks," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 5131–5139.

[29] Z. Jie, Y. Wei, X. Jin, J. Feng, and W. Liu, "Deep self-taught learning for weakly supervised object localization," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 4294–4302.

[30] Y. Li, L. Liu, C. Shen, and A. van den Hengel, "Image co-localization by mimicking a good detectors confidence score distribution," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2016, pp. 19–34.

[31] F. Wan, P. Wei, J. Jiao, Z. Han, and Q. Ye, "Min-entropy latent model for weakly supervised object detection," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 1297–1306.

[32] X. Wang, Z. Zhu, C. Yao, and X. Bai, "Relaxed multiple-instance svm with application to object discovery," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 1224–1232.

[33] J. Wu, Y. Yu, C. Huang, and K. Yu, "Deep multiple instance learning for image classification and auto-annotation," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 3460–3469.

[34] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Is object localization for free? weakly supervised learning with convolutional neural networks," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 685–694.

[35] M. Shi and V. Ferrari, "Weakly supervised object localization using size estimates," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2016, pp. 105–121.

[36] D. Hoiem, Y. Chodpathumwan, and Q. Dai, "Diagnosing error in object detectors," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2012, pp. 340–353.

[37] Y. Zhu, Y. Zhou, Q. Ye, Q. Qiu, and J. Jiao, "Soft proposal networks for weakly supervised object localization," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017, pp. 1841–1850.

[38] M. Shi, H. Caesar, and V. Ferrari, "Weakly supervised object localization using things and stuff transfer," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017.

[39] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017.

[40] K. Kumar Singh and Y. Jae Lee, "Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017.

[41] D. Kim, D. Cho, D. Yoo, and I. S. Kweon, "Two-phase learning for weakly supervised object localization," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017.

[42] A. J. Bency, H. Kwon, H. Lee, S. Karthikeyan, and B. Manjunath, "Weakly supervised localization using deep feature maps," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2016, pp. 714–731.

[43] P. Tang, X. Wang, S. Bai, W. Shen, X. Bai, W. Liu, and A. L. Yuille, "Pcl: Proposal cluster learning for weakly supervised object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018.

[44] W. Ke, J. Chen, J. Jiao, G. Zhao, and Q. Ye, "SRN: side-output residual network for object symmetry detection in the wild," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017, pp. 302–310.

[45] J. Aczél and Z. Daróczy, "Charakterisierung der entropien positiver ordnung und der shannonschen entropie," *Acta Mathematica Academiae Scientiarum Hungarica*, vol. 14, no. 1-2, pp. 95–121, 1963.

[46] D. Bouchacourt, S. Nowozin, and M. P. Kumar, "Entropy-based latent structured output prediction," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2015.

[47] H. Bagherinezhad, M. Horton, M. Rastegari, and A. Farhadi, "Label refinery: Improving imagenet classification through label progression," *arXiv preprint arXiv:1805.02641*, 2018.

[48] E. L. Allgower and K. Georg, *Numerical Continuation Methods*, 1990.

[49] Y. Bengio, R. Collobert, and J. Weston, "Curriculum learning," in *Proc. ACM Int. Conf. Mach. Learn.*, 2009, pp. 41–48.

[50] C. Gulcehre, M. Moczulski, F. Visin, and Y. Bengio, "Mollifying networks," in *Int. Conf. Learn. Repres.*, 2017.

[51] U. J. RR, V. de Sande Koen EA, G. Theo, and S. A. WM, "Selective search for object recognition," *Int. J. Comput. Vis*, vol. 104, no. 2, pp. 154–171, 2013.

[52] Z. C. Lawrence and D. Piotr, "Edge boxes: Locating object proposals from edges," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2014, pp. 391–405.

[53] E. Mark, V. G. Luc, W. C. KI, W. John, and Z. Andrew, "The pascal visual object classes (voc) challenge," *Int. J. Comput. Vis*, vol. 88, no. 2, pp. 303–338, 2010.

[54] D. Jia, D. Wei, S. R., L. L. Jia, L. Kai, and L. F. Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2009, pp. 248–255.

[55] T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2014, pp. 740–755.

[56] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," *arXiv:1405.3531*, 2014.

[57] K. Alex, S. Ilya, and H. G. E, "Imagenet classification with deep convolutional neural networks," in *Adv. in Neural Inf. Process. Syst. (NIPS)*, 2012, pp. 1097–1105.

[58] S. Karen and Z. Andrew, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.

[59] Y. Wei, Z. Shen, B. Cheng, H. Shi, J. Xiong, J. Feng, and T. Huang, "Ts2c:tight box mining with surrounding segmentation context for weakly supervised object detection," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2018, pp. 434–450.

[60] P. Tang, X. Wang, A. Wang, Y. Yan, W. Liu, J. Huang, and A. Yuille, "Weakly supervised region proposal network and object detection," in *Proc. Europ. Conf. Comput. Vis. (ECCV)*, 2018, pp. 352–368.

[61] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 2921–2929.

[62] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778.

[63] C. Sun, M. Paluri, R. Collobert, R. Nevatia, and L. Bourdev, "Pronet: Learning to propose object-specific boxes for cascaded neural networks," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 3485–3493.

[64] X.-Z. Wu and Z.-H. Zhou, "A unified view of multi-label performance measures," *arXiv preprint arXiv:1609.00288*, 2016.

**Fang Wan** received the B.S. degree from Wuhan University, Wuhan, China, in 2013. Since 2013, he has been a Ph.D student in the School of Electronic, Electrical and Communication Engineering, University of Chinese Academy of Sciences, Beijing, China. His research interests include computer vision and machine learning, specifically for weakly supervised learning and visual object detection.
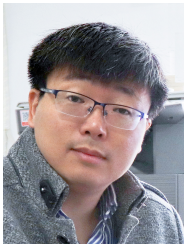
**Pengxu Wei** received the B.S. degree in computer science from the China University of Mining and Technology, Beijing, China, in 2011, and the Ph.D. degree from University of Chinese Academy of Sciences in 2018. Since 2018, she has been a research scientist at Sun Yat-sen University, Guangzhou, China. Her research interests include computer vision and machine learning, specifically for data-driven vision and scene image recognition.

**Zhenjun Han** (M'12) received the B.S. degree from Tianjin University, Tianjin, China, in 2006 and the M.S. and Ph.D. degrees from University of Chinese Academy of Sciences (UCAS), Beijing, China, in 2009 and 2012, respectively. Since 2013, he has been an Associate Professor of UCAS. His research interests include visual object detection, tracking and recognition. He has published about 40 papers in refereed conferences and journals.

**Jianbin Jiao** (M'10) received the B.S., M.S., and Ph.D. degrees from the Harbin Institute of Technology (HIT), China, in 1989, 1992, and 1995, respectively. From 1997 to 2005, he was an Associate Professor with HIT. Since 2006, he has been a Professor with the University of the Chinese Academy of Sciences, Beijing, China. In the research areas about image processing and pattern recognition. He has authored over 50 papers in refereed conferences and journals.

**Qixiang Ye** (M'10-SM'15) received the B.S. and M.S. degrees from Harbin Institute of Technology, China, in 1999 and 2001, respectively, and the Ph.D. degree from the Institute of Computing Technology, Chinese Academy of Sciences in 2006. He has been a professor with the University of Chinese Academy of Sciences (UCAS) since 2009, and was a visiting assistant professor with the University of Maryland, College Park until 2013. His research interests include visual object detection and machine learning. He has published more than 80 papers in refereed conferences and journals including IEEE CVPR, ICCV, ECCV, and PAMI, and received the Sony Outstanding Paper Award.