

Self-supervised Photographic Image Layout Representation Learning

Zhaoran Zhao, Peng Lu*, Xujun Peng, Wenhao Guo

Abstract—In the domain of image layout representation learning, the critical process of translating image layouts into succinct vector forms is increasingly significant across diverse applications, such as image retrieval, manipulation, and generation. Most approaches in this area heavily rely on costly labeled datasets and notably lack in adapting their modeling and learning methods to the specific nuances of photographic image layouts. This shortfall makes the learning process for photographic image layouts suboptimal. In our research, we directly address these challenges. We innovate by defining basic layout primitives that encapsulate various levels of layout information and by mapping these, along with their interconnections, onto a heterogeneous graph structure. This graph is meticulously engineered to capture the intricate layout information within the pixel domain explicitly. Advancing further, we introduce novel pretext tasks coupled with customized loss functions, strategically designed for effective self-supervised learning of these layout graphs. Building on this foundation, we develop an autoencoder-based network architecture skilled in compressing these heterogeneous layout graphs into precise, dimensionally-reduced layout representations. Additionally, we introduce the LODB dataset, which features a broader range of layout categories and richer semantics, serving as a comprehensive benchmark for evaluating the effectiveness of layout representation learning methods. Our extensive experimentation on this dataset demonstrates the superior performance of our approach in the realm of photographic image layout representation learning.

Index Terms—Image layouts, Graph, Representation learning.

I. INTRODUCTION

Image layout, the strategic arrangement of elements within an image, is a critical determinant of how we perceive and convey visual contents [1]. It governs the information flow, directs the viewer’s gaze, and profoundly influences our comprehension and emotional response to an image [2]. In this context, layout representation learning, which seeks to encode image layouts into compact, low-dimensional vectors within an embedding space, has emerged as a pivotal area in multimedia. It finds applications across image layout retrieval, manipulation, generation, and scene comprehension, underscoring its significance. [3]–[9]

Methods for extracting photographic image layout representations can be categorized into two primary approaches: supervised and weakly supervised methods. Supervised methods, such as those found in [10], [11], rely on labeled

layout category information to train networks. These networks generate final layer features as image layout representations. However, the creation of labeled datasets for these methods often demands expertise in aesthetics, making it cost-prohibitive. Consequently, the limited availability of labeled data restricts the generalization of supervised layout representation learning.

In response to these limitations, a growing trend favors weakly supervised methods. These approaches leverage tasks aligned with image layout learning, such as aesthetic assessment [12]–[14] and aesthetic strengthening [15], [16]. They incorporate corresponding losses from these tasks during training, enabling the acquisition of embedded representations enriched with layout information from extensive datasets. This effectively addresses the challenge posed by the scarcity of labeled data. Nevertheless, it is crucial to note that weakly supervised methods inherently capture semantic elements and their relationships, which may limit their generalizability across varied scenes in photographic images.

In recent times, self-supervised methods [17]–[19] tailored specifically for representing graphic design layouts have emerged as promising alternatives. These methods demonstrate a notable capability for modeling the latent layout information concealed within pixel domain through the utilization of graph data structures, thereby benefiting from human-derived insights. Importantly, they possess an inherent capacity to effectively filter out extraneous information, enabling concentrated modeling of layout information within the visual content. However, the direct application of these methods to layout representation learning for photographic images encounters significant challenges, primarily attributable to two key factors.

Firstly, graphic design layouts typically demonstrate a consistent granularity at the object level for layout primitives, which can be effectively represented using a homogeneous graph that models both the primitives and their interrelations. In contrast, photographic image layouts, as illustrated in Fig. 1, display a significantly higher level of diversity and cannot be adequately represented using merely object-level primitives. This necessitates the incorporation of more intricate layout primitives and their complex relationships.

Secondly, the effectiveness of self-supervised learning heavily depends on the design of pretext tasks, which are artificial tasks devised to facilitate learning in the absence of labeled data. In the realm of self-supervised learning aimed at capturing the complex layout primitives and their interrelations in photographic images, there is a significant lack of suitable pretext tasks that can be directly applied. As a result, it is crucial to redevelop pretext tasks that can accurately recover diverse layout primitives and their relationships within photo-

* Corresponding author.

Z. Zhao, P. Lu and W. Guo are with the School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, China (e-mail: zhaorzhaoran@bupt.edu.cn; lupeng@bupt.edu.cn; whguo@bupt.edu.cn).

X. Peng is with Amazon Annapurna Labs (e-mail: penxunjun@amazon.com).

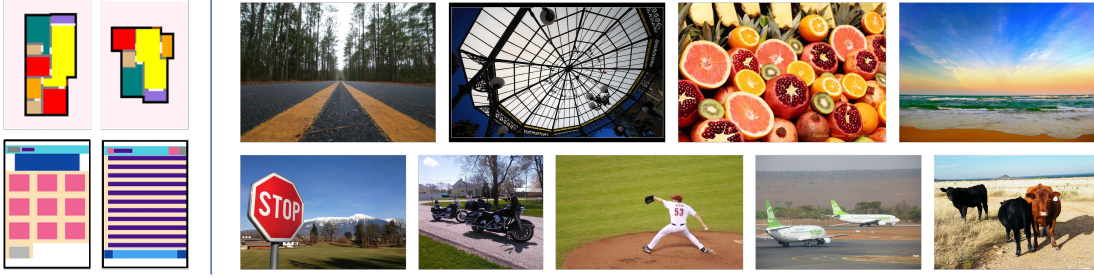


Fig. 1: Graphic design images, as seen in left side images (samples of floorplan and UI design from RPLAN [22] and RICO [23]), comprise object-level elements and their relationships. Photographic images often have a more disorganized subject arrangement or lack clear main subjects but feature distinct visual guiding lines, as shown in the images on the right side, which are sampled from proposed LODB.

graphic images.

To address the diversity inherent in modeling the layout information hidden in the pixel domain of photographic images, our study introduces multi-granularity layout primitives tailored to their intricate and varied characteristics, in line with Gestalt theory [20], [21]. Consequently, we develop a heterogeneous layout graph model to effectively encapsulate these layout primitives and their interconnections, thereby explicitly and computationally modeling the layout information hidden within the pixel domain.

To achieve effective compression of this complex heterogeneous graph into precise, low-dimensional layout representation vectors, we design specific pretext tasks and corresponding loss functions. Building upon this foundation, we construct a layout representation learning network utilizing an autoencoder-based architecture to precisely and effectively embed the information carried by the layout graph into fixed-dimension, low-dimensional vectors.

Furthermore, we introduce the LODB dataset equipped with more detailed layout labels and a broader spectrum of semantic content. This diversity renders the dataset highly conducive to applications in image layout retrieval and classification evaluation. The performance of our approach on LODB is demonstrative of its cutting-edge capabilities, highlighting its efficacy and potential in advancing the field of layout representation learning, thereby fostering the development of a wide range of tasks in the multimedia domain.

In summary: 1) A heterogeneous graph is constructed to model photographic image layouts. 2) Multiple pretext tasks and associated losses are designed for self-supervised learning of heterogeneous layout graphs. 3) An evaluation dataset named LODB is constructed to assess the effectiveness of image layout representation learning. 4) Our approach achieves state-of-the-art retrieval performance on LODB.

II. RELATED WORK

This section presents a concise yet thorough review of the literature on image layout representation learning, encompassing supervised layout classification, weakly-supervised layout representation learning, and self-supervised approaches for graphic design images.

A. Supervised Image Layout Classification

Pioneering work in supervised image layout classification includes Lee et al.’s 2018 study [10], which utilized a deep learning model based on AlexNet to learn layout information via classification labels. However, this method lacked a specialized structure for layout learning, leading to suboptimal representation capabilities. Another significant contribution is SampNet [11], which implemented a multi-pattern pooling module aligned with established composition patterns for layout analysis. Despite its innovation, SampNet’s reliance on predefined templates limited its adaptability to a wide variety of image layouts. These approaches, while groundbreaking, are constrained by their dependence on manual annotations, which can impede performance due to insufficiently diverse or large training datasets.

B. Weakly-supervised Image Layout Representation Learning

The advancement in weakly-supervised image layout representation learning has been marked by notable efforts to efficiently embed layout information, a key factor in image assessment and aesthetic strengthening [12], [13], [15], [16], [24]. The representative works, [12], [13], utilized graph neural networks for embedding layout through inter-patch relationships, exemplifies this trend. Similarly, Ke et al. [14] adopted Transformers to enhance patch interaction modeling. While these studies represent significant strides in embedding layouts through weak supervision, they tend to emphasize feature map patch relationships, potentially introducing semantic biases into the learning process.

C. Self-supervised Graphic Design Layout Representation Learning

In graphic design layout representation learning [23], [25]–[27], recent investigations have focused on Graph Convolutional Networks (GCNs) for simulating interactions in user interface (UI) components, thereby improving UI layout modeling [17], [18]. A notable advancement by [19] introduced hierarchical graphs and auto-encoders to capture semantic and positional relationships within UI elements. While these methodologies have proven effective in specific applications like UI design and floorplan generation, their applicability to

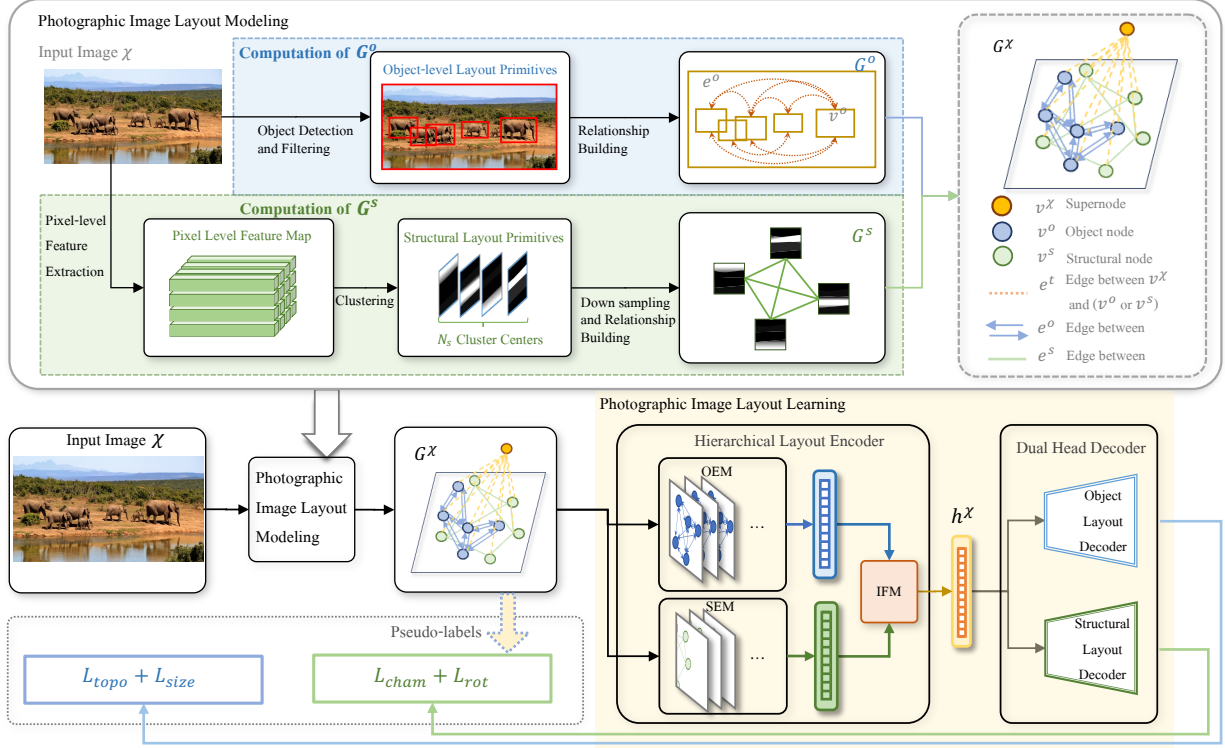


Fig. 2: In the illustrated workflow, for a given input image χ , two layout subgraphs, G^o and G^s , are initially constructed. Based on these, relationships connecting heterogeneous nodes with the supernode can be established to form a heterogeneous graph G^X . G^X is processed through an autoencoder, which has been developed according to pretext tasks designed in this paper. Comprising a Hierarchical Layout Encoder and a Dual Head Decoder, this autoencoder is designed to integrate four tailored loss functions. This integration facilitates the computing of a comprehensive embedding h^X at the full-graph level, effectively representing the layout of image χ .

photographic images, particularly those with complex, overlapping objects, remains limited due to the inherent constraints of task design in graphic design-oriented approaches.

D. Datasets for Supervised Image Layout Classification and Assessment

The development of datasets has been crucial for supervised image layout classification. The KU-PCP dataset [10], a pioneering large-scale dataset, categorizes photographic composition into nine classes but is primarily focused on outdoor scenes, leading to a semantic-layout coupling issue. Its reliance on coarse-grained labels limits its utility in nuanced image layout retrieval. Similarly, the CADB dataset [11], derived from the AADB, faces challenges in label granularity and suitability for layout retrieval assessments, as it includes images with non-typical layouts.

III. METHODOLOGY

This section introduces the proposed framework for image layout modeling and representation learning. As illustrated in Fig. 2, our approach consists of two main phases: (1) Image layout modeling involves transitioning from the pixel domain to a heterogeneous graph, encompassing the definition of layout primitives, computation of their attributes, and mapping of their interrelations. (2) Representation learning from the

heterogeneous layout graph starts with establishing pretext tasks and corresponding loss functions for accurate and effective layout embedding and concludes with the development of encoder-decoder architectures based on pretext tasks.

A. Photographic Image Layout Modeling

A substantial number of photographic images lack a clear subject, or the subject is not the decisive factor influencing the layout of the image (as illustrated in the first row of images in the right part of Fig. 1). This presents a challenge in designing effective primitives for layout modeling. Specifically, the fine granularity of these primitives introduces computational challenges, while a coarser granularity could result in the loss of critical layout information. To address this issue, we have drawn upon Gestalt theory [20], [21], which posits that people tend to perceive visually similar elements as a unified whole. Accordingly, we propose a novel approach to structural layout primitives, utilizing these primitives and their co-occurrence relationships to construct a structural layout subgraph.

While this approach effectively captures the overall layout of images, it falls short in accurately modeling the precise relative positions and sizes of objects. To address this issue, we introduce object-level layout primitives. These primitives enhance the composition of image layouts by considering

interactions that include relative positions and sizes among objects.

Integrating these two types of primitives, a heterogeneous graph is constructed to computationally represent image layout. For a specific image χ , we construct its heterogeneous graph G^χ , which includes two subgraphs and a virtual super-node. Subgraph $G^s(V^s, E^s)$, consists of $N^s = |V^s|$ nodes representing structural layout primitives and their relationships. Subgraph $G^o(V^o, E^o)$ includes $N^o = |V^o|$ nodes for object-level layout primitives. The super-node v^χ , aggregating information h^χ from all nodes, connects to each via undirected edges $e_j^t \in E^t$. This setup facilitates adaptive learning and information exchange among heterogeneous nodes.

Structural Layout Modeling: According to Gestalt theory [20], [21], primitives with visual similarity and proximity tend to be perceived as unified wholes. This raises the question: how can we effectively capture these structurally significant unified wholes?

An initial strategy is to use masks generated by deep neural semantic segmentation networks as structural primitives. However, this approach may heavily depend on the performance of the segmentation network, which poses challenges in scenarios involving unseen semantic categories. To overcome this limitation, we adopt a clustering method based on image features for segmentation. The regions obtained through this clustering serve as our primitives. We start by employing a deep learning network for pixel-wise feature extraction, using HRNet [28] to obtain feature maps. We then apply clustering techniques to these feature maps, specifically KMeans [29] and Gaussian Mixture Models (GMM) [30]. Upon comparison, we find GMM superior to KMeans due to its probabilistic nature and soft segmentation capabilities, which allow for a more nuanced and flexible delineation of image regions. This enhances the robustness of our primitive computation method in a wider variety of scenes.

In G^s , each node $v_i^s \in V^s$ represents a structural primitive, with undirected fully connected edges $e_{ij}^s \in E^s$ establishing co-occurrence relationships between node pairs v_i^s and v_j^s . This structure allows the network to adaptively adjust weights, reflecting the co-occurrence relationships among nodes.

Object-Level Layout Modeling: Images often feature one or more salient objects, and understanding their interactions is crucial for capturing object-level layout information [31]. We represent these objects and their interconnections using a graph G^o , with bidirectional edges $e_{ij}^o \in E^o$ capturing the geometric relationships between object nodes v_i^o and v_j^o , encoding their relative size and positional relationships.

Building upon the approach in [18], we leverage information regarding the size and position of bounding boxes, obtained through object detection in [32]. These bounding box attributes serve as essential information encapsulated within the nodes v^o , while the edges of the graph convey critical aspects such as relative size, distance, and directional relationships among identified objects.

In more detail, we establish specific filtering criteria for photographic images, applying these to object primitives based on the detection network’s confidence, area, and nesting properties. Our goal is to include a significant number of bounding

boxes without overly prioritizing semantic correctness. Thus, we set a confidence threshold of 0.4 for detection boxes, excluding those below this threshold, smaller than one percent of the total image area, or fully nested within larger boxes.

B. Learning Photographic Image Layout

Once we have modeled the layout information, denoted as G^χ , from the pixel domain of image χ , our primary goal is to derive a low-dimensional vector, h^χ , associated with the supernode v^χ . To mitigate our dependence on extensive labeled training datasets, we propose the adoption of a self-supervised methodology.

In graph-level embedding acquisition, contrast learning [33] is a prevalent method. However, it falls short in intuitiveness and efficiency, particularly when applied to the proposed layout graph embedding, where it exhibits limited interpretability and is unsuitable for visual recognition tasks. In the realm of graphic design, employing semantic segmentation as a pretext task is a common strategy [17], [19]. These methods involve rasterizing semantic regions according to semantic categories and supervising the regression of the rasterized images by Mean Square Error (MSE). Such physical constraints are instrumental in fulfilling the requirement of permutation invariance [34] during the graph-level embedding process.

Despite their advantages, these approaches are not suitable for photographic image layout representation learning for the following reasons: The pretext tasks fail to adequately differentiate between instances within the same semantic category, leading to a poor representation of their topological relationships. Moreover, accurately capturing the topological relationships between different semantic categories remains challenging. Crucially, for enhancing the generalization of layout representation, it is essential to decouple semantic information from the layout information embedding process. In this context, a significant unresolved issue is ensuring permutation invariance in graph-level embedding, especially during the design of pretext tasks.

Therefore, to more effectively address the embedding of heterogeneous layout graphs, we break down the task and design pretext tasks specifically tailored to the attributes of nodes and relationships. The following detailed discussion will offer insights into the design of these pretext tasks and their corresponding losses, along with the autoencoder network designed based on the pretext tasks.

Pretext Task Designed for Structural Layout Information

Reconstruction: In the context of recovering structural layout information, E^s serves as the medium for information propagation, and V^s is the target for reconstruction. It is worth emphasizing that, due to the permutation invariance property and the nature of layouts, the order in which graph nodes are reconstructed does not impact the final reconstruction outcome. Successful reconstruction is achieved as long as there is a one-to-one correspondence between primitives in the input node set and primitives in the output node set, with each output node accurately mapping back to its corresponding input node. This principle is similar to point cloud reconstruction, where the arrangement of points does not affect the accuracy of the

reconstruction, provided the original positions of the points can be accurately restored.

Hence, the reconstruction of structural layout information is akin to the high-dimensional point cloud reconstruction task. To quantify the distance between V^s and the reconstructed \hat{V}^s , we employ the widely used Chamfer Loss in point cloud reconstruction, with the modified equation as follows:

$$L_{cham} = \sum_{i=1}^{|V^s|} \min_j \|h_i^s - \hat{h}_j^s\|^2 + \sum_{j=1}^{|V^s|} \min_i \|\hat{h}_j^s - h_i^s\|^2 \quad (1)$$

However, unlike the 3D point cloud reconstruction task, our reconstruction target has higher spatial dimensions and is more sparsely distributed. Therefore, using only L_{cham} would result in information loss or redundancy in the reconstruction. To address this issue and ensure the accurate reconstruction of information for each node, we introduce a new loss term denoted as L_{rot} . Specifically, the tensor composed of information carried by all nodes from the set V_s is represented as H_s . Assuming the successful reconstruction of information for each node in the input graph G^s , the difference between H^s and the reconstructed \hat{H}^s is an orthogonal matrix, yielding the identity matrix when multiplied by its transpose. Under this concept, the formula for L_{rot} is as follows:

$$L_{rot} = |f(P) - I| \quad (2)$$

where $P = H^s \hat{H}^s^{-1}$ and $f(x) = \text{softmax}(x/\varepsilon)$ ensures that each row of tensor x approximates a one-hot vector with $\varepsilon = 0.01$. Both L_{rot} and L_{cham} are employed simultaneously for recovering structural layout information.

Pretext Task Designed for Object Level Layout Information Reconstruction: Object layout information encompasses tasks involving center positions, relative relationships, and sizes. To address this, we introduce two pretext tasks: keypoint detection and size regression. The keypoint detection task is employed to reconstruct center positions and topological relationships. Due to significant area imbalances between central and non-central regions, using an L2 loss may disproportionately favor larger non-central regions, thereby undermining the detection of smaller central regions. To mitigate this issue, we adopt Focal Loss, which addresses the imbalance of positive and negative samples. Following the methodology described in [32], we map the center positions of object level layout primitives in χ to a low-resolution heatmap $\mathcal{Y} \in [0, 1]^{W \times H}$, where \mathcal{Y}_{xy} represents the probability of an object center existing at a specific position (x, y) . The Focal Loss used to reconstruct center positions and topological relationships can be expressed as:

$$L_{topo} = \frac{-1}{|V^o|} \sum_{xy} \begin{cases} (1 - \hat{\mathcal{Y}}_{xy})^\alpha \log(\hat{\mathcal{Y}}_{xy}) & \text{if } (\mathcal{Y}_{xy} = 1) \\ (1 - \mathcal{Y}_{xy})^\beta (\hat{\mathcal{Y}}_{xy})^\alpha \log(1 - \hat{\mathcal{Y}}_{xy}) & \text{else} \end{cases} \quad (3)$$

where α and β are hyper-parameters of the focal loss, we use $\alpha = 2$ and $\beta = 4$ following [32].

In the size regression task, we predict the area ratio when the target is present. To accommodate continuous values, we introduce *Quality Focal Loss* [35] to reconstruct L_{size} , enhancing performance in size regression tasks. Similar to the computation of \mathcal{Y}_{xy} , we map the area proportion of objects onto a low-resolution heatmap $\mathcal{S} \in (0, 1)^{W \times H}$, where \mathcal{S}_{xy} denotes the area proportion of an object with position (x, y) . L_{size} can be expressed as:

$$L_{size} = - \sum_{xy} |\mathcal{S}_{xy} - \hat{\mathcal{S}}_{xy}|^\gamma ((1 - \mathcal{S}_{xy}) \log(1 - \hat{\mathcal{S}}_{xy}) + \mathcal{S}_{xy} \log(\hat{\mathcal{S}}_{xy})) \quad (4)$$

The total training loss is expressed as:

$$L_{total} = k_1 \cdot L_{cham} + k_2 \cdot L_{rot} + k_3 \cdot L_{size} + k_4 \cdot L_{topo} \quad (5)$$

where k_1, k_2, k_3, k_4 represent the weights of different loss terms.

C. Encoder and Decoder Design

The input consists of a heterogeneous graph denoted as G^x , comprising two subgraphs, G^o and G^s . To achieve the compressed encoding of this graph and capture the information h^x carried by v^x , a Hierarchical Layout Encoder is introduced. This encoder is composed of three modules: the Structural Layout Information Encoding Module (SEM), the Object-based Layout Information Encoding Module (OEM), and the Information Fusion Module (IFM). In the SEM, three layers of Graph Attention Networks replicate information exchange among ESPs. An Attention Pooling layer is also integrated to capture co-occurrence relationships. Following a strategy from [17], the OEM creates triplets and employs Graph Convolutional Network layers to compute interactions between object primitives, thus obtaining embedding information for G^o . The IFM uses a fully connected network to simulate information fusion between heterogeneous nodes, resulting in the embedding vector h^x for G^x .

To fully reconstruct the information in G^x , a Dual Head Decoder is designed consisting of a Structural Layout Decoder and an Object Layout Decoder. Convolution and transposed convolution are combined for embedded vector upsampling. The Object Layout Decoder generates a $C \times W \times H$ grid heatmap, with $C = 2$ and $W = H = 80$, encapsulating topological and size details of the object layout. In the Structural Layout Decoder, the reconstruction target is a tensor with shape $|V^s| \cdot X \cdot Y$ and $X = Y = 20$ in our implementation.

D. Implementation Details

The proposed network is supervised with our designed loss, where $k_1 = 2, k_2 = 40, k_3 = 0.5, k_4 = 0.1$ for network training with *Adam*. The initial learning rate is set to 0.001 for the first 50 epochs and dampened to 0.001 for the rest 150 epochs.

TABLE I: Abbreviations and image counts for various image layout categories in LODB

Abbre.	Class	Num.
Cent.	Rule of Center	947
RoT.L	Rule of Three (left)	214
RoT.R	Rule of Three (right)	275
O2Dia.L	Two objects lie diagonally from top-left corner	301
O2Dia.R	Two objects lie diagonally from the top-right corner	293
O2Hor.	Horizontal Layout for two objects	386
O3Li.	Three Objects on the same line	157
O3Tri.	Triangle Layout (Three objects)	115
Oline.	More than three objects arranged on a same line	110
Pat.	Image full of patterns	255
DiaL.	Diagonal structure (from the top-left corner)	446
DiaR.	Diagonal structure (from the top-right corner)	236
Hor.	Horizontal Structure	574
Tri.	Triangle Structure	451
Ver.	Vertical Structure	455
Radi.	Radial Structure	187
DiaX.	Bi-diagonal Structure	861

IV. LODB DATASET

To achieve effective image layout representation, algorithms should be capable of accurately encoding layout information across a wide array of semantic scenes, while distinctly differentiating between various image layouts. However, current datasets for layout classification and evaluation commonly encounter the following problems: Firstly, the granularity of labels is relatively coarse, with considerable variations in layout present under the same label, impeding the precise assessment of layout representation learning. Additionally, the validation set sizes of these datasets are typically limited, coupled with an inadequate representation of semantic diversity, thus hindering a comprehensive evaluation of the algorithm’s effectiveness in diverse scenarios. To address these issues, we introduce a novel evaluation dataset – LODB¹.

LODB exhibits enhanced granularity in labeling, an extensive spectrum of semantic scenes, and a broader scale, comprising 17 diverse categories with a cumulative total of 6029 images designated for validation purposes. Representative sample images from this dataset are illustrated in the right section of Fig. 1. Comprehensive descriptions, abbreviations, and the respective image counts for each category are systematically delineated in Table I. Such advancements in the LODB dataset are crucial for enabling a more thorough and precise evaluation of algorithms for image layout representation.

A. Dataset Construction

We collaborated with five experts in the fields of photography and painting to establish standards and selection

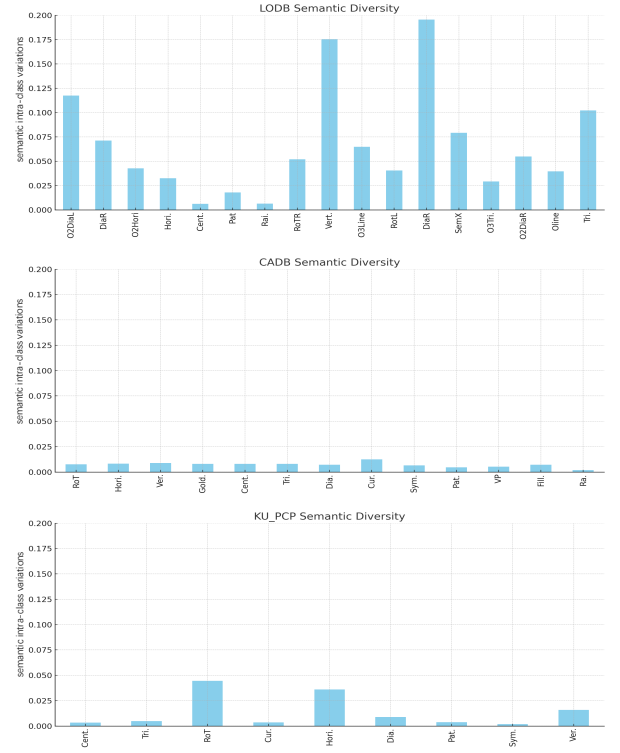


Fig. 3: The intra-class semantics variations comparison for the proposed LODB, CADB, and KU-PCP.

criteria for our dataset categories. Specifically, grounded in the fundamental principles of photographic composition, we meticulously selected 17 categories from typical compositions. For annotation, each image underwent a rigorous categorization process by these experts, determining its alignment with either a single or multiple layout types. A consensus among at least three experts was required to confirm an image’s classification. Any image failing to achieve this agreement threshold was omitted from the dataset. Additionally, we employed specialized scripts to consistently monitor and ensure semantic diversity within each category, as further elucidated in Eq. 6, ensuring a robust and varied dataset.

The images were sourced from various open-source datasets, including PARA [36], AADB [37], AVA-test [38], ADE20k [39], COCO2017 [40], FLICKR-AES [41], TAD66k [42], VOC2007 [43], and VOC2012 [44].

B. Qualitative and quantitative analysis of LODB

Our selection process emphasized reliability and semantic diversity by excluding atypical images and minimizing semantic repetition within each layout category. The proposed dataset exhibits significantly higher semantic intra-class variations compared to CADB [11] and KU-PCP [10]. Moreover, the labels are more detailed in LODB, we distinguish not only the quantity and arrangement of main layout elements but also variations in the orientation of similar structural information. The layout categories and semantic diversity of the dataset are depicted in Fig. 3.

¹<https://github.com/CV-xueba/Image-Layout-Learning>

We employ the magnitude of intra-class variations [45] to assess the semantic diversity of images within different layout categories. A larger intra-class variation indicates a higher dispersion of semantics, implying richer semantic information. Specifically, for a layout category c , the calculation of the semantic intra-class variations d^c is formulated as shown in the equation below:

$$d^c = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i^c - \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^c \right\|_2 \quad (6)$$

where N denotes the number of images within layout category c . The notation \mathbf{x}_i^c represents the semantic vector corresponding to an image, which is obtained using the method in [46].

V. EXPERIMENTS

We evaluate the effectiveness of layout representation learning methods through retrieval tasks. These tasks allow for a more detailed and intuitive assessment of the effectiveness and precision of layout representation learning by examining the retrieval results and their order. This evaluation is conducted using the proposed LODB dataset, which offers a variety of semantic scene categories and layout types.

A. Datasets & Metrics

Training Set: A subset of 80,000 images randomly sampled from the training set of AVA dataset [38].

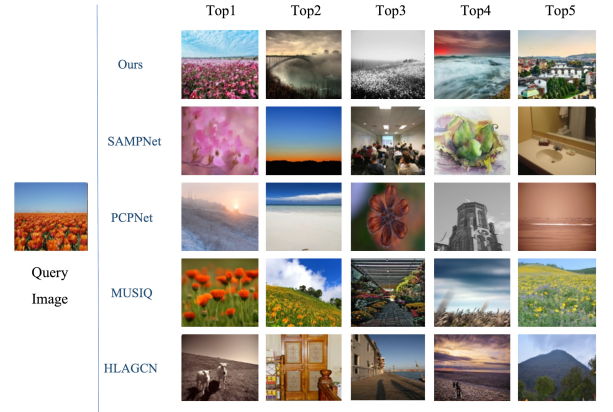
Evaluation Set: The proposed LODB dataset.

Metrics In the provided LODB dataset, an evaluation is conducted using a consistent set of query and gallery images for all baselines. More specifically, within each category, around ten percent of the images are randomly selected and used as query images, while the rest are designated as gallery images. We then followed [47], [48] and compute the average $F1@k$ and the $mAP@k$ for all query images (where $k = 1, 5, 10, 20$). Furthermore, we compute the $mAP@k$ (where $k=20$) [49] for each layout category present within the query images, facilitating a more comprehensive comparative analysis.

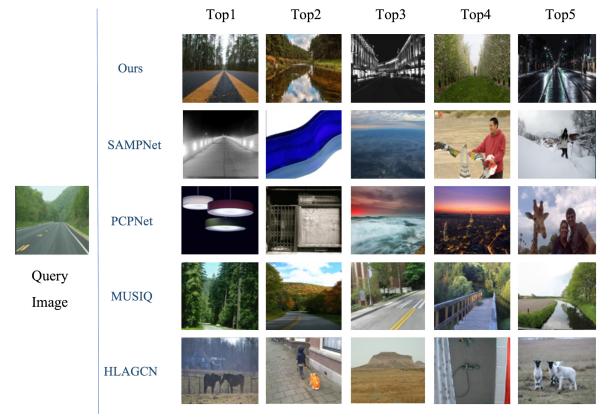
B. Comparison with SOTA Methods

We conducted a comparative analysis of our proposed network's performance in learning image layout representation against five baseline layout extractors. These baseline extractors were constructed based on state-of-the-art methods as follows:

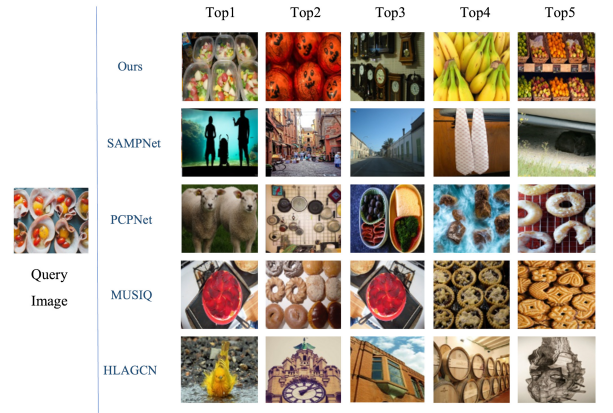
- In the context of supervised approaches, we employed the network architecture preceding the classification layer from SampNet [11] and PCPNet (implemented as per Lee et al. [10]) to create layout extractors.
- We utilized the weakly supervised methods MUSIQ [14] and HLAGCN [13], which is explicitly tailored for image assessment and is closely aligned with photographic image layout learning, to establish layout extractors. The components preceding the classification layer in these methods served as our layout extractors.
- To effectively compare our proposed method with self-supervised graphic design approaches [17], [19], we



(a) Query image with layout category *Hor.* from LODB



(b) Query image with layout category *DiaX.* from LODB



(c) Query image with layout category *Pat.* from LODB

Fig. 4: Top 5 retrieval results for some representative query images from different categories of LODB.

adopted the segmentation pretext task and the corresponding MSE loss from these studies and replaced the original ones in our method. This model, named as GDNet in this paper, was trained on the same dataset as ours. And the encoder of GDNet was used as layout extractors. Given the limitations of these self-supervised methods

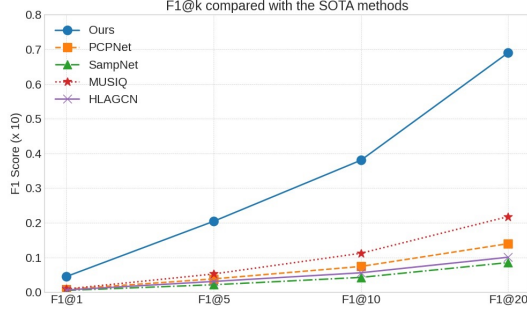


Fig. 5: The $F1@k$ of the proposed method, the supervised methods, and the weakly supervised methods.

in encoding structured layout information, we ensured fairness in our comparison by focusing exclusively on retrieval performance within layout categories with well-defined subjects.

TABLE II: Comparisons with supervised and weakly supervised methods. Each row represents $mAP@20$ for a specific layout category and the last row indicates the average $mAP@20$ of all categories.

Method	Ours	Samp-Net	PCP-Net	MUSIQ	HLA-GCN
Cent.	0.899	0.272	0.314	0.361	0.310
RoT.L	0.877	0.093	0.116	0.116	0.155
RoT.R	0.919	0.189	0.120	0.181	0.163
O2Dia.L.	0.741	0.073	0.177	0.200	0.205
O2Dia.R.	0.710	0.152	0.185	0.215	0.145
O2Hor.	0.751	0.105	0.172	0.151	0.160
O3Li.	0.834	0.025	0.090	0.137	0.108
O3Tri.	0.725	0.000	0.000	0.080	0.300
Oline.	0.750	0.500	0.000	0.000	0.042
Pat.	0.880	0.079	0.230	0.343	0.117
Dia.L.	0.539	0.201	0.141	0.106	0.223
Dia.R.	0.528	0.109	0.147	0.161	0.148
Hor.	0.851	0.320	0.430	0.746	0.309
Tri.	0.594	0.225	0.148	0.153	0.111
Ver.	0.662	0.266	0.244	0.278	0.230
Radi.	0.870	0.252	0.426	0.357	0.258
DiaX.	0.679	0.252	0.389	0.605	0.434
avg.	0.820	0.203	0.262	0.339	0.234

Comparisons with Supervised Methods: Table II and Fig. 5 reveal that in retrieval experiments conducted on previously unobserved datasets, supervised methods demonstrate significantly lower values for both $mAP@20$ and $F1@k$, thus underscoring their limited generalization capabilities. In Fig. 4, we present a qualitative comparison of our approach with supervised methods in the second and third rows in (a)(b)(c). Notably, when evaluated on the unseen dataset, LODB, supervised methods face challenges due to limited generalization capability, resulting in suboptimal retrieval performance among all baselines.

Comparisons with Weakly-supervised Methods: Weakly supervised approaches, including MUSIQ and HLAGCN, exhibit limitations in retrieval performance when compared to our method, as shown in Table II and Fig. 5. These limitations

highlight their challenges in effectively capturing and understanding image layout knowledge. The core issue with these methods stems from their reliance on Graph Neural Networks (GNN) or Transformers for direct information extraction from feature maps, blending semantic and layout information and thus disrupting layout representation learning. Our observations, evident in the fourth and fifth rows of (a)(b)(c) in Fig. 4—particularly in the retrieval results of MUSIQ—clearly illustrate how semantic interference affects layout information. This underscores the critical importance of minimizing semantic disruption during the layout representation learning process.

TABLE III: Comparisons with self-supervised methods. Each row represents $mAP@20$ for a specific layout category and the last row indicates the average.

	Obj.	Ours	GDNet
Cent.	1	0.899	0.851
RoT.L	1	0.877	0.770
RoT.R	1	0.919	0.896
O2Dia.L.	2	0.741	0.634
O2Dia.R.	2	0.710	0.667
O2Hor.	2	0.751	0.729
O3Li.	3	0.834	0.691
O3Tri.	3	0.725	0.551
Oline.	>4	0.750	0.505
avg.		0.826	0.755

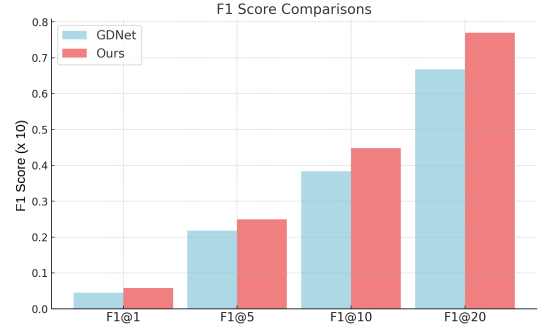


Fig. 6: The $F1@k$ ($k = 1, 5, 10, 20$) for the proposed method and the GDNet.

Comparisons with Self-supervised Methods: As shown in Fig. 6, our method outperforms GDNet across all layout categories, as evidenced by higher $F1@20$ scores. This highlights the precision of our approach in capturing layout information. Additionally, performance analysis in Table III reveals that our method excels in identifying the positions of individual objects (*Cent.*, *RoT.L*, *RoT.R*), as indicated by the second column's count of layout elements. However, our method's retrieval effectiveness diminishes in increasingly complex layouts, especially in scenarios with multiple objects (two, three, or more).

In terms of quality, the first to third rows in Fig. 7 demonstrate an incremental increase in the number of objects (noted as Obj. in Table III). In the first row, which represents the centered composition type, both our method and GDNet successfully retrieve this layout. However, our approach distinctly preserves structural information in the



Fig. 7: The retrieval results for the proposed method and the GDNNet. The leftmost query images are from LODB with the type of *Cent.*, *O2Hor.* and *O3Line*.

retrieval results, especially capturing curves and large circular regions. Moving to the second row, which depicts two objects arranged in parallel, and the third row, featuring three objects arranged in line, the precision of our retrieval results is notably higher. In particular, in the second row, our retrieval maintains background structural features, such as the horizontal layout, more effectively. In contrast, due to its insufficient modeling of structural layout specific to photographic images, GDNNet faces challenges in these scenarios. Furthermore, as indicated in Table III, with an increase in the number of objects, the decline in retrieval performance is more pronounced for GDNNet compared to our method.

C. Ablation Studies

To demonstrate the rationality of pretext task design and the necessity of losses, we conducted six ablation experiments compared with the full model, including G1: Training the network branch for Structural Layout Information Reconstruction along with its specific losses. G2: Similarly, training for Object Level Layout Information Reconstruction. Additionally, we have G3: Training without L_{cham} , G4: without L_{rot} , G5: without L_{size} , and G6: without L_{topo} . The results of these ablation experiments are summarized in Table IV.

For detailed ablation studies on pretext task design, we disabled the corresponding branch and trained the network accordingly. In the case of network losses, we set the corresponding loss weights to zero, keeping all other aspects consistent with the full model.

Effectiveness of Pretext Tasks Design: Ablation experiments G1 and G2 show that combining structural graph reconstruction with topology and size information recovery significantly boosts overall model performance. The complete model achieves the highest mAP scores across diverse retrieval metrics, highlighting that separately modeling and learning

TABLE IV: F1@k (k = 1, 5, 10, 20) and $mAP@20$ of pretext tasks and losses

F1@k ($\times 10$)	k=1	k=5	k=10	k=20	mAP@20
Full	0.045	0.204	0.380	0.690	0.820
G1	0.035	0.158	0.291	0.518	0.667
G2	0.036	0.171	0.313	0.573	0.698
G3	0.012	0.057	0.107	0.200	0.322
G4	0.037	0.163	0.294	0.517	0.682
G5	0.037	0.162	0.300	0.541	0.685
G6	0.035	0.157	0.285	0.513	0.667

specific types of information is insufficient for effectively representing image layouts.

Effectiveness of L_{rot} and L_{cham} : Comparing G3 and G4 with the full model reveals a clear decrease in mAP , emphasizing the importance of including both L_{rot} and L_{cham} to capture geometric interactions among layout elements effectively. Also, the inclusion of L_{rot} is vital, as it works in tandem with L_{cham} to jointly supervise the reconstruction of structural layout information.

Effectiveness of L_{size} and L_{topo} : Ablation experiments G5 and G6 highlight the importance of L_{size} and L_{topo} in bolstering the model’s capacity to restore size information and topological relationships. The omission of L_{size} has a more pronounced impact on the outcomes compared to the omission of L_{topo} . This discrepancy arises because L_{size} also contributes to recovering topological information, utilizing area data from object center positions as positive samples for supervision. Nevertheless, the necessity of L_{topo} is evident, as the synergy between L_{size} and L_{topo} amplifies the network’s ability to learn object level layout.

VI. CONCLUSION AND FUTURE WORK

Our research marks a significant advancement in photographic image layout representation. We developed a unique graph model and an autoencoder-based network to effectively handle the complexities of this task. The introduction of the LODB enhances our work, offering a benchmark to assess layout representation learning. Our method shows state-of-the-art performance on LODB, highlighting its potential in advancing the field, thereby fostering the development of a wide range of tasks in the multimedia domain.

Although our current approach introduces an innovative method for modeling and learning photographic image layouts without needing large-scale labeled datasets, it relies on a two-stage network process. We aim to refine this into a more efficient, end-to-end network framework in our future research.

REFERENCES

- [1] G. Ambrose and P. Harris, *Basics Design 02: Layout*. Ava Publishing, 2011.
- [2] J. Wang and S. Langer, “A brief review of human perception factors in digital displays for picture archiving and communications systems,” *Journal of Digital Imaging*, vol. 10, pp. 158–168, 1997.
- [3] X. Wujiang, Y. Xu, G. Sang, L. Li, A. Wang, P. Wei, and L. Zhu, “Recursive multi-relational graph convolutional network for automatic photo selection,” *IEEE Transactions on Multimedia*, 2022.
- [4] M. Zhang, M. Li, J. Yu, and L. Chen, “Aesthetic photo collage with deep reinforcement learning,” *IEEE Transactions on Multimedia*, 2022.

- [5] H. Tan, B. Yin, K. Wei, X. Liu, and X. Li, "Alr-gan: Adaptive layout refinement for text-to-image synthesis," *IEEE Transactions on Multimedia*, 2023.
- [6] C. Li, P. Zhang, and C. Wang, "Harmonious textual layout generation over natural images via deep aesthetics learning," *IEEE Transactions on Multimedia*, vol. 24, pp. 3416–3428, 2021.
- [7] J. Cheng, F. Wu, L. Liu, Q. Zhang, L. Rutkowski, and D. Tao, "Indecgan: Learning to generate complex images from captions via independent object-level decomposition and enhancement," *IEEE Transactions on Multimedia*, 2023.
- [8] M. Song, G.-M. Um, H. K. Lee, J. Seo, and W. Kim, "Dynamic residual filtering with laplacian pyramid for instance segmentation," *IEEE Transactions on Multimedia*, 2022.
- [9] F.-L. Zhang, M. Wang, and S.-M. Hu, "Aesthetic image enhancement by dependence-aware object recomposition," *IEEE Transactions on Multimedia*, vol. 15, no. 7, pp. 1480–1490, 2013.
- [10] J.-T. Lee, H.-U. Kim, C. Lee, and C.-S. Kim, "Photographic composition classification and dominant geometric element detection for outdoor scenes," *Journal of Visual Communication and Image Representation*, vol. 55, pp. 91–105, 2018.
- [11] B. Zhang, L. Niu, and L. Zhang, "Image composition assessment with saliency-augmented multi-pattern pooling," *arXiv preprint arXiv:2104.03133*, 2021.
- [12] J. Hou, S. Yang, and W. Lin, "Object-level attention for aesthetic rating distribution prediction," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 816–824.
- [13] D. She, Y.-K. Lai, G. Yi, and K. Xu, "Hierarchical layout-aware graph convolutional network for unified aesthetics assessment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8475–8484.
- [14] J. Ke, Q. Wang, Y. Wang, P. Milanfar, and F. Yang, "Musiq: Multi-scale image quality transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5148–5157.
- [15] P. Lu, H. Zhang, X. Peng, and X. Jin, "Learning the relation between interested objects and aesthetic region for image cropping," *IEEE Transactions on Multimedia*, vol. 23, pp. 3618–3630, 2020.
- [16] S. Ni, F. Shao, X. Chai, H. Chen, and Y.-S. Ho, "Composition-guided neural network for image cropping aesthetic assessment," *IEEE Transactions on Multimedia*, 2022.
- [17] D. Manandhar, D. Ruta, and J. Collomosse, "Learning structural similarity of user interface layouts using graph networks," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*. Springer, 2020, pp. 730–746.
- [18] A. G. Patil, M. Li, M. Fisher, M. Savva, and H. Zhang, "Layoutgmn: Neural graph matching for structural layout similarity," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 048–11 057.
- [19] Y. Bai, D. Manandhar, Z. Wang, J. Collomosse, and Y. Fu, "Layout representation learning with spatial and structural hierarchies," 2023.
- [20] P. T. Quinlan and R. N. Wilton, "Grouping by proximity or similarity? competition between the gestalt principles in vision," *Perception*, vol. 27, no. 4, pp. 417–430, 1998.
- [21] P. Kałamała, A. Sadowska, W. Ordziniak, and A. Chuderski, "Gestalt effects in visual working memory: Whole-part similarity works, symmetry does not," *Experimental Psychology*, vol. 64, no. 1, p. 5, 2017.
- [22] W. Wu, X.-M. Fu, R. Tang, Y. Wang, Y.-H. Qi, and L. Liu, "Data-driven interior plan generation for residential buildings," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 6, pp. 1–12, 2019.
- [23] B. Deka, Z. Huang, C. Franzen, J. Hirschman, D. Afargan, Y. Li, J. Nichols, and R. Kumar, "Rico: A mobile app dataset for building data-driven design applications," in *Proceedings of the 30th annual ACM symposium on user interface software and technology*, 2017, pp. 845–854.
- [24] L. Zhang, M. Song, Y. Yang, Q. Zhao, C. Zhao, and N. Sebe, "Weakly supervised photo cropping," *IEEE Transactions on Multimedia*, vol. 16, no. 1, pp. 94–107, 2013.
- [25] T. F. Liu, M. Craft, J. Situ, E. Yumer, R. Mech, and R. Kumar, "Learning design semantics for mobile apps," in *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*, 2018, pp. 569–579.
- [26] D. Manandhar, H. Jin, and J. Collomosse, "Magic layouts: Structural prior for component detection in user interface designs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 809–15 818.
- [27] J. Jin, Z. Xue, and B. Leng, "Shrag: Semantic hierarchical graph for floorplan representation," in *2022 International Conference on 3D Vision (3DV)*. IEEE, 2022, pp. 271–279.
- [28] Y. Yuan, X. Chen, and J. Wang, "Object-contextual representations for semantic segmentation," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*. Springer, 2020, pp. 173–190.
- [29] D. Arthur and S. Vassilvitskii, "K-means++ the advantages of careful seeding," in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 2007, pp. 1027–1035.
- [30] J. D. Banfield and A. E. Raftery, "Model-based gaussian and non-gaussian clustering," *Biometrics*, pp. 803–821, 1993.
- [31] E. S. Spelke, K. Breinlinger, K. Jacobson, and A. Phillips, "Gestalt relations and object perception: A developmental study," *Perception*, vol. 22, no. 12, pp. 1483–1501, 1993.
- [32] L. Cai, Z. Zhang, Y. Zhu, L. Zhang, M. Li, and X. Xue, "Bigdetection: A large-scale benchmark for improved object detector pre-training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4777–4787.
- [33] Y. Liu, M. Jin, S. Pan, C. Zhou, Y. Zheng, F. Xia, and S. Y. Philip, "Graph self-supervised learning: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 6, pp. 5879–5900, 2022.
- [34] R. Winter, F. Noé, and D.-A. Clevert, "Permutation-invariant variational autoencoder for graph-level representation learning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 9559–9573, 2021.
- [35] X. Li, W. Wang, L. Wu, S. Chen, X. Hu, J. Li, J. Tang, and J. Yang, "Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 002–21 012, 2020.
- [36] Y. Yang, L. Xu, L. Li, N. Qie, Y. Li, P. Zhang, and Y. Guo, "Personalized image aesthetics assessment with rich attributes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 19 861–19 869.
- [37] S. Kong, X. Shen, Z. Lin, R. Mech, and C. Fowlkes, "Photo aesthetics ranking network with attributes and content adaptation," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. Springer, 2016, pp. 662–679.
- [38] N. Murray, L. Marchesotti, and F. Perronnin, "Ava: A large-scale database for aesthetic visual analysis," in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 2408–2415.
- [39] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ade20k dataset," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 633–641.
- [40] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [41] J. Ren, X. Shen, Z. Lin, R. Mech, and D. J. Foran, "Personalized image aesthetics," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 638–647.
- [42] S. He, Y. Zhang, R. Xie, D. Jiang, and A. Ming, "Rethinking image aesthetics assessment: Models, datasets and benchmarks," in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, 2022, pp. 942–948.
- [43] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, pp. 303–338, 2010.
- [44] M. Everingham and J. Winn, "The pascal visual object classes challenge 2012 (voc2012) development kit," *Pattern Anal. Stat. Model. Comput. Learn., Tech. Rep.*, vol. 2007, no. 1–45, p. 5, 2012.
- [45] Z. Fu, Z. Mao, B. Hu, A.-A. Liu, and Y. Zhang, "Intra-class adaptive augmentation with neighbor correction for deep metric learning," *IEEE Transactions on Multimedia*, 2022.
- [46] FoamLiu, "A github repository for scene classification," <https://github.com/foamliu/Scene-Classification>, 2023, accessed: August 18, 2023.
- [47] D. Huynh and E. Elhamifar, "A shared multi-attention framework for multi-label zero-shot learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8776–8786.
- [48] L. Yang, B. Geng, Y. Cai, A. Hanjalic, and X.-S. Hua, "Object retrieval using visual query context," *IEEE Transactions on multimedia*, vol. 13, no. 6, pp. 1295–1307, 2011.
- [49] S. Li, Z. Chen, J. Lu, X. Li, and J. Zhou, "Neighborhood preserving hashing for scalable video retrieval," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8212–8221.