

MDFF: A Method for Fine-Grained UFZ Mapping With Multimodal Geographic Data and Deep Network

Song Ouyang¹, Shihong Du¹, Xiuyuan Zhang¹, Shouhang Du¹, and Lubin Bai¹

Abstract—As basic units of urban areas, urban functional zones (UFZs) are fundamental to urban planning, management, and renewal. UFZs are mainly determined by human activities, economic behaviors, and geographical factors, but existing methods 1) do not fully use multimodal geographic data owing to a lack of semantic modeling and feature fusion of geographic objects and 2) are composed of multiple stages, which lead to the accumulation of errors through multiple stages and increase the mapping complexity. Accordingly, this study designs a multimodal data fusion framework (MDFF) to map fine-grained UFZs end-to-end, which effectively integrates very-high-resolution remote sensing images and social sensing data. The MDFF extracts physical attributes from remote sensing images and models socioeconomic semantics of geographic objects from social sensing data, and then fuses multimodal information to classify UFZs where object semantics guide the fine-grained classification. Experimental results in Beijing and Shanghai, two major cities of China, show that the MDFF greatly improves the quality of UFZ mapping with the accuracy about 5% higher than state-of-the-art methods. The proposed method significantly reduces the complexity of UFZ mapping to complete the urban structure analysis conveniently.

Index Terms—Deep learning, image classification, multimodal geographic data fusion, remote sensing, urban functional zone (UFZ) mapping.

I. INTRODUCTION

URBAN functional zones (UFZs) are spatially represented by zones with the same social functions, such as commercial, residential, and industrial zones [1]. UFZs are often employed as basic units to analyze urban spatial structures and social characteristics [2], which have become increasingly crucial for urban planning [3], [4], e.g., transportation planning, factory relocation, environmental protection, and sustainable development [5], [7]. However, fine-grained UFZ maps with large areas are hardly available [8] due to 1) the lack of methods for fine-grained UFZ mapping and 2) unreliable interpretation

of UFZs. The rapid development of urbanization and economy has brought a huge demand for timely updating UFZ maps. Therefore, urban studies urgently need to break through the existing difficulties to obtain large-scale and fine-grained UFZ maps that keep up with the pace of urban development in a timely and accurate manner [9].

As functions of urban zones are mainly determined by human activities, economic behaviors, and geographical factors, UFZs present complex physical structures and socioeconomic properties, which makes the UFZ mapping full of challenges. With the development of earth observation technologies, lots of high-quality remote sensing (RS) images are produced, especially for very-high-resolution satellite images (VHR) which hold advantages in representing UFZs, because of the large spatial coverage, detailed information, and wide availability [10], [11]. However, it is a challenge to accurately and reliably interpret UFZs relying on a single source data. For example, with only the satellite image (see Fig. 1), it is difficult to decide whether the blue area in (a) is a residential or a commercial area and which category the shadow area of (b) belongs to, because UFZs are not only related to physical attributes, but also to socioeconomic attributes.

There have been substantial multimodal geographic data to represent UFZs, such as image data (VHR, synthetic aperture radar images, nightlight images, etc.) and social sensing data (e.g., points of interest, mobile signal data, social media data), which constitute complementary information sources [12]. RS images provide pixel-level characteristics but lack information on the characterization of geographic objects, while social sensing data contain socioeconomic attributes and correlations of geographic objects. Therefore, these two types of data are complementary. Comprehensive utilization of multimodal information can help to reduce the uncertainty of interpretation. With socioeconomic properties of points of interest (POIs) attached to the satellite image in Fig. 1, the blue area in (a) and the shadow area in (b) are easily assigned as the commercial and the transport, respectively. However, existing studies do not fully fuse multimodal semantics to assist the classification and are complex for multistage processes, which significantly limits the quality and reliability of UFZ mapping. To map fine-grained UFZs end-to-end with multimodal data, three challenging issues must be resolved, i.e., multimodal feature representation, multimodal feature fusion, and urban-function classification.

Manuscript received 18 August 2023; revised 3 October 2023; accepted 14 October 2023. Date of publication 20 October 2023; date of current version 6 November 2023. This work was supported by the Chinese National Nature Science Foundation under Grant 42001327. (Corresponding author: Xiuyuan Zhang.)

The authors are with the Institution of Remote Sensing and Geographic Information System, School of Earth and Space Sciences, Peking University, Beijing 100871, China (e-mail: song_ouyang@outlook.com; shdu@pku.edu.cn.com; xy_zhang@pku.edu.cn; dsjcug@163.com; lbbai@stu.pku.edu.cn).

Digital Object Identifier 10.1109/JSTARS.2023.3326160

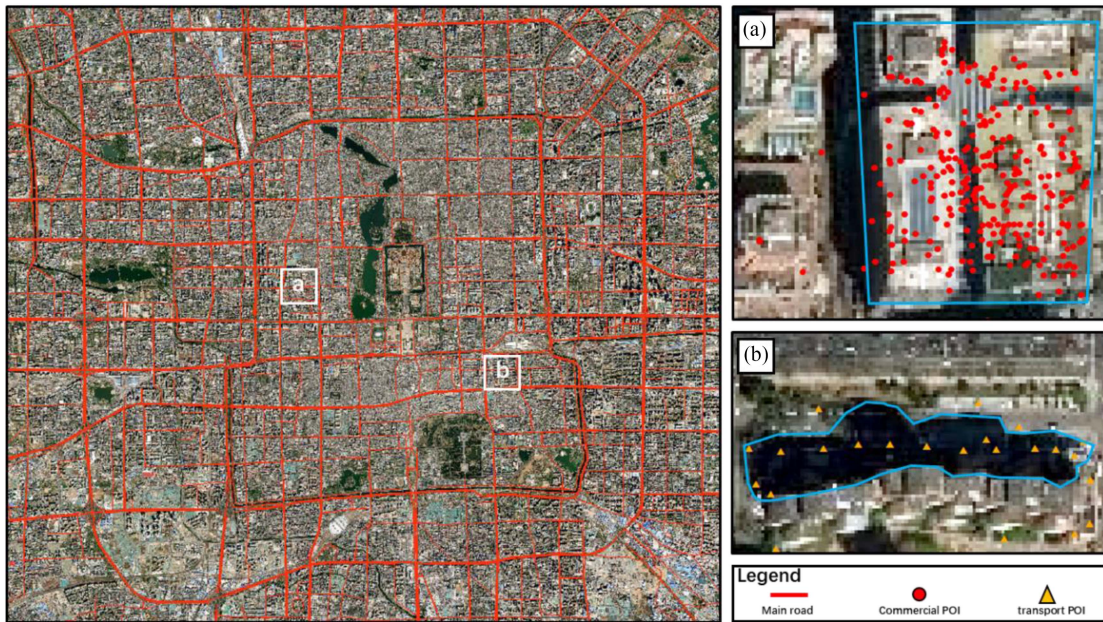


Fig. 1. Examples of interpreting UFZs in satellite images.

A. Multimodal Feature Representation

Appearance features extracted from images, such as spectrums, textures, and geometry, are adopted in the early land cover and urban function classification [13], [14], [15]. These features are shallow-level features because of their weak expressiveness and robustness. For better discrimination, middle-level cues are generated from shallow-level features by the developed methods such as the principal component analysis, the probability latent semantic analysis, the support vector machine, and the latent Dirichlet allocation [16], [19]. With the rapid progress of deep learning, convolutional neural networks (CNNs) have great advantages in extracting deep-level features from images for their powerful learning ability [20], [21], [22]. Zhou et al. [23] employed a trained CNN to assign functional attributes with RS images. Therefore, deep features from CNNs can effectively characterize UFZs' physical properties. Unlike RS images, social sensing data contain socioeconomic semantics which are important for UFZ interpretation. However, socioeconomic semantics are difficult to extract due to discrete arrangement of social sensing data and the lack of modeling methods. Existing studies mainly used social sensing data as training labels of classifiers, samples for accuracy evaluation, or shallow-level features [24], [25]. It not only failed to mine the high-level information of social sensing data, but also discarded semantic relationships between objects. Thus, there is an urgent need for a modeling approach within social sensing data to extract UFZs' socioeconomic semantics.

B. Multimodal Feature Fusion

It is challenging to fuse multimodal features especially for RS images and social sensing data as these two types of data are arranged in completely different ways. The fusion methods are divided into two types: the pixel-wise overlay approaches and

the object-based methods [27]. The former directly concatenates images and the density maps of social sensing data pixel by pixel. Some studies follow this way [28], but these pixel-wise approaches only fuse at the pixel level and lose the object semantics [29], [30]. The latter extracts attributes of social sensing data to characterize objects for classifiers [31], [32]. Based on the units of road blocks, Zhang et al. [12] calculated discrete features from POIs and social media data, then combined them with spectral attributes to map land uses. Du et al. [24] merged spectral features and semantic features at the object level for large-scale UFZ mapping. However, existing methods mainly fuse features at a single scale (pixel scale or object scale), and shallow fusion at the attribute level rather than semantic level cannot mine objects' semantic information from multimodal geographic data. The pixel-wise approaches provide detailed features in pixel space, and the object-based methods are suitable for expressing object semantics [33]. Therefore, integrating information at both pixel and object levels to characterize UFZs is critical for UFZ mapping. Therefore, it is urgent to develop fusion methods that consider joint features and the high-level semantics of multimodal geographic data under the two levels [34], [35], [36].

C. Urban-Function Classification

Due to its importance to urban space planning, urban-function classification has attracted considerable attention in the past decade [9]. On one hand, RS images and social sensing data are employed to map UFZs in early attempts. However, multimodal geographic data were not fully utilized in UFZ mapping due to the insufficient feature expression. On the other hand, each UFZ is composed of geographic objects [37], such as the commercial zone in Fig. 1(a) including buildings and roads. Therefore, it has become the mainstream method that segments land cover

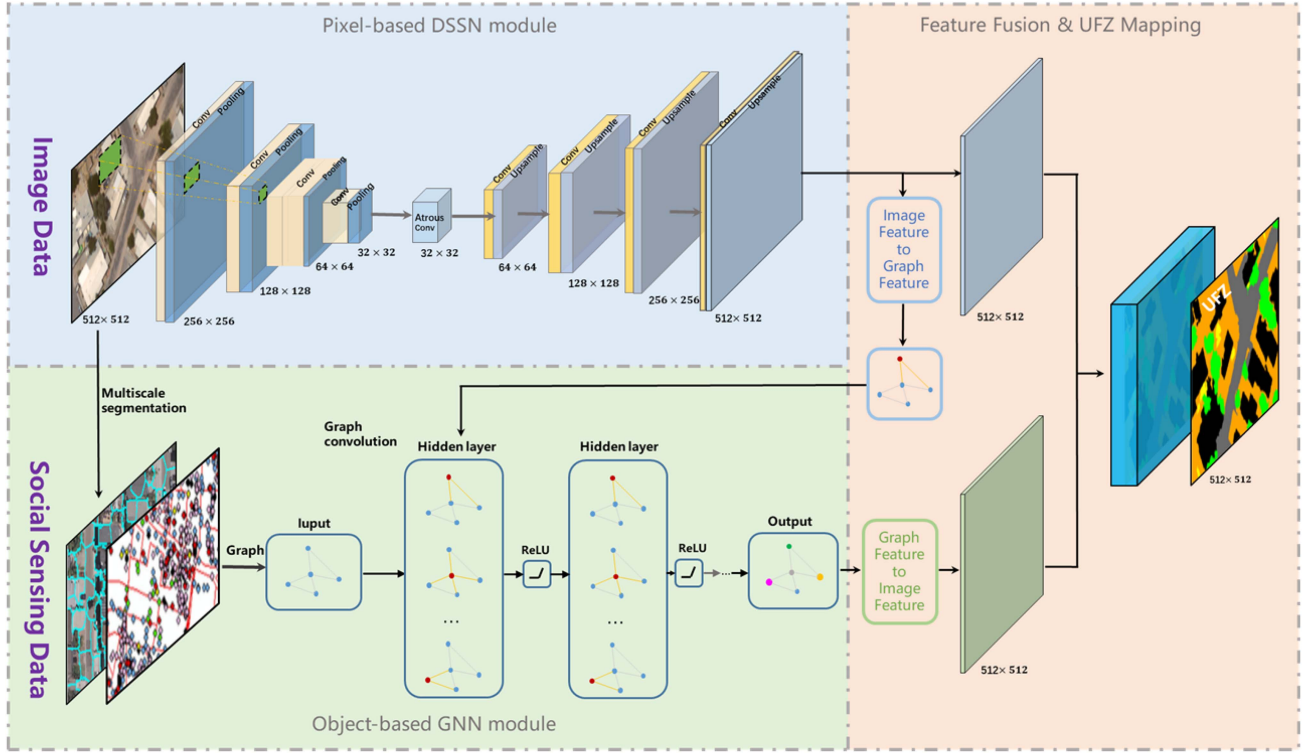


Fig. 2. Workflow of the proposed MDFF.

objects at first and then aggregates them into UFZs. In the hierarchical semantic cognition model [12], objects in the model are segmented with the multiresolution segmentation method [38], and then developed the inverse hierarchical semantic cognition model for refining the results at the land-cover layer [39]. Du et al. [9] designed a CNN to extract pixel-wised predictions for initializing UFZ categories of geographic objects, then used the conditional random field with road block constraints to regroup the objects into the fine-grained UFZ map. In general, these methods are multi-stage. The lack of the end-to-end mechanism leads to the accumulation of errors from different stages and increases the complexity of UFZ mapping. Therefore, it is necessary to develop an end-to-end method to interpret UFZs with multimodal geographic data.

In summary, existing studies have three limitations. First, making no full use of multimodal geographic data, i.e., insufficient expression of image features and the lack of socioeconomic semantic modeling of social sensing data, results in failing to interpret UFZs accurately and reliably. Second, the feature fusion of multimodal geographic data is still at the shallow level rather than the semantic level, which limits the semantic representation of UFZs. Third, existing UFZ mapping methods are multi-stage, e.g., CNN + postprocessing. There is a lack of methods fusing multimodal geographic data to map UFZs end-to-end. To solve these issues, this study designs a multimodal data fusion framework (MDFF) to map fine-grained UFZs end-to-end, which organically integrates multimodal geographic data, i.e., RS images and social sensing data. The MDFF learns the deep features from images at the pixel scale and models

socioeconomic semantics of UFZs with social sensing data at the object scale. With the designed deep feature fusion module, the MDFF fuses multimodal features at both pixel- and object scales to take advantage of the two scales in expressing detailed features and object semantics for UFZs. The training or operation of the MDFF is completed in one step. The following three contributions have been made in this study.

- 1) A method (MDFF) fusing multimodal geographic data is proposed for mapping fine-grained UFZs end-to-end. It automatically learns physical attributes from RS images and models socioeconomic semantics from social sensing data.
- 2) The MDFF fuses multimodal information both at pixel- and object scale with two fusion modules [i.e., the image feature to graph feature (IF2GF) and the graph feature to image feature (GF2IF) in Section II-B] for learning urban semantics where object semantics guide the fine-grained UFZ mapping.
- 3) Experiments conducting on two typical modal data, i.e., RS images and social sensing data, show that the MDFF can significantly improve the quality of UFZ mapping and reduce the complexity of UFZ mapping.

II. METHODOLOGY

A. Overview

Fig. 2 illustrates the workflow of the proposed MDFF.

- 1) *Multimodal feature extraction*: In the MDFF, the pixel-based module extracts deep features from RS images,

which provides detailed information for fine-grained mapping, while the object-based module generates characteristics of objects from social sensing data, which models semantics relationships between objects within UFZs.

- 2) *Multimodal feature fusion*: With the proposed feature fusion modules, i.e., the IF2GF and the GF2IF, the pixel-wise deep features and the object-wise semantic characteristics are fused at pixel and object scales. Integrating multimodal information under the “pixel-object-UFZ” semantic hierarchy is critical for mapping fine-grained UFZs and improving the reliability of the interpretation.
- 3) *UFZ mapping*: The MDFF is trained end-to-end with pixel samples and object samples, where the pixel samples are collected manually in the form of polygons, and the object samples are generated with the first category of POI data and the pixel samples. After the training, the MDFF is exploited to map UFZs with multimodal geographic data in the study area.

B. Multimodal Feature Extraction

In this Big Data era, image and social sensing data are important components of multimodal geographic data. They bring rich information for eliminating the uncertainty of the interpretation. The former provides pixel-level characteristics, while the latter contains socioeconomic properties of and mutual relationships between objects. Images and social sensing data need to be processed separately to obtain representative features because of their different characteristics. The MDFF is proposed to take full advantage of RS images and social sensing data. Therefore, the pixel-based module and the object-based module are designed in the MDFF to extract multimodal features. Each UFZ is composed of objects and pixels constitutes each object. Pixel-scale features and object-scale features are exploited to characterize UFZs in physical and socioeconomic terms.

1) *Pixel-Based Module*: The pixel-based module is designed to learn image features, which consists of an encoder and a decoder. The encoder extracts deep features from RS images with the stacked layers. Each layer is composed of the convolution and the pooling operations. The convolution scans the image step by step and generates feature representations of a local region in each step. The pooling expands the receptive field of the network to reduce information redundancy and obtain global information. The deep features with abstract expression are generated through the stacked layers. f_i^e denotes the i th layer feature from the encoder function φ_i^e which is composed of the convolution $\text{Conv}(\cdot)$ and the pooling $\text{Pool}(\cdot)$. w_i^e and b_i^e , which are learnable parameters, represent weights and bias of the convolution kernel of φ_i^e

$$f_{i+1}^e = \varphi_i^e (w_i^e \cdot f_i^e + b_i^e) \quad (1)$$

$$\varphi_i^e = \text{Pool}(\text{Conv}(\cdot)). \quad (2)$$

The encoder reduces the size of the features with the convolution and the pooling operations. The features are too small to hold fine-grained information, which hardly meets the requirements of intensive prediction tasks (e.g., UFZ mapping). Therefore, the features are enlarged to the original size with the upsampling of

the decoder. The decoder function φ_j^d generates the j th layer feature f_j^d . φ_j^d consists of the convolution $\text{Conv}(\cdot)$ and the upsampling $\text{Ups}(\cdot)$ with the learnable parameters w_j^d and b_j^d

$$f_{j+1}^d = \varphi_j^d (w_j^d \cdot f_j^d + b_j^d) \quad (3)$$

$$\varphi_j^d = \text{Conv}(\text{Ups}(\cdot)). \quad (4)$$

2) *Object-Based Module*: UFZs are aggregated from objects and the compositional patterns of objects affect the categories of UFZs [40]. To construct the graph, image objects are regarded as the graph nodes and the spatial relationships between objects as the graph edges. In this way, the object compositional patterns of UFZs are modeled in the form of the graph. The objects are characterized with object features ($f^o = [f^{\text{sp}}, f^t, f^g, f^s]$) in Table I. Spectral features (f^{sp}), textural features (f^t), and geometrical features (f^g) have been widely used to measure objects from different aspects in most object-based image analyses [25]. Features from social sensing data present the composition of socioeconomic attributes of UFZs, as social sensing data are closely related to socioeconomic activities. The vector $f^s = [f_{\text{poi}}^s, f_{\text{sm}}^s, f_{\text{ms}}^s]$ is used to express primary semantics of each object, where f_{poi}^s , f_{sm}^s , and f_{ms}^s denote features from POIs, social media data, and mobile signal data, respectively. POIs and social media data are generated by socioeconomic activities, which contain socioeconomic semantics. Mobile signal data reflect the flow of urban population and is closely related to socioeconomic activities. Therefore, these three kinds of social sensing data are selected to extract socioeconomic semantics for UFZ mapping. $f_{\text{poi}}^s = [f_{\text{density}}^s, f_1^s, f_2^s, \dots, f_c^s]$, where c denotes the second category of POIs and f_i^s represents the proportion of the i th category within an object. Similarly, f_{sm}^s can be obtained from social media data. The average population mobility of objects constitutes f_{ms}^s .

We first segment RS images into image objects using multiresolution segmentation method [38]. Then, object features f_i^o in Table I and deep features f_j^d are adopted to characterize objects. The feature vector of each node is $x_n = [f_1^o, f_2^o, \dots, f_m^o, f_j^d]$ and the graph feature is denoted as $\mathbf{X} = [x_1, x_2, \dots, x_N]^T \in \mathbb{R}^{K \times D}$, where N is the number of the nodes and D is the dimension of x_k . Finally, the first-order adjacency relationships (with common edge) between objects are adopted as graph edges to take the topological spatial relationship into consideration. The adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ quantifies the graph edges where the strength of each edge is determined by the spectral similarity and the spatial location of objects. The connecting edge between node v_i and node v_j has the weight $a_{ij} \in \mathbf{A}$. if v_i is adjacent to v_j , $a_{ij} = \exp(-\frac{\|\text{lab}_i - \text{lab}_j\|}{\sigma_w^2})$, otherwise, $a_{ij} = 0$. lab denotes the average color value of the node in the CIE LAB color space and σ_w^2 controls the range of the weight.

Graph convolution has powerful ability in modeling semantic relationships of discrete objects. Thus, the object-based module is designed with the graph convolution to extract socioeconomic semantics from social sensing data and its structure is shown in Table II. This module is a hierarchical network with L layers ($L = 3$). Each layer consists of a graph convolution function and a nonlinear activation function.

TABLE I
 OBJECT FEATURES FROM IMAGES AND POI

Types	Names	Meanings
Socioeconomic	POI Semantics	The proportion of each semantic category of POIs The density of POIs in an object
	Population Mobility	Within an object, the population mobility average over time
	Semantics of Social Media Data	The proportion of each semantic category of social media data within an object
Spectral	Mean	Average spectrum of pixels
	Std. Dev	Gray standard deviation of pixels in an object
	Skewness	Skewness of spectral histogram
	Border Contrast	The average differences between border pixels and their neighborhood
Textural	GLCM Homogeneity	The homogeneity derived from GLCM
	GLCM Dissimilarity	The heterogeneity parameters derived from GLCM
	GLCM Entropy	Information entropy derived from GLCM
	GLCM Correlation	Correlation of pixels which is derived from GLCM
	GLDV	The vector composed of diagonal elements of GLCM
	GLDV Mean	Average value of GLCM
	GLDV Entropy	Information entropy derived from GLDV
	GLDV Contrast	The Contrast of GLDV
Geometrical	GLDV Ang. 2nd moment	The Ang. 2nd moment of GLDV
	Area	The number of pixels within image objects
	Length/Width	Length-width ratio of the object's MBR
	Eclipse Fit	The fitting degree of eclipse fit
	Main Direction	Eigenvectors of covariance matrix
	Shape Index	The ratio of perimeter to four times side length

 TABLE II
 LAYER STRUCTURE OF THE OBJECT-BASED MODULE

Layers	Components	Input Size	Output Size
Input	input	$N \times 100$	–
Layer-1	graph convolution ReLU	$N \times 100$	$N \times 18$
Layer-2	graph convolution ReLU	$N \times 18$	$N \times 18$
Layer-3	graph convolution ReLU	$N \times 18$	$N \times 9$
Output	output	–	$N \times 9$

Note: . N represents the number of objects in the graph.

The object-based module implements the graph convolution [41] to extract objects' features \mathbf{X} . $\mathbf{X}^{(l)}$ is the feature from the l th layer of the graphic neural network (GNN). After the graph convolution, new features $\mathbf{X}^{(l+1)}$ are activated by the nonlinear activation function σ . The core idea of the graph convolution is that graph nodes with edge connections are aggregated for generating new representations of the nodes where semantic information flows between the nodes with edge connections, so the object-based module can model the semantic relationships between image objects and build the spatial object patterns of UFZs. On the other hand, graph nodes have the socioeconomic features from POI data, thus the object-based module

can learn both the socioeconomic attributes and relationships within $\mathbf{X}^{(l+1)}$ to restore the real state of UFZs

$$\mathbf{X}^{(l+1)} = \sigma \left(\mathbf{L}_G \mathbf{X}^{(l)} \mathbf{W}^{(l)} \right) \quad (5)$$

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I} \quad (6)$$

$$\mathbf{L}_G = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \quad (7)$$

where \mathbf{W} and \mathbf{L}_G are the Laplacian matrix and the learnable parameters, respectively, \mathbf{I} is the identity matrix and $\tilde{\mathbf{D}}$ is the degree matrix of $\tilde{\mathbf{A}}$.

C. Multimodal Feature Fusion

Features from images are pixel-regular and vision-related, while graph features from social sensing data are discrete and semantic-related, thus different modalities of features make it difficult to fuse these two kinds of features. One solution is to transform the form of the two kinds of features. Therefore, the IF2GF and the GF2IF are designed to conduct this task in the MDFF (see Fig. 3).

Due to the gradient vanishing [42], GNNs are usually limited to shallow layers, which restricts GNN's performance in feature learning. Therefore, the deep features from the pixel-based module's decoder are transformed by the IF2GF and then added to the GNN-base module (see Fig. 3), which helps GNN focus on modeling semantic relationships. In the IF2GF, deep features \mathbf{f}_1^d are transformed to graph features \mathbf{f}_1^g with the transforming

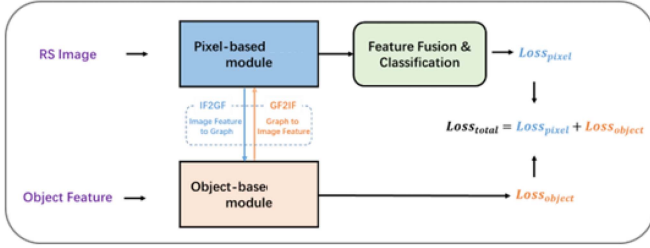


Fig. 3. Multimodal feature fusion and network training of the proposed MDF.

function T_1 . Within each object (the addressing function A), the average value on each channel of the deep features is taken as the graph feature (the assignment function S). In the object-based module, the graph features and the attribute features presented in Table II are concatenated together at object scale, and are used as the features of graph nodes

$$f_1^g = T_1 (f_1^d) \quad (8)$$

$$T_1 = S(A(\text{Mean}(\cdot))). \quad (9)$$

The object-based module outputs graph features f_1^g . In the GF2IF, the graph features are transformed to image features f_2^d with the transforming function T_2 , where the pixels in each object are assigned the same features from the graph feature of the object. Then, the GF2IF aligns the graph features with the deep features pixel to pixel. To adjust value distribution of features, the transformed features are input into a module Con with a convolution layer (3×3 kernel, 1 stride, and 1 padding), a BatchNorm layer, and a ReLU activation function. Then, both the transformed features and the deep features are concatenated together at the pixel scale, where the graph features provide additional semantics of socioeconomic relationships for the UFZ mapping

$$f_2^d = T_2 (f_2^g) \quad (10)$$

$$T_2 = \text{Con}(S(A(\cdot))). \quad (11)$$

The deep features provide detailed features in pixel space, while the graph features contain object semantics and relationships. With the feature fusion modules, i.e., the IF2GF and the GF2IF, the deep feature and the graph feature are fused at both the object and the pixel scales successively, which further help to improve features' expressiveness in depth or semantics and generate richer and more comprehensive features for UFZs. Therefore, multimodal information is integrated for better details and semantics to interpret UFZs.

D. UFZ Mapping

1) *Study Area and Data Collection*: The two largest and most famous cities in China, Beijing and Shanghai, are chosen in this study. Beijing (see Fig. 4) is the capital of China and Shanghai is the economic center of China. Beijing and Shanghai are rapidly developing at a high urbanization level and contains diverse archaic and modern urban zones. The evolution of UFZs is crucial for revealing urban sustainable development. UFZ maps are produced in the study area (within the fifth ring road of Beijing

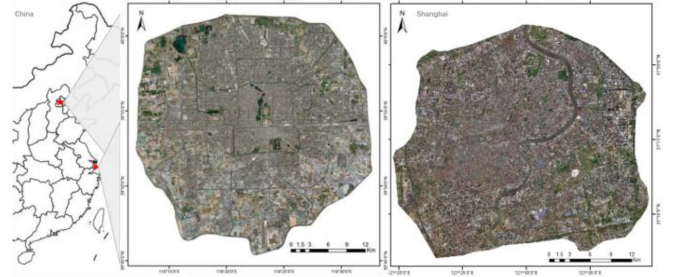


Fig. 4. Study areas in Beijing and Shanghai.

TABLE III
CATEGORY SYSTEM

Du's category system	Ours category system
Commercial	Commercial
Residential-1	Residential
Residential-2	
Residential-3	Shantytown
Institutional	Institutional
Industrial	Industrial
Transport	Transport
Urban green	Urban green
Undeveloped	Undeveloped
Woodland	–
Farmland	–
Water	Water

and the outer ring road of Shanghai, main urban area of Beijing and Shanghai) with the proposed method. For image data, the level-3 VHR satellite image (see Fig. 4) from ArcGIS World Imagery in 2019 is used as the experimental area, which cover 667 km^2 of Beijing and 664 km^2 of Shanghai with the spatial resolution of 2.4 m and the RGB band. As shown in Table III, the class system proposed in the paper [9] is adopted to map UFZs in this study. As there are no woodlands and farmlands in the main urban area of Beijing and Shanghai, the two categories are removed. The Residential-1 and the Residential-2 zones are merged into one class, Residential area, because of their similar functions.

Social sensing data are common in cities and easily accessible. It reflects socioeconomic attributes which are closely related to UFZs. POIs, mobile signal data, and social media data are typical social sensing data. Take POIs as an example, it is discrete and irregularly arranged in the geospatial space (see Figs. 5 and 6).

a) *POIs*: About 267 000 POIs of Beijing and 356 000 POIs of Shanghai are adopted in this study. The categories of POIs are divided into two levels: the first category and the second category (see Table IV). The POIs are labeled with six categories, i.e., commercial points (55.8% of Beijing), residential points (20.6%), shantytown points (3.1%), institutional points (11.4%), industrial points (8.6%), and urban green points (0.5%). There is a large class imbalance within the POI categories because of the POI's correlation with socioeconomic activities.

b) *Mobile signal data*: The data record the spatial locations of the mobile phone users when exchanging information

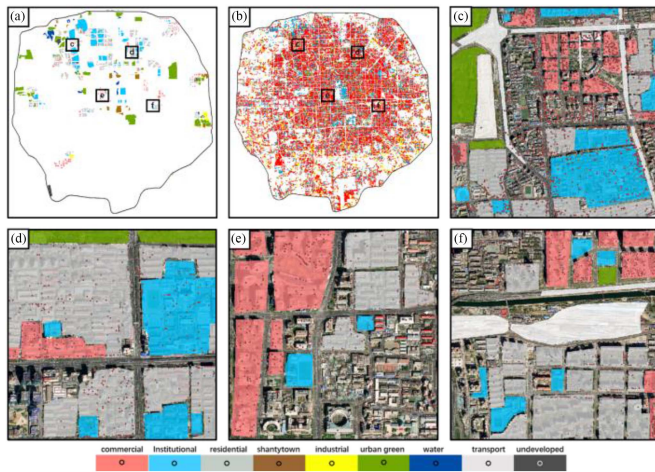


Fig. 5. UFZ samples (a) and POIs (b) of Beijing. (c)–(f) Detailed illustrations of four sub-regions.

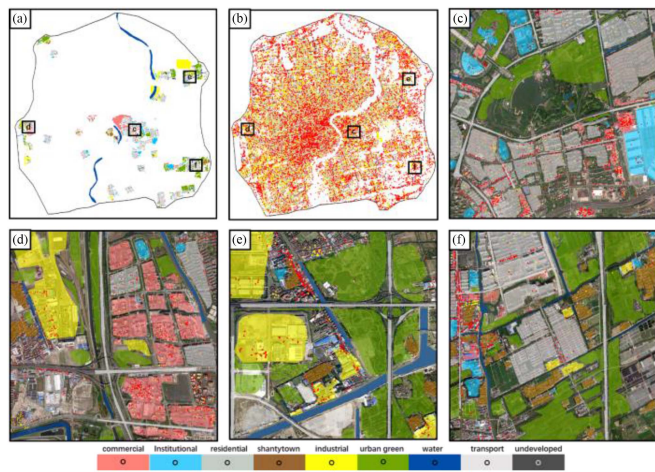


Fig. 6. UFZ samples (a) and POIs (b) of Shanghai. (c)–(f) Detailed illustrations of four sub-regions.

with the base station. It presents the spatial and temporal flow of people’s movement on a large scale. The base stations are distributed in the grid of $500\text{ m} \times 300\text{ m}$. Within a base station, local population, working population, residential population, flow population, and transit population are recorded, which are closely related to socioeconomic activities. This study adopts mobile signal data in Beijing and Shanghai from January to February 2020 (33 days in total).

c) Social media data: The data are obtained from the web of Dianping, showing commercial activities of merchants which are divided into 15 categories (see Table IV). About 268 000 records of Beijing and 688 000 records of Shanghai are used to extract socioeconomic attributes of UFZs.

The MDFF is trained end-to-end with both pixel and object samples. In the MDFF, pixel samples guide the pixel-based module to extract deep vision features and the object samples help the object-based module to learn semantic features for reliable UFZ mapping. Pixel samples of UFZs are manually selected in polygonal form, and polygons are labeled according

TABLE IV
CATEGORIES OF POIS AND SOCIAL MEDIA DATA

First category	POI		Social media data
	Second category		
Commercial	Car shop and parking lot	Gourmet shop	
Residential	Company	Shopping store	
Shantytown	Restaurant	Leisure and entertainment	
Institutional	Government agency	Scenic spot	
Industrial	Hospital	Sports venue	
Urban green	Public place	Barber and beauty shop	
–	Hotel and residential	Wedding spot	
–	Shop	Hotel	
–	Sport center	Car shop	
–	Scenic spot	Parent-child spot	
–	–	Business education institution	
–	–	Life service	
–	–	Medical health	
–	–	Home decoration	
–	–	Pet care store	

to the new class system (see Figs. 5 and 6). These samples are randomly divided into three parts with the ratio of 8:1:1. Then, 80% samples are exploited for training the MDFF, 10% for evaluation, and 10% for test. By the proposed method in Section II-C.2, the object samples are automatically generated with the first category of POI data and the pixel samples.

2) Generating Object Samples: Categories of objects are required for supervising the MDFF to learn graph features. Objects (graph nodes) can be described by its main attributes. Therefore, we count POIs in each object and assign the most frequent class of POIs as the category of the object. However, POIs are distributed unevenly in space as shown in Figs. 5 and 6. Categories of objects without POIs are unknown. In addition, categories of objects may be incorrect as there is a large class imbalance in POI data. Incorrect samples will seriously affect the learning process of deep networks. Thus, object samples must be optimized. Pixel samples are the true label of the real world and thus can be used to optimize the object samples. If a graph object intersects with a sample polygon by more than 50%, its category will be changed to the category of the polygon; otherwise, its category will be maintained.

As described above, object samples need to be optimized by pixel samples. As the example shown in Fig. 7, the white circle areas in (b) are corrected in (c) with true labels of (a) and object samples become better in (c), which means the object samples have been optimized. The better samples will facilitate the network to learn more discriminative feature representations.

3) Generating UFZ Maps End-to-End: UFZs are the results of urbanization and are closely related to human activities, economic behaviors, and geographical factors. Meanwhile, land covers constitute UFZs in RS images. As a result, physical features and semantic cues are critical for mapping UFZs. The multimodal features from images and social sensing data are

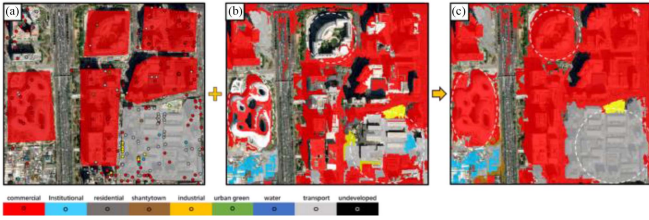


Fig. 7. Examples of optimizing object samples with pixel samples. (a) Pixel samples and POIs. (b) Object samples generated from POIs where the areas of each object is shown with its class. (c) is object samples optimized from (b) with the image polygon sample of (a).

more helpful for the interpretation. In the MDFF, the fused multimodal features from Section II-C are input into the last convolution layer (3×3 kernel, 1 stride, and 1 padding) to output the confidence map. The UFZ map is generated from the confidence map according to the maximum confidence. $\text{Loss}_{\text{pixel}}$ and $\text{Loss}_{\text{object}}$ are calculated by the cross-entropy loss function with pixel samples and object samples, respectively. The MDFF is trained end-to-end with the total loss $\text{Loss}_{\text{total}}$ which is composed of the pixel- and the object-based losses (see Fig. 3)

$$\text{Loss}_{\text{total}} = \text{Loss}_{\text{pixel}} + \text{Loss}_{\text{object}} \quad (12)$$

$$\text{Loss}_{\text{pixel}} = - \sum_{i=1}^w \sum_{j=1}^h \sum_{c=1}^n y_{ij}^c \log(p_{ij}^c) \quad (13)$$

$$\text{Loss}_{\text{object}} = - \sum_{k=1}^k \sum_{c=1}^n y_k^c \log(p_k^c). \quad (14)$$

The feature map has width of w and height of h . There are c categories that the pixel (i, j) or the object k belongs to. Prediction from the forward propagation of network is $Y_{ij} \in \mathcal{Y}$, if $Y_{ij} = c$, then $y_{ij}^c = 1$, otherwise, $y_{ij}^c = 0$.

III. EXPERIMENTS

A. Implementation Details and Evaluation Metrics

As the nodes in the graph, image objects are segmented from RS images using the multiresolution segmentation method [38] with the multiple scale of 60. The multiresolution segmentation algorithm is a commonly used image processing technique that can obtain more comprehensive and accurate segmentation results by processing images at different scales. The basic principle is achieved through steps of scale transformation, segmentation algorithms, and scale fusion. For multimodal feature representation, the advanced deep network DeepLab v3+ [43] is used as the baseline and its backbone is adopted as the backbone of the pixel-based module to extract deep features from images, and the object-based module with three layers ($L = 3$) is designed for social sensing data modeling. The RS image and training polygon labels are clipped into samples with size of 512×512 in the step of 256 pixels, where each sample has a clipped image and the corresponding label. After preparation of data, the proposed MDFF is trained with the learning rate, the batch size, and the epoch being 0.0001, 4, and 150, respectively. The

stochastic gradient descent method (SGD) and the cross entropy are adopted as the optimizer and the loss function, respectively.

In this article, the overall accuracy (OA), the intersection over union (IoU), the mean intersection over union ($MIoU$) and the frequency weighted intersection over union ($FWIoU$) are adopted as the evaluation metrics. IoU , $MIoU$, and $FWIoU$ are classic metrics, which are often used to validate intensive predictions such as semantic segmentation

$$OA = (\text{TP} + \text{TN}) / (\text{TP} + \text{FP} + \text{TN} + \text{FN}) \quad (15)$$

$$IoU_i = \frac{\text{TP}_i}{\text{TP}_i + \text{FP}_i + \text{FN}_i}, \quad i = 1, 2, \dots, n \quad (16)$$

$$MIoU = \frac{1}{n} \sum_{i=1}^n IoU_i, \quad i = 1, 2, \dots, n \quad (17)$$

$$FWIoU = \sum_{i=1}^n \left(IoU_i \cdot \frac{\text{TP}_i + \text{FN}_i}{\text{TP}_i + \text{FP}_i + \text{TN}_i + \text{FN}_i} \right) \quad (18)$$

where TP, TN, FP, and FN refer to the number of true positive points, true negative points, false positive points, and false negative points, respectively, and n is the number of classes.

B. Classification Accuracy of UFZ

The comparative experiments are set up to verify the effectiveness of the proposed MDFF. As the MDFF makes full use of UFZs' physical and socioeconomic information from multimodal geographic data, thus the MDFF is significantly better than the baseline (see Tables V–VIII), which demonstrates the advance of the proposed method. The MDFF has the reliable performance in classifying complex UFZs, e.g., commercial, institutional, and residential zones, as there are many POIs in these UFZs. It illustrates that the representation and the fusion of multimodal information improve the discriminative ability of the network. From the results of the third and the fourth rows, the MDFF performs better as it has learned the better features with the optimized labels. In addition, when trained with the Graph Loss, the MDFF achieves the better accuracy, which shows that the Graph Loss assists the MDFF to make full use of social sensing data. Therefore, the accurate label and the loss constraint guides the learning process of the graph model. With the Graph Label-OP and the Graph Loss, the MDFF obtains the best accuracy in OA , $MIoU$, and $FWIoU$.

C. UFZ Mapping Results

Figs. 8 and 9 show the final UFZ maps of Beijing and Shanghai, respectively, which is produced by the proposed MDFF with the Graph Label-OP and the Graph Loss. There are different UFZ categories: commercial, institutional, urban green, industrial, residential, shantytown, undeveloped, and transport.

In Beijing, it is shown in the UFZ map those large areas of commercial zones are in the city center. Institutional zones are mainly clustered in the center and northwest of the city for the distribution of government departments and campuses. Industrial zones and urban-green zones are located on the urban outskirts, which reflects the distribution of highly urbanized

TABLE V
OA (%) OF THE UFZ MAPPING OF BEIJING

Model	Graph Label-OP	Graph Loss	Com.	Insti.	Urgr.	Indu.	Resi.	Shto.	Water	Unde.	Tran.	Overall (OA)
Baseline	×	×	87.63	92.93	98.37	95.39	96.05	98.56	99.90	83.27	95.87	95.33
MDFF			91.37	95.57	98.78	97.93	97.88	98.37	99.79	85.55	97.23	96.95
MDFF		✓	91.77	95.25	99.15	96.80	98.09	99.19	99.98	87.81	97.76	97.18
MDFF	✓	✓	95.65	96.77	99.23	99.34	98.22	98.05	99.91	87.47	97.50	97.79

Note: The graph label-OP means that the object samples have been optimized with the object. The graph loss represents that the graph loss is used as a part of the total loss. (Com: Commercial, Insti: Institutional, Urgr: Urban Green, Indu: Industrial, Resi: Residential, Shto: Shantytown, Unde: Undeveloped, Tran: Transport.)

The significance of bold values represent the best values of the experimental results in the columns.

TABLE VI
OA (%) OF THE UFZ MAPPING OF SHANGHAI

Model	Graph Label-OP	Graph Loss	Com.	Insti.	Urgr.	Indu.	Resi.	Shto.	Water	Unde.	Tran.	Overall (OA)
Baseline	×	×	88.19	80.43	93.63	87.27	95.69	89.95	91.94	80.93	92.64	92.01
MDFF			88.61	82.86	95.02	92.38	96.37	96.27	91.81	79.23	93.85	93.51
MDFF		✓	89.88	82.27	94.66	92.43	96.78	96.20	92.64	81.65	93.43	93.69
MDFF	✓	✓	91.94	87.04	95.30	92.02	97.29	93.41	92.26	86.59	94.36	94.45

Note: The graph label-OP means that the object samples have been optimized with the object. The graph loss represents that the graph loss is used as a part of the total loss. (Com: Commercial, Insti: Institutional, Urgr: Urban Green, Indu: Industrial, Resi: Residential, Shto: Shantytown, Unde: Undeveloped, Tran: Transport.)

The significance of bold values represent the best values of the experimental results in the columns.

TABLE VII
IOU (%) OF THE UFZ MAPPING OF BEIJING

Model	Graph Label-OP	Graph Loss	Com.	Insti.	Urgr.	Indu.	Resi.	Shto.	Water	Unde.	Tran.	MIoU	FWIoU
Baseline	×	×	79.56	88.20	97.13	82.75	90.57	95.36	99.68	76.78	94.36	89.38	91.22
MDFF			87.05	92.10	97.98	88.29	93.66	96.59	99.75	81.21	96.56	92.58	94.14
MDFF		✓	87.70	92.66	98.23	90.12	94.14	96.43	99.89	83.16	96.77	93.24	94.57
MDFF	✓	✓	91.56	95.10	98.61	87.19	95.59	96.65	99.68	82.84	96.86	93.79	95.73

The significance of bold values represent the best values of the experimental results in the columns.

TABLE VIII
IOU (%) OF THE UFZ MAPPING OF SHANGHAI

Model	Graph Label-OP	Graph Loss	Com.	Insti.	Urgr.	Indu.	Resi.	Shto.	Water	Unde.	Tran.	MIoU	FWIoU
Baseline	×	×	80.69	70.83	89.50	82.80	89.35	76.57	87.92	65.70	84.91	80.92	85.34
MDFF			83.59	74.16	90.66	89.86	91.02	78.31	89.41	70.07	88.34	83.94	87.94
MDFF		✓	83.61	75.53	91.27	89.44	91.05	81.46	89.61	67.74	88.51	84.25	88.23
MDFF	✓	✓	86.74	82.03	91.47	89.91	92.85	84.87	88.16	69.60	88.29	85.99	89.56

The significance of bold values represent the best values of the experimental results in the columns.

areas. In addition, shantytown zones cluster in the urban center because of the historical legacy (Hutong). Hutong is the unique building of Beijing.

Shanghai is the economic center of China, with a highly developed economy. Commercial zones are mainly clustered in the urban center in Shanghai. In addition, Shanghai’s industry is highly developed, especially the automobile industry and the semiconductor industry. It is shown in Fig. 9 that there are many industrial zones in Shanghai. Same with Beijing, shantytown

zones in Shanghai cluster in the urban center for the historical legacy.

To evaluate the fine-grained UFZ mapping of the MDFF, the detailed UFZ maps of eight areas in Beijing and Shanghai are shown in Figs. 10(a)–(h) and 11(a)–(h), respectively. These maps overlap on the original RS images with 40% transparency. The proposed MDFF interprets multimodal geographic data accurately and achieves good performance in classify UFZs, especially for commercial, residential, institutional, industrial,

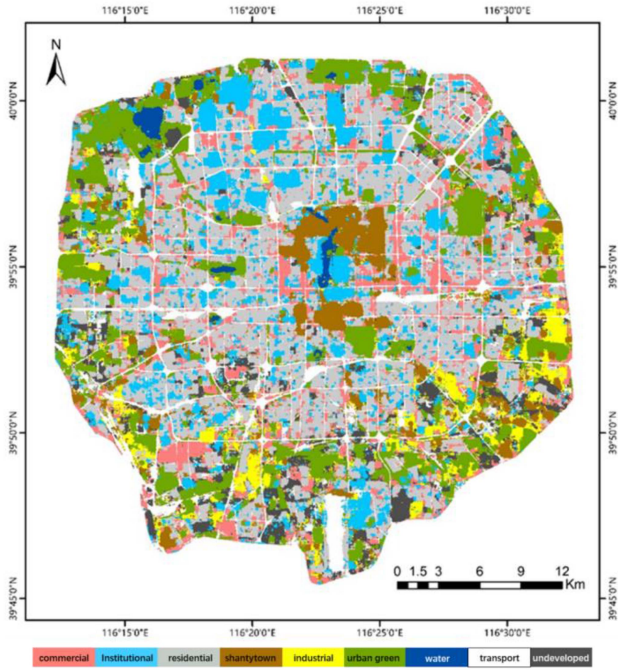


Fig. 8. UFZ map of Beijing.

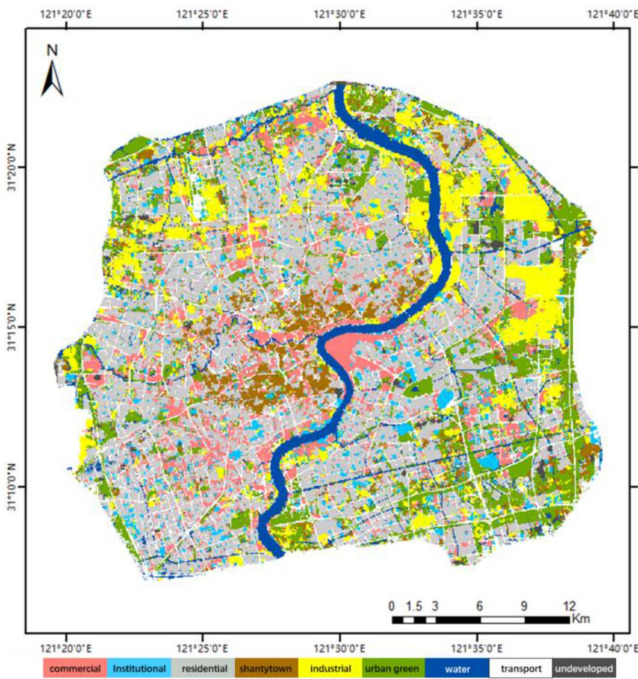


Fig. 9. UFZ map of Shanghai.

shantytown, urban green, and water. The classification results match well the real status of UFZs.

D. Urban Functional Structure Analysis

The urban functional structures reflect the development status of cities [44]. Proportions of diverse UFZs in the main urban area of Beijing and Shanghai are calculated from the UFZ results of the proposed MDFF. As shown in Fig. 12, residential districts of

the two city have the largest proportion about 30%. This is a normal phenomenon in big cities, because the population brought by urbanization requires many dwellings to carry. As the capital of China, Beijing is the political and cultural center of China, where there are many government agencies and research institutes (12% institutional zones). As the economic center of China, Shanghai has the thriving economy reflected by the commercial zones with 11%. The proportion of transport shows that ground transportation of the two city is well developed. However, urban greening rate (15% of Beijing and 13% of Shanghai) needs to be improved. In addition, there is 7% undeveloped area left in Beijing for further construction. Industrial zones are mainly distributed in the suburbs, so there are few industrial zones in the urban area. Due to the developed automobile and semiconductor industries in Shanghai, Shanghai has more industrial zones than Beijing in downtown. As a legacy of history, shanty towns still account for 5% in Beijing and Shanghai.

IV. DISCUSSIONS

A. Comparing With Existing Methods

Some representative studies for mapping UFZs are listed in Table IX, such as traditional methods based on image tiles and blocks [8], [44], [45], hierarchical semantic cognition models [25], [39], and CNN-based approaches [9], [46]. Although the adopted evaluation methods are different, these studies are more likely consistent in the study area for comparison. Traditional methods are often limited because of insufficient feature representation. The CNN-based methods have obvious advantages in UFZ mapping. Du et al. [9] designed a CNN for feature extraction and adopted the conditional random field to regroup objects for the large-scale and fine-grained UFZ maps. These methods achieve good performance but still have limitations: 1) multimodal geographic data are not fused to eliminate uncertainty of the interpretation, and 2) the two-stage methods (CNN + postprocessing) are not the end-to-end mechanism, which increases the complexity of interpretation and can lead to accumulation of classification errors. The proposed MDFF overcomes the above limitations and improves the accuracy (about 5% in OA) and reliability of UFZ mapping.

B. Advantages of Fusing Multimodal Geographic Data

Data-driven learning strategies are often unstable, which makes deep networks susceptible to noise, especially for RS images. The complex spectrum and spatial structure greatly increase the uncertainty of classification. Multimodal geographic data contain information (e.g., physical, and socioeconomic information) obtained from different perspectives, which can give the classifier more discriminative basis to enhance the reliability of the classification. With all the multimodal geographic data shown in Tables X and XI, the best classification accuracy is achieved (about 3% higher than that of single modality with images, Tables X and XI), and POIs are more important than mobile signal data and social media data in UFZ mapping, because POIs are more closely related to socioeconomic semantics

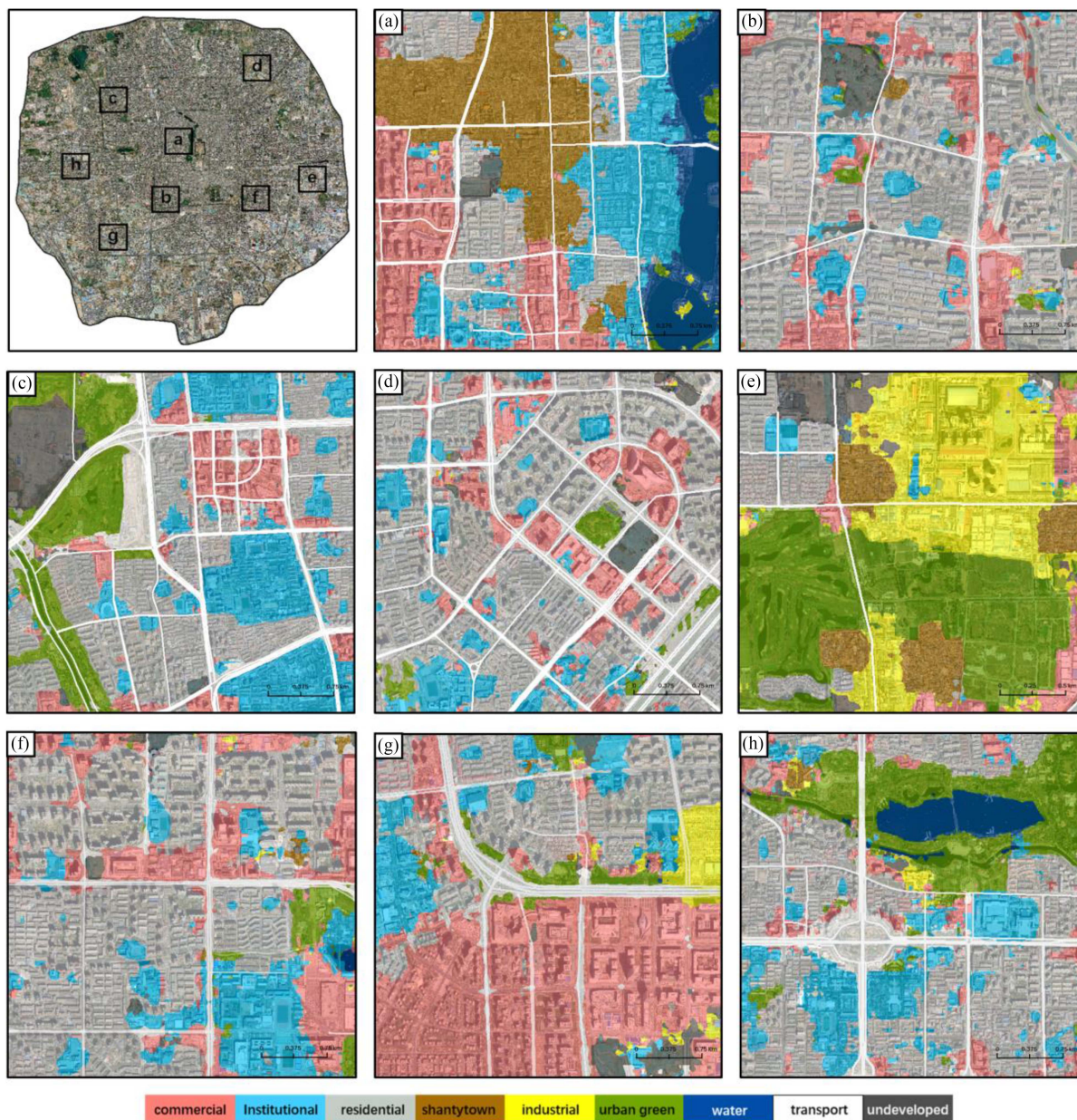


Fig. 10. Detailed UFZ maps of eight areas in Beijing.

of UFZs. In addition, mobile signal data and social media data are sparser, resulting in local lack of semantics.

However, effectively fusing multimodal geographic data and features at the high level is difficult due to the different distributions and structures of the data. In the MDFF, two feature fusion structures are designed to fuse image data and social sensing data at pixel scale and object scale. The MDFF effectively fuses multimodal geographic data and generates accurate and reliable classification results. Compared with (a) and (b) in Fig. 13, the UFZ results from the MDFF with images and social sensing data are more refined and accurate. The detailed results of (c) and (d) shows that the MDFF accurately identifies UFZs, especially in commercial, shantytown, industrial, water, transport. As shown in the red circles of (c) and (d), the MDFF corrects misclassifications with the help of socioeconomic semantics of UFZs from

POI data, which demonstrates the effectiveness of multimodal geographic data fusion.

C. How to Select Multimodal Geographic Data for UFZ Mapping

RS images carry the spectral and spatial properties of ground objects, which characterize UFZs in detail. In addition, RS images hold advantages in representing UFZs, because of the large spatial coverage and wide availability. UFZs are closely related to socioeconomic activities, thus socioeconomic attributes are very important for UFZ mapping. Socioeconomic semantics from social sensing data are helpful for interpreting UFZs. Three typical social sensing data, i.e., POIs, mobile signal data, and social media data, have different characteristics. POIs record



Fig. 11. Detailed UFZ maps of eight areas in Shanghai.

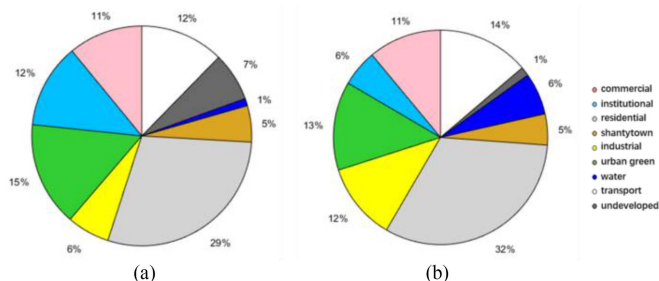


Fig. 12. Proportions of diverse UFZs. (a) For Beijing and (b) for Shanghai.

the semantics of various subjects in socioeconomic activities, while social media data generally provide static information for a single subject and mobile signal data present the spatiotemporal flow of people’s movement. Therefore, POIs provide richer and more representative socioeconomic semantics. RS images and the three typical social sensing data are selected for experiments. Experimental results in Tables X and XI show that RS images contribute the most to the interpretation of UFZs, which demonstrates that RS images are essential to the interpretation. In addition, it can be seen from Tables X and XI that POIs are more important than mobile signal data and social media data for UFZ mapping. In general, RS images and POIs are vital to the high-precision mapping of UFZs.

TABLE IX
COMPARISONS WITH EXISTING STUDIES

Model	Method	Data	Study Area	Mapping unit	Reported Accuracy
Hu et al. [8]	Traditional method	MRI POI	Beijing (16 410 km ²)	road blocks	81.0%
Simwanda and Murayama [45]	Traditional method	MHRI	Lusaka (420 km ²)	road blocks	85.0%
Zhang et al. [12]	HSC	MHRI POI	Beijing (67.1 km ²)	road blocks	90.8%
Huang et al. [44]	Traditional method	MHRI	Hong Kong (170.8 km ²)	road blocks	91.3%
Du et al. [9]	CNN+CRF	MHRI	Beijing (2267 km ²), Shanghai (2967 km ²)	Objects and road blocks	91.6% (Beijing) 89.1% (Shanghai)
Rosier et al. [46]	CNN	MHRI POI Road network	Netherlands	Pixels	85.0%
The proposed MDFF	CNN	MHRI POI Mobile signal data Social media data	Beijing (667 km ²), Shanghai (664 km ²)	Pixels and objects	97.8% (Beijing) 94.5% (Shanghai)

Note: MRI: Medium-resolution images; SHRI: Submeter-level high-resolution images; MHRI: Meter-level high-resolution images; MPPD: Mobile phone positioning data; POI: Point of interest.

TABLE X
UFZ MAPPING ACCURACY (%) OF THE MDFF UNDER DIFFERENT DATA CONFIGURATIONS IN BEIJING

POIs	social media data	mobile signal data	images	OA (%)
✓				26.58
			✓	95.33
		✓	✓	96.87
	✓		✓	96.90
✓			✓	97.77
✓	✓	✓	✓	97.79

The significance of bold values represent the best values of the experimental results in the columns.

TABLE XI
UFZ MAPPING ACCURACY (%) OF THE MDFF UNDER DIFFERENT DATA CONFIGURATIONS IN SHANGHAI

POIs	social media data	mobile signal data	images	OA (%)
✓				26.60
			✓	92.01
		✓	✓	93.74
	✓		✓	93.80
✓			✓	94.32
✓	✓	✓	✓	94.45

The significance of bold values represent the best values of the experimental results in the columns.

D. MDFF's Extensibility

As demonstrated above, the MDFF serves as a general framework for mapping UFZs, and we have evaluated the framework on two typical modal data, i.e., RS images and POI, but it has different forms according to different modal data. Lots of different modal data are available in cities, such as images (VHR, SAR, DEM, nightlight images, etc.) and social sensing data (POIs, GNSS trajectories, mobile signal data, social media data, etc.). These data characterize the physical and socioeconomic attributes of urban zones from different perspectives. For example, nightlight images provide strength of nightlight related to economic development, and mobile signal data reflect population mobility. Therefore, the MDFF can be extended with more data for better performance. There are two approaches: 1) using another data, such as replacing POIs with mobile signal

data, and 2) integrating all data (see Fig. 14) by adding branches (pixel-based module or object-based module) to extract features from data (image data or social sensing data). Then, training the MDFF according to different types of input data. Thus, the MDFF is a flexible framework which can be easily extended according to the input data.

E. Pros and Cons of the MDFF

The MDFF achieves good performance in interpreting UFZs due to the following three aspects. First, the MDFF effertely integrates multimodal geographic data and fuses multimodal features, which makes full use of physical compositions and socioeconomic characteristics of UFZs to reduce the uncertainty of interpretation. Second, unlike two-stage methods, the MDFF has only one stage and works end-to-end. It avoids error ac-

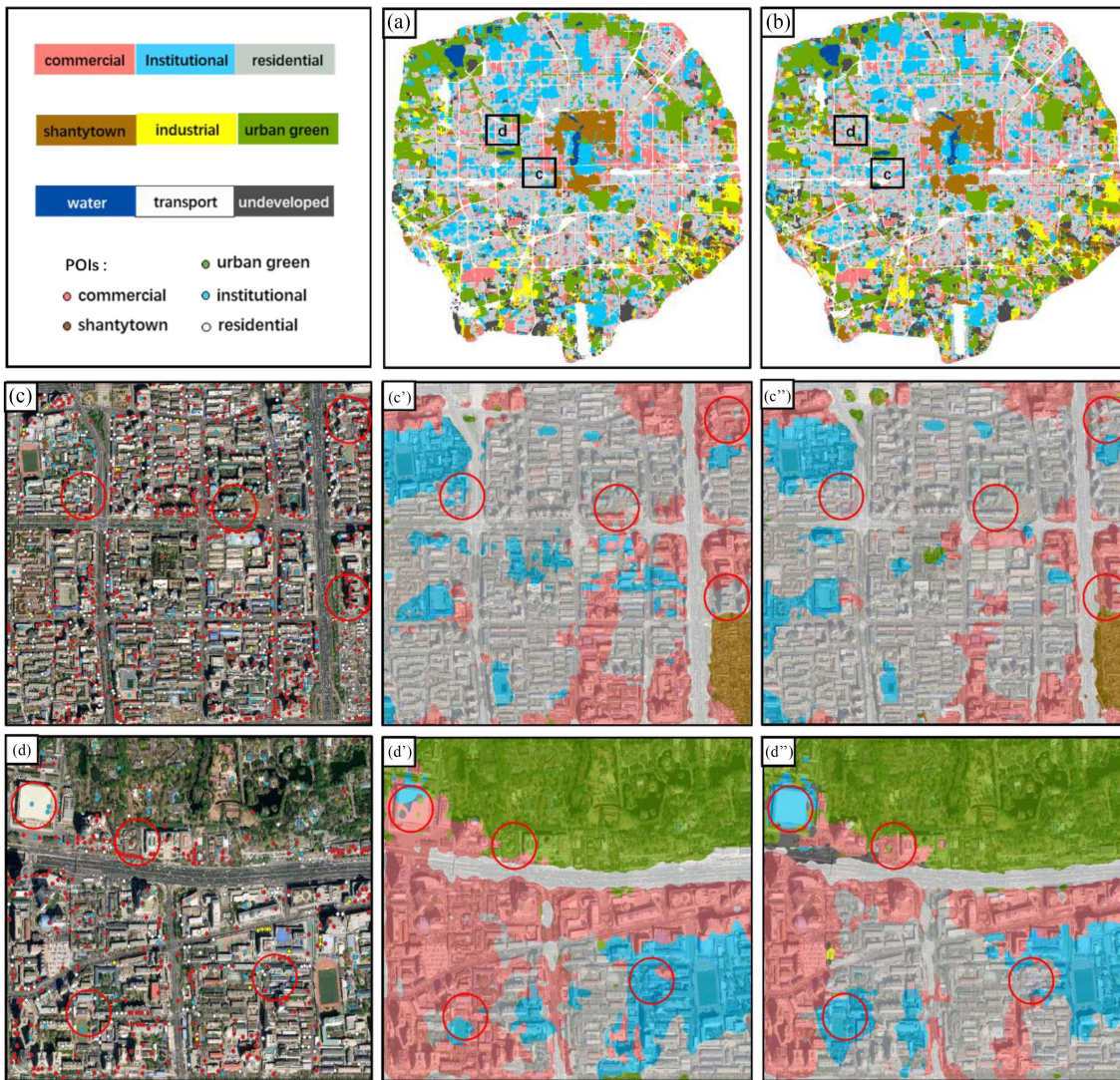


Fig. 13. Visual comparison of classification results. (a) Is the UFZ map from the MDFF only with images. (b) is generated by the MDFF with images and social sensing data. (c) and (d) are sub areas of (a) and (b). (c') and (d') are UFZ results from (a). (c'') and (d'') are UFZ results from (b).

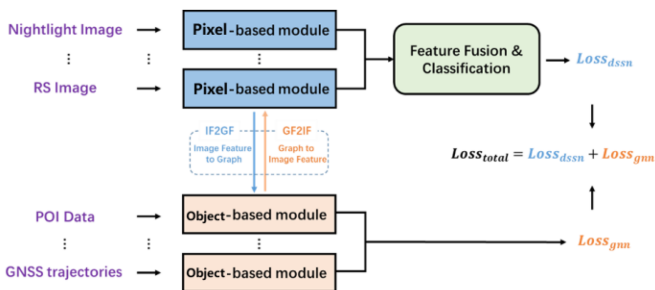


Fig. 14. Extended MDFF for integrating more data.

cumulation in stages and reduces mapping complexity, which facilitates fine-grained mapping in large urban areas. Third, the MDFF is a flexible framework which can be easily extended according to the input data, which is beneficial for dealing with unstable supply of urban data. However, the MDFF still has some limitations. First, objects generated by the multiresolution segmentation method are inconsistent with objects of the

real world, which reduces the integrity of the mapping results. Second, UFZs are usually composed of land-cover objects, while the MDFF does not take semantics of land-covers into consideration. Consequently, these two issues need to be further studied.

V. CONCLUSION

UFZ maps are essential for urban studies and applications. It is a major open problem in the field of urban geography and survey to map fine-grained UFZs accurately and reliably. In order to solve the problem, the MDFF is proposed to fuse multimodal geographic data with deep network to map fine-grained UFZs end-to-end, which effectively integrates image data and social sensing data. With two designed fusion modules (i.e., the IF2GF and the GF2IF), the MDFF fuses multimodal features to learn urban semantics. The semantic features of geographic objects assist the classification of UFZs. Experiments are conducted on two typical modal data, i.e., VHR satellite images and social

sensing data. First, the MDFF significantly improve the quality of UFZ mapping with the accuracy about 5% higher than the state-of-the-art method, which shows advantages of the end-to-end mechanism and multimodal geographic data fusion of the MDFF. Second, socioeconomic semantics from social sensing data effectively enhance the interpretation's reliability. Third, the MDFF is a flexible framework which can be easily extended according to the input data, which is beneficial for dealing with unstable supply of urban data. Fourth, RS images and POIs are vital to the high-precision mapping of UFZs. Finally, the MDFF significantly reduces the complexity of UFZ mapping and fine-grained UFZ maps from the MDFF has effectively contributed to urban structure analysis.

REFERENCES

- [1] X. Zhang and S. Du, "A linear Dirichlet mixture model for decomposing scenes: Application to analyzing urban functional zonings," *Remote Sens. Environ.*, vol. 169, pp. 37–49, 2015.
- [2] X. Zhang, S. Du, S. Du, and B. Liu, "How do land-use patterns influence residential environment quality? A multiscale geographic survey in Beijing," *Remote Sens. Environ.*, vol. 249, Nov. 2020, Art. no. 112014.
- [3] R. H. Matsuoka and R. Kaplan, "People needs in the urban landscape: Analysis of landscape and urban planning contributions," *Landscape Urban Plan.*, vol. 84, no. 1, pp. 7–19, 2008.
- [4] A. P. Montanges, G. Moser, H. Taubenböck, M. Wurm, and D. Tuia, "Classification of urban structural types with multisource data and structured models," in *Proc. Joint Urban Remote Sens. Event*, 2015, pp. 1–4.
- [5] J. C. Castella, S. P. Kam, D. D. Quang, P. H. Verburg, and C. T. Hoanh, "Combining top-down and bottom-up modelling approaches of land use/cover change to support public policies: Application to sustainable management of natural resources in northern Vietnam," *Land Use Policy*, vol. 24, no. 3, pp. 531–545, 2007.
- [6] U. Heiden, W. Heldens, S. Roessner, K. Segl, T. Esch, and A. Mueller, "Urban structure type characterization using hyperspectral remote sensing and height information," *Landscape Urban Plan.*, vol. 105, no. 4, pp. 361–375, 2012.
- [7] H. B. Shin, "Residential redevelopment and the entrepreneurial local state: The implications of Beijing's shifting emphasis on urban redevelopment policies," *Urban Stud.*, vol. 46, no. 13, pp. 2815–2839, 2009.
- [8] T. Hu, J. Yang, X. Li, and P. Gong, "Mapping urban land use by using Landsat images and open social data," *Remote Sens.*, vol. 8, no. 2, 2016, Art. no. 151.
- [9] S. Du, S. Du, B. Liu, and X. Zhang, "Mapping large-scale and fine-grained urban UFZs from VHR images using a multi-scale semantic segmentation network and object-based approach," *Remote Sens. Environ.*, vol. 261, no. 2, 2021, Art. no. 112480.
- [10] M. M. Nielsen, "Remote sensing for urban planning and management: The use of window-independent context segmentation to extract urban features in Stockholm," *Comput., Environ. Urban Syst.*, vol. 52, pp. 1–9, 2015.
- [11] M. Amani et al., "Remote sensing systems for ocean: A review (Part 1: Passive systems)," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 15, pp. 210–234, 2022.
- [12] Y. Zhang, Q. Li, H. Huang, W. Wu, X. Du, and H. Wang, "The combined use of remote sensing and social sensing data in fine-grained urban land use mapping: A case study in Beijing, China," *Remote Sens.*, vol. 9, 2017, Art. no. 865.
- [13] J. K. Steeves, G. K. Humphrey, J. C. Culham, R. S. Menon, A. D. Milner, and M. A. Goodale, "Behavioral and neuroimaging evidence for a contribution of color and texture information to scene classification in a patient with visual form Agnosia," *J. Cogn. Neurosci.*, vol. 16, no. 6, pp. 955–965, 2004.
- [14] M. Voltersen, C. Berger, S. Hese, and C. Schmillius, "Object-based land cover mapping and comprehensive feature calculation for an automated derivation of urban structure types at block level," *Remote Sens. Environ.*, vol. 154, pp. 192–201, 2014.
- [15] X. Zhang, S. Du, and Y.-C. Wang, "Semantic classification of heterogeneous urban scenes using intrascene feature similarity and interscene semantic dependency," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 8, no. 5, pp. 2005–2014, May 2015.
- [16] R. Kusumaningrum, H. Wei, R. Manurung, and A. Murni, "Integrated visual vocabulary in latent Dirichlet allocation-based scene classification for IKONOS image," *J. Appl. Remote. Sens.*, vol. 8, 2014, Art. no. 083690.
- [17] M. Lienou, H. Maître, and M. Datcu, "Semantic annotation of satellite images using latent Dirichlet allocation," *IEEE Geosci. Remote Sens. Lett.*, vol. 7, no. 1, pp. 28–32, Jan. 2010.
- [18] G. Moser and S. B. Serpico, "Combining support vector machines and Markov random fields in an integrated framework for contextual image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 5, pp. 2734–2752, May 2013.
- [19] B. Zhao, Y. Zhong, G.-S. Xia, and L. Zhang, "Dirichlet-derived multiple topic scene classification model for high spatial resolution remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 4, pp. 2108–2123, Apr. 2016.
- [20] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, 2015.
- [21] L. Ma, Y. Liu, X. Zhang, and Y. Ye, "Deep learning in remote sensing applications: A meta-analysis and review," *ISPRS J. Photogrammetry Remote Sens.*, vol. 152, pp. 166–177, 2019.
- [22] G. Mountrakis, J. Li, X. Lu, and O. Hellwich, "Deep learning for remotely sensed data," *J. Photogrammetry Remote Sens.*, vol. 145, pp. 1–2, 2018.
- [23] W. Zhou, D. Ming, X. Lv, K. Zhou, H. Bao, and Z. Hong, "SO-CNN based urban functional zone fine division with VHR remote sensing image," *Remote Sens. Environ.*, vol. 236, 2020, Art. no. 111458.
- [24] S. Du, S. Du, B. Liu, X. Zhang, and Z. Zheng, "Large-scale urban functional zone mapping by integrating remote sensing images and open social data," *GIScience Remote Sens.*, vol. 9, pp. 1–20, 2020.
- [25] X. Zhang, S. Du, and Q. Wang, "Hierarchical semantic cognition for urban UFZs with VHR satellite images and POI data," *J. Photogrammetry Remote Sens.*, vol. 132, pp. 170–184, 2017.
- [26] S. Ouyang and Y. Li, "Combining deep semantic segmentation network and graph convolutional neural network for semantic segmentation of remote sensing imagery," *Remote Sens.*, vol. 13, no. 1, 2020, Art. no. 119.
- [27] S. W. Myint, P. Gober, A. Brazel, S. Grossman-Clarke, and Q. Weng, "Per-pixel vs. object-based classification of urban land cover extraction using high spatial resolution imagery," *Remote Sens. Environ.*, vol. 115, pp. 1145–1161, 2011.
- [28] H. Bao, D. Ming, Y. Guo, K. Zhang, and S. Du, "DFCNN-based semantic recognition of urban UFZs by integrating remote sensing data and POI data," *Remote Sens.*, vol. 12, no. 7, 2020, Art. no. 1088.
- [29] T. Blaschke, "Object based image analysis for remote sensing," *ISPRS J. Photogrammetry Remote Sens.*, vol. 65, no. 1, pp. 2–16, 2010.
- [30] C. Cleve, M. Kelly, F. R. Kearns, and M. Moritz, "Classification of the wildland-urban interface: A comparison of pixel-and object-based classifications using high-resolution aerial photography," *Comput. Environ. Urban Syst.*, vol. 32, pp. 317–326, 2008.
- [31] J. M. Peña-Barragán, M. K. Ngugi, R. E. Plant, and J. Six, "Object-based crop identification using multiple vegetation indices, textural features and crop phenology," *Remote Sens. Environ.*, vol. 115, no. 6, pp. 1301–1316, 2011.
- [32] S. C. Zhu, R. Zhang, and Z. Tu, "Integrating bottom-up/top-down for object recognition by data driven Markov chain Monte Carlo," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2000, vol. 1, pp. 738–745.
- [33] H. Rhinane, A. Hilali, A. Berrada, and M. Hakdaoui, "Detecting slums from SPOT data in Casablanca Morocco using an object-based approach," *J. Geographic Inf. Syst.*, vol. 3, no. 03, 2011, Art. no. 217.
- [34] E. Farahzadeh, T. J. Cham, and W. Li, "Semantic and spatial content fusion for scene recognition," in *New Development in Robot Vision*. Berlin, Germany: Springer, 2015, pp. 33–53.
- [35] J. Yang, C. Wu, B. Du, and L. Zhang, "Enhanced multiscale feature fusion network for HSI classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 12, pp. 10328–10347, Dec. 2021.
- [36] J. Stuckens, P. R. Coppin, and M. E. Bauer, "Integrating contextual information with per-pixel classification for improved land cover classification," *Remote Sens. Environ.*, vol. 71, no. 3, pp. 282–296, 2000.
- [37] S. Du, S. Du, B. Liu, and X. Zhang, "Context-enabled extraction of large-scale urban UFZs from very-high-resolution images: A multiscale segmentation approach," *Remote Sens.*, vol. 11, 2019, Art. no. 1902.
- [38] M. Baatz and A. Schape, "Multiresolution segmentation: An optimization approach for high quality multi-scale image segmentation," *Angewandte Geographische Informations-Verarbeitung*, vol. 12, pp. 12–23, 2000.
- [39] X. Zhang, S. Du, and Q. Wang, "Integrating bottom-up classification and top-down feedback for improving urban land-cover and functional-zone mapping," *Remote Sens. Environ.*, vol. 212, pp. 231–248, 2018.
- [40] G. J. Hay, T. Blaschke, D. J. Marceau, and A. Bouchard, "A comparison of three image-object methods for the multiscale analysis of landscape structure," *ISPRS J. Photogrammetry Remote Sens.*, vol. 57, pp. 327–345, 2003.

- [41] T. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *Proc. Int. Conf. Learn. Representations*, Toulon, France, Apr. 2017.
- [42] P. Veličković, G. Cucurull, and A. Casanova, "Graph attention networks," in *Proc. Int. Conf. Learn. Representations*, Vancouver, BC, Canada, Apr./May 2018.
- [43] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, vol. 11211, pp. 833–851.
- [44] B. Huang, B. Zhao, and Y. Song, "Urban land-use mapping using a deep convolutional neural network with high spatial resolution multispectral remote sensing imagery," *Remote Sens. Environ.*, vol. 214, pp. 73–86, 2018.
- [45] M. Simwanda and Y. Murayama, "Integrating geospatial techniques for urban land use classification in the developing sub-Saharan African city of Lusaka, Zambia," *ISPRS Int. J. Geo-Inf.*, vol. 6, 2017, Art. no. 102.
- [46] J. Rosier, H. Taubenbock, P. Verburg, and J. Vliet, "Fusing Earth observation and socioeconomic data to increase the transferability of large-scale urban land use classification," *Remote Sens. Environ.*, vol. 278, 2022, Art. no. 113076.



Song Ouyang received the B.S. and M.S. degrees from the School of Remote Sensing and Information Engineering, Wuhan University, Wuhan, China, in 2018 and 2021, respectively. He is currently working toward the Ph.D. degree in cartography and geographical information system in Peking University, Beijing, China.

His current research interests include knowledge reasoning, urban landscape modeling, remote sensing, and ontology modeling.



Shihong Du received the B.S. and M.S. degrees in cartography and geographic information system from Wuhan University, Hubei, China, in 1998 and 2001, respectively, and the Ph.D. degree in cartography and geographic information system from the Institute of Remote Sensing Applications, Chinese Academy of Sciences, Beijing, China, in 2004.

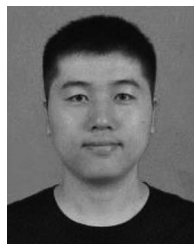
He is currently a Professor with Peking University, Beijing. His research interests include urban remote sensing, machine learning, object-based image analysis, time series analysis, change detection and dynamic analysis of land cover/use, urban environmental information extraction, and applications of big earth data in support of landscape change analysis and urban sustainable development.

and applications of big earth data in support of landscape change analysis and urban sustainable development.



Xiuyuan Zhang received the B.S. degree in remote sensing from the China University of Geosciences, Wuhan, China, in 2014, and the Ph.D. degree in geographic information system from Peking University, Beijing, China, in 2019.

He is currently a Research Associate with the Institute of Remote Sensing and GIS, Peking University. His research interests include remote sensing and GIS for urban landscape modeling and urban sustainable development monitoring.



Shouhang Du received the B.S. degree in surveying and mapping engineering from the China University of Geosciences (Wuhan), Wuhan, China, in 2014, the M.S. degree in photogrammetry and remote sensing from the Central South University, Changsha, China, in 2017, and the Ph.D. degree in cartography and geographic information system from the Peking University, Beijing, China, in 2021.

He is currently a Lecturer with the China University of Mining and Technology (Beijing), Beijing, China. His research interests include intelligent understand-

ing of spatial data including GIS and remote sensing data.



Lubin Bai received the bachelor's degree in survey and mapping engineering from Wuhan University, Wuhan, China, in 2020. He is currently working toward the master's degree in cartography and geographic information system in Peking University, Beijing, China.

His research interests include remote sensing image processing, self-supervised learning, and representation learning.