

Article



A Complete Reinforcement-Learning-Based Framework for Urban-Safety Perception

Yaxuan Wang 1,+, Zhixin Zeng 2,+, Qiushan Li 3,* and Yingrui Deng 1

- ¹ School of Computer Science, Sichuan University, Chengdu 610207, China
- ² School of Cyber Science and Engineering, Sichuan University, Chengdu 610207, China
- ³ Institute for Disaster Management and Reconstruction, Sichuan University, Chengdu 610207, China
- * Correspondence: liqiushan@scu.edu.cn; Tel.: +86-156-809-97814

+ These authors contributed equally to this work.

Abstract: Urban-safety perception is crucial for urban planning and pedestrian street preference studies. With the development of deep learning and the availability of high-resolution street images, the use of artificial intelligence methods to deal with urban-safety perception has been considered adequate by many researchers. However, most current methods are based on the feature-extraction capability of convolutional neural networks (CNNs) with large-scale annotated data for training, mainly aimed at providing a regression or classification model. There remains a lack of interpretable and complete evaluation systems for urban-safety perception. To improve the interpretability of evaluation models and achieve human-like safety perception, we proposed a complete decisionmaking framework based on reinforcement learning (RL). We developed a novel feature-extraction module, a scalable visual computational model based on visual semantic and functional features that could fully exploit the knowledge of domain experts. Furthermore, we designed the RL module-comprising a combination of a Markov decision process (MDP)-based street-view observation environment and an intelligent agent trained using a deep reinforcement-learning (DRL) algorithm-to achieve human-level perception abilities. Experimental results using our crowdsourced dataset showed that the framework achieved satisfactory prediction performance and excellent visual interpretability.

Keywords: urban-safety perception; reinforcement learning; scene understanding; artificial intelligence; interpretability; feature extraction

1. Introduction

Safety perception is a prerequisite for happiness, health, and a high quality of life [1,2]. Low perceptions of safety are likely to influence walking preferences and psychological awareness and may trigger criminal behaviour [3–5]. Several studies have shown that the visual appearance of a city is vital to human perception and responses to the surrounding environment [6–9]. Moreover, the visual qualities of urban spaces affect the psychological state of its inhabitants [10], making it critical to understand people's safety perceptions and evaluations of urban spaces.

With the popularisation of online street views and the extensive use of machine learning techniques, previous methods of interviewing city residents and manually reviewing photographs and videotapes of a city for analysis can now be executed automatically. A typical approach is to obtain large-scale street-view data using online street-view websites, make annotations, and then use machine-learning techniques to regress and predict street-view images to obtain safety scores [7,11,12]. This approach usually requires crowdsourcing websites [6] to collect the absolute safety scores of images or the occurrence of crimes on government websites to acquire appropriate labels. However, the labels obtained through crowdsourcing often result in considerable noise, owing to factors

Citation: Wang, Y.; Zeng, Z.; Li, Q.; Deng, Y. A Complete Reinforcement-Learning-Based Framework for Urban-Safety Perception. *ISPRS Int. J. Geo-Inf.* 2022, *11*, 465. https://doi.org/ 10.3390/ijgj11090465

Academic Editors: Wei Huang and Wolfgang Kainz

Received: 20 July 2022 Accepted: 25 August 2022 Published: 29 August 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https://creativecommons.org/licenses/by/4.0/).

2 of 16

such as an insufficient number of comparisons [13] and the random actions of users. Consequently, they cannot be treated as actual labels, which means that the subjectivity of the prediction results of the above models cannot be avoided, possibly leading to low confidence in the results.

Unlike previous studies, this study proposes a complete evaluation framework to predict human perceptions of safety from geotagged images of urban spaces. This is the first study to make high-level judgements on urban safety by combining visual–semantic features and reinforcement learning (RL) methods. Referring to [14–16], we modelled the problem of urban-safety perception as a Markov decision process (MDP) problem [17] and then used the deep reinforcement-learning (DRL) algorithm to solve it to obtain a human-like perception policy on urban safety. It was experimentally established that the RL-based framework outperformed previous regression and convolutional neural network (CNN)-based methods. By simulating the sensory basis of humans in decision making, we could visualise safety perception evaluation patterns and form a reward function for the RL process to guide the output of the model. Using this new framework, the construction process of subjective perception could be analysed at a deeper level. In addition, our framework permitted the quantitative analysis of the image's visual–semantic features, which correlated with the concept of perceived safety.

2. Review of Related Fields

This chapter describes critical work in three related research areas—the safety perception of the city, interpretable scene understanding of street views, and decision making.

2.1. Safety Perception of the City

The perception of the city—which refers to the recognition and feeling of a particular place—forms part of the city's imagery. Lynch proposed the concept in The Image of the City [18] in 1960, that is, imagery results from the interaction between an observer and the objects being observed, which refers to people's overall perceptions and impressions of the city. In her book, Jane Jacobs proposed the "street eye" theory [19], where she discussed the street's role as a city's main visual scene and its impact on human perception. The theory states that safety is an essential function of streets. As a pioneer in the study of city perception, her work inspired later researchers in the field of city planning and safety perception. In contrast to objective safety, safety perception is a subjective feeling that is usually associated with the fear of crime and potential danger [20]. A low level of safety perception affects human behaviour, which may lead to further negative consequences. Li et al. [21] found that a low level of safety perception could lead to a greater reluctance to engage in physical activity among older adults. The fear of crime has also been reported [22,23] to be negatively associated with psychological and physical health, making the study of city-safety perceptions of great importance.

In recent years, several studies have analysed the physical elements of street views that have different levels of impact on safety perception. Liu et al. [24] found that geometric scenes in street views generally produce a higher level of safety perception than naturalistic scenes. Li et al. [25] found that, for different types of terrain, the visibility of green vegetation plays a vital role in improving the safety perception in city areas. Recently, several studies have proposed metrics and descriptors to understand and represent the visual perception of streets. Cheng et al. [8] proposed four inducers to characterise streets' visual perception—namely, significant area saturation, visual entropy, the green view index, and the sky openness index. Zhang et al. [9] proposed a local representation framework using scene elements consisting of street scene ontology and street visual descriptors for the qualitative understanding and quantitative representation of street scenes, respectively.

3 of 16

Currently, most studies focus only on the impact of one or two of the perceptual factors associated with city construction, residential areas, or crime [3]. However, they do not propose a general approach that considers the impact of both urban functional areas and street features on safety perception. To the best of our knowledge, ours is the first quantitative representation model for street-scene perception.

2.2. Interpretable Scene Understanding

Before deep learning was widely applied, researchers often conducted mathematical and statistical analyses based on existing statistics [26] or designed questionnaires to collect data [27] to understand scenes. Both approaches are costly and do not directly examine the visual scenes. Advances in computer technology have made it feasible to conduct research directly based on the perceptual features of vision. Moreover, computer vision methods have been shown to emulate human perception and reliably predict safety judgements using pictures of urban scenes [7,28,29]. Visual features have been widely used to analyse their impact on security perception [11,12,30].

Quercia et al. [6] used a crowdsourcing approach to annotate street-view images and investigate the role of visual features in perceptual predictions. Arietta et al. [3] proposed a support vector regression model to predict non-visual attributes (crime rate or house prices) from images and automatically determined the visual elements associated with the predicted attributes. Naik et al. [7] designed the Streetscore algorithm (training it on the Place Pulse 1.0 [31] dataset), which could predict the safety perception score of street-view images. However, the algorithm focused on only a few specific features, and the limitations of the dataset made it less effective, with it failing to predict non-typical street-view images with unusual elements.

Naik et al. [11] proposed an Streetscore-CNN and Ranking Streetscore-CNN to study the safety perception problem on a larger scale, achieving more accurate results than Streetscore. Since then, CNNs have been increasingly applied in this field. For instance, Liu et al. [12] proposed a unified framework to quantify the perceived attributes (e.g., the level of safety) of a city's physical environment and introduced CNNs to parameterise instance-level scoring functions. This method could generate region-level safety scores to interpret the perception process.

However, the problem of identifying visual elements that correlate with high-level semantic visual attributes has rarely been addressed, even though image semantic segmentation is critical for understanding scenes [9]. A recent study [4] used a semantic segmentation technique to obtain visual–semantic features before assessing the impact of micro-built environmental variables on drug activities. However, they did not analyse the effect on safety perception. These studies used either CNN [11,28] or regression methods [7,12] to analyse or make predictions. Such an approach can be prone to overfitting on smaller datasets [32]; it can also fail to properly understand the inherent patterns of perception evaluation performed by humans under such conditions. In addition, the data obtained through crowdsourcing have considerable noise, which can negatively impact the building of a theoretical streetscape safety perception framework.

2.3. Decision Making

Decision making lies at the core of control theory [14] and is widely used in autonomous driving [15], robot navigation, and human-level control in gaming [16]. In these applications, agents are designed to interact with the environment to observe and execute specific actions to fulfil predefined goals. The RL method [33–36] is well-suited for handling decision-making tasks. You et al. [15] used RL to achieve the desired driving behaviour in autonomous-vehicle planning problems and successfully applied the RL method to decision-making problems.

To the best of our knowledge, no decision-making framework has been applied to city perception or interpretable street-view analysis. This study proposes a training method using DRL algorithms driven by visual features, functional features, and expert knowledge. Our method performed well in this evaluation.

3. Methods and Data Processing

In our study, the problem of urban-safety perception was abstracted as a decisionmaking problem and represented using the MDP. In this process, we form a computational model to represent the elements of the MDP and design the reward function. DRL algorithms can then be used to solve the problem of obtaining a safety evaluation strategy for street-view images, thus forming an urban-safety perception decision-making framework (Figure 1) based on visual–semantic features, functional features, and objective safety scores.



Figure 1. The proposed decision-making framework for urban-safety perception prediction consists of two main modules: a feature-extraction module to extract visual semantic and functional features, as well as safety score obtained from crowdsourcing; and a reinforcement learning module to obtain the safety perception prediction.

This section introduces a novel semantic-segmentation-based image perceptual feature computational model combined with functional and visual–semantic features. The exact system-modelling method is presented in Section 4. First, we introduce the process of collecting Dujiangyan street images and designing a crowdsourcing website to obtain the safety score of the images using the learning-to-rank (LTR) method (Section 3.1). We then describe the application of expert knowledge in this study (Section 3.2).

3.1. Building the Dataset

Following Naik et al. [11], a crowdsourcing data-collection website was designed to obtain annotations for the datasets used in this study.

Building the Dataset. The original dataset used in this study was based on street-view images of Dujiangyan City provided by the Baidu street-view API. The dataset included 11,584 panoramas, each panorama being sliced into four images corresponding to the four angles of observation (-90°, 90°, 180°, 270°). We also obtained the latitude–longitude values of the images through the GPS satellite positioning of the street car that took the pictures. Owing to the data-collection method of the Baidu street-view platform, the original dataset had a high level of redundancy, which was not convenient for annotation and analysis. Therefore, we filtered and constructed a smaller dataset for this study. Based on the latitude–longitude values, we obtained the specific street locations corresponding to each image and built the dataset by randomly filtering the images according to the street

locations (Figure 2a). In this process, we ensured that the images were geographically averaged, which meant that our dataset reflected the overall appearance of Dujiangyan City. In contrast to [11], we concatenated the images of the four orientations in a two-by-two sequence, thus simulating the horizontal view of the human eye [37]. We used panoramas for the crowdsourcing data annotation to obtain better subjective safety indicators.



Figure 2. (a) The distribution of locations of selected street view images in Dujiangyan city. (b) Examples of normalised safety scores (represented as Q) using Trueskill and LambdaRank.

Design of the Crowdsourced Website. We designed a website to collect pairwise comparisons of pictures. When visiting our website, users were shown a random pair of images from our dataset and asked to click one image in response to the question, "Which picture leaves you with a safer impression?" The user choices were collected for comparison.

Obtaining Safety Scores. Gathering relative comparisons is a more efficient and accurate method of obtaining human rankings than obtaining numerical scores from each user [11]. Based on this, we recorded the number of user-clicks as a pairwise comparison of a set of images. From a single click, we used a 3-tuple (-1, 0, +1) to record the user's choices, +1 denoting a "win" for the image, -1 denoting a "fail", and 0 denoting a "tie" for both images. We also recorded a number of randomly selected images and used these to dynamically adjust the extraction strategy to ensure that the number of annotations for each image was comparatively balanced. TrueSkill [13] and LambdaRank [38] were used to convert the participants' clicks (preferences) into objective scores. We used the average results of the two algorithms as the final safety scores for the images (Figure 2b). The Trueskill algorithm is designed based on a Bayesian graphical model to evaluate the skill level of game players. In our case, a sufficient number of user clicks could determine the "superior" image (i.e., the "winner"). The LambdaRank algorithm transforms a ranking problem into a pairwise classification problem, directly defining the gradient using its physical meaning in the ranking problem and then solving the gradients for the final ranking. LambdaRank can fully utilise an image's "win" or "fail" to rank it as a feature.

3.2. Expert Rating

To model this problem, the selection of visual features in the dataset and the construction of an expert system for safe reasoning require reliable expert knowledge, the use of which is a simple and efficient method [39]. Although subjectivity and bias cannot be eliminated, the reliability of expert knowledge is acceptable. It is practical when the actual properties of an environment can be consulted and the objectivity of the data evaluation can be adequately guaranteed [40].

To obtain expert knowledge, 15 researchers and professors in related fields were invited, i.e., sociology, urban planning, landscape architecture, and computer vision. Specifically, as a panel of experts for this project, they were asked to complete the following four tasks on the dataset.

- Low-quality images were removed, such as those with tunnels or few elements. Following [41], the most representative features in the street-view images were selected.
- Criteria were proposed for zoning based on the features and actual geographical location of each functional area in the city. We describe zoning in Section 3.3.
- A corresponding expert system for perceptual safety prediction based on different functional areas was designed; this expert system was the basis of our RL method (for both the reward function and the state definition).
- The safety scores were amended to reduce the uncertainty and set the score threshold—that is, images above the threshold were considered safe and labelled "1". otherwise, the label was "0". It is worth noting that, even after correction by experts, noise still existed in the labels.

Due to the small size of the filtered dataset and to further incorporate the objectivity and comprehensiveness of expert knowledge into the model, we invited experts to evaluate the entire dataset as much as possible. Fifteen experts were needed to consider 1142 images each. We used the voting method to determine the final labels for those images with different opinions from experts. We also used statistical methods to analyse the data to reduce the potential bias for the different types of features proposed by the experts. After the expert rating, 651 images could be used for further analysis

3.3. Feature Extraction

Based on expert knowledge (Section 3.2), we used semantic segmentation techniques to extract features from street-view images. To model the problem, we defined two types of features, that is, visual and functional features [42]. In addition, we used the safety score obtained in Section 3.1 as a third type of feature.

Visual features. In this study, we define visual features as those that can represent the overall characteristics of the street-view image (Table 1). Combining previous literature and expert analyses, we selected 11 visual features in 3 categories, that is, (1) Field of view (FoV), which is the cover ratio of a specific scene element in the field of view [9]; (2) Visual entropy, which is the result of combining the concepts of information entropy with the characteristics of the human visual system [8], reflecting images' visual complexity and richness; and (3) Tiny objects, which are the small-scale objects in the street-view images after semantic segmentation, such as poles, electric wires, and the number of vehicles and people.

Feature Name	Description		
Sky-FoV	The cover ratio of the sky in the field of view		
Greenery-FoV	The cover ratio of the terrain and vegetation in the field of view		
Wall-FoV	The cover ratio of the wall in the field of view		
Sidewalk-FoV	The cover ratio of the sidewalk in the field of view		
Building-FoV	The cover ratio of the building in the field of view		
Traffic_light-FoV	The cover ratio of the traffic light in the field of view		
Traffic_sign-FoV	The cover ratio of the traffic sign in the field of view		
Visual Entropy	The magnitude of the visual entropy value can reflect the visual complexity and		
	richness of an image		

Table 1. Specific definitions of street view visual features.

Vehicle Number	The number of the vehicles
Person Number	The number of the person
Electric Wire	Whether there is any electric wire. The value is 0 or 1.

To calculate the FoV, we used the semantic segmentation technique to segment the city scenes into scene elements and calculate their coverage ratio, that is, the pixel proportion (Figure 3). The scene elements included the sky, vegetation, terrain, poles, buildings, walls, fences, roads, sidewalks, traffic signs, and traffic lights. SegFormer [43], OCRNet [44], and DeeplabV3 + [45] were used to obtain the best results. All techniques were trained using the Cityscapes dataset [46]. We eventually applied the SegFormer model, which combines a transformer with a lightweight multilayer perceptron (MLP) decoder. No positional coding was required, thus avoiding interpolation, which can lead to performance degradation when the test resolution differs from the training. It was also proposed that the MLP decoder aggregate information from different layers, combining local and global attention for semantic representation.

Original Image

OCRNet

OCRNet

SegFormer

Figure 3. Results of different semantic segmentation algorithms.

Functional features. The city can be seen as a combination of many different functional areas [42], while a specific geographic location can impact safety perception. For different functional areas (Figure 4), we defined functional features as visual features that could fully express the characteristics of their locations. It should be noted that, in the feature vector representation, the functional features were reflected in ordinal numbers, with different numbers representing different functional areas.



Industrial Area



Business Area Suburban Area

Feature vector representation. Considering that the original visual features have different ranges of values and that using normalisation can lead to a loss of their original meaning [42], we calculated the threshold value for each visual feature individually and then mapped the features based on the threshold to obtain binary values. The following expression can be used to construct binary values, x_i :





$$x_i = \theta \cdot \max(\operatorname{sgn}(f_i - t_i), 0) + (1 - \theta) \cdot \min(1 - \operatorname{sgn}(f_i - t_i), 1).$$
(1)

In the above expression, x_i denotes the *i*-th binary visual feature value, f_i denotes the *i*-th visual feature value, and t_i denotes the threshold of the *i*-th feature. The value of θ denotes whether the feature value is positively correlated with perceived safety, that is, the value selection of θ is 0 or 1.

Unlike in [42], we use *X* to represent both functional and visual features after mapping to binary values. *X* can be represented as a feature vector obtained from an image, as follows

$$X = \left\{ y_j, x_j^1, x_j^2, ..., x_j^{f(y_j)} \right\},$$
 (2)

where y_j denotes the *j*-th functional area in Table 2, and $f(y_j)$ denotes the number of selected visual features in the functional area y_j .

Table 2. The visual features correlated to different functional areas. Each area also has a safety score.

No. (<i>j</i> -th)	Functional Area	Features
1	Business Area	GVI, Wall-FoV, Traffic_light-FoV, Traffic_sign-FoV, Electric Wire, Sky-FoV,
		Building-FoV
2	Cultural Area	GVI, Wall-FoV, Electric Wire, Sky-FoV, Person Number
3	Residential Area	Wall-FoV, Electric Wire, Building-FoV, Sky-FoV, Visual Entropy, GVI
4	Industrial Area	Electric Wire, Wall-FoV, Sidewalk-FoV, Vehicle Number, GVI
5	Suburban Area	Electric Wire, Wall-FoV, Sky-FoV, Visual Entropy
6	Others	Visual Entropy, Electric Wire, Building-FoV, GVI, Wall-FoV, Sky-FoV

4. System Modelling

Combining RL and deep learning (DL), DRL can be used to solve complex decisionmaking problems, making it suitable for simulating human perceptual patterns and obtaining safety perception predictions of street-view images. Specifically, we modelled the interaction between experts and street-view images as a stochastic MDP and used DRL methods to solve the MDP problem for the optimal policy that could predict the perceptual safety of given images.

In this section, we first briefly introduce the MDP as the necessary theoretical support for modelling RL. We then describe how to model the evaluation strategy using visual and functional features. Finally, we demonstrate how to solve the decision-making problem using the DRL method. Our framework can be applied to any dataset, outputting human-like perceived safety predictions, which could help improve urban design and planning efficiency.

4.1. Markov Decision Process

The Markov decision process (MDP) is a discrete-time stochastic process that provides a mathematical framework for modelling decision-making problems. The MDP framework can model the interaction between the agent and environment for most RL problems. In this interaction process, the agent performs actions based on policies to obtain as many rewards as possible, which is the goal of the MDP. RL methods need to be applied to find the best policy.

A typical MDP model can be represented by a quadruplet (S, A, P_a, R_a) , where *S* denotes the set of all possible states in the environment, that is, the state space (finite set); *A* denotes the set of all possible actions, that is, the action space (finite set); P_a denotes the state transition probability; and R_a denotes the immediate reward after the state transition. In the MDP case, the cumulative reward sum that the agent can obtain is

 $E\left[\sum_{t=0}^{\infty} \gamma^{t} R_{at}(s_{t}, s_{t+1})\right]$, where $\gamma \in [0,1]$ denotes the discount factor that enables the policymaker (agent) to consider the impact of current and future rewards on the overall reward. We use the function $\pi(s)$ to denotes the action that the agent performs in that state and π^{*} to denote the best RL strategy. The following equation describes the goal of RL.

$$\pi^* = \arg\max_{\pi} E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t))\right]$$
(3)

4.2. Image-Based Modelling

Inspired by [47] and the traditional model of an expert system, we designed an image-based model to predict the safety perception of a city. Our decision-making model can be regarded as an agent that interacts with the environment, the environment being a given image (*I*), the feature vector of *I* defining the state, and the action being to predict whether the street-view image is safe or not as $A \triangleq \{0,1\}$.

State Definition. Considering both the visual and functional features of each streetview image, we defined the state of the MDP, as shown in Figure 5. According to the expert criteria, each image was classified into the corresponding functional area. The states of the MDP each represent one of the six conditions listed in Table 2. It is worth noting that the features corresponding to each functional area are selected based on expert knowledge, so subjectivity cannot be entirely avoided. For a given functional area y_i , the

total number of possible states for that given functional area is $2^{f(y_i)}$. Thus, the total number of states for all functional areas is 336.



Figure 5. The Markov decision process of our framework. The neural network agent (**left**) learns to exploit the visual state (**bottom**) to predict the action, which is the perceived safety indicator (**top**). Then, the agent explores the environment (**right**) to maximise the total reward (**middle**).

Reward Function Design. The design of the reward function was based on expert knowledge. The reward function directly determines the performance of a model. In our case, this is represented by mapping manually selected features. These features depend on the environmental state. A linear combination of features can be used to represent the reward function [15], which can be expressed as follows.

$$R(s,a) = \omega^T (\Phi(s,a) + Q), \qquad (5)$$

where ω denotes the weight vector, Q denotes the safety score of the image, and $\Phi(s,a)$ denotes the feature vector, with each component representing a single feature point in the state-action space. In this study, we define $\Phi(s,a)$ as follows.

$$\Phi(s,a) = \operatorname{sgn}(X \oplus a - 0.5) \tag{6}$$

The conditions are as follows: (1) the impact of features on safety perception varies across functional areas, and (2) there are different features for different functional areas. A higher reward was given for the correct step to ensure that the agent learned an effective strategy. We designed the weight vector ω to reward or penalise specific features in the reward function. Consequently, the DRL algorithms could be applied to learn the optimal policy by maximising the total reward. According to [8,9], safety perception is highly associated with features that have a clear physical meaning, so selecting features and designing weight vectors is based on expert knowledge.

Policy Learning. After designing the reward function of the corresponding MDP, we determined the optimal decision policy for the agent using DRL, which can solve complex, previously intractable, real-world decision-making problems.

In our environment, the actions are discrete. Consequently, we used the DRL algorithm D3QN, which can be applied to a discrete action space. The algorithm combines the advantages of a Double Deep Q-Network (DDQN) [48] and a duelling DQN [49]. A DDQN can train two Q networks simultaneously to reduce overestimation bias; a Duelling DQN uses the advantage function to accurately estimate the Q values of the states. The D3QN is a model-free DRL algorithm that iteratively solves the Bellman equation. The following equation formulates the algorithm

$$Q_{\pi}(s,a) = R_{s}^{a} + \gamma \sum_{s' \in S} P_{ss'}^{a} \left[\sum_{a' \in A} \pi(a' \mid s') Q_{\pi}(s',a') \right],$$
(7)

where $Q_{\pi}(s, a)$ denotes the state-action-value function, P_{ss}^{a} denotes the probability that the agent reaches *s*' from *s* by performing action *a*, and $\pi(\cdot)$ denotes the agent's policy.

Similar to the DDQN, the dual-network structure of the D3QN consists of an evaluation network Q_1 and a target network Q_2 . Q_1 selects an action for the next step using the argmax function, and Q_2 is used to reduce the overestimation of the argmax function. D3QN trains the networks simultaneously, and a smaller Q value is selected to calculate the temporal difference error, thus reducing the overestimation bias.

$$\begin{cases} a_{\max} = \arg\max_{a} Q_1(s', a; \theta) \\ y = r + \gamma Q_2(s', a_{\max}; \theta^-) \end{cases}$$
(8)

Based on the Duelling DQN [49], D3QN separates the action-value function $Q_{\pi}(s,a)$ into two parts—namely, a state-value function $V_{\pi}(s)$ (which is only related to the state) and an advantage function $A_{\pi}(s,a)$ (which is related to both states and actions).

$$Q_{\pi}(s,a;\theta,\alpha,\beta) = V_{\pi}(s;\theta,\alpha) + A_{\pi}(s,a;\theta,\beta)$$
⁽⁹⁾

Summary of DRL Contributions. In this study, the existing D3QN algorithm was adapted to real-world applications, bridging the gap between simulation-based DRL applications and challenging real-world tasks. Moreover, a novel reward function that could

consider the visual and functional features for evaluating city-safety perception was developed to guide D3QN-based model training.

5. Experimental Results

5.1. Experimental Setup

RL Environment. Based on the visual and functional features mentioned in Section 3.3 and the designed reward function (5), we could design an environment for training the agent. It is worth noting that the safety scores computed using Trueskill and LambdaRank were also used as features for the reward function. The gym library was used to design an RL environment.

Implementation Details. We used the DRL algorithms PPO, A2C, SAC, and D3QN for training on the open-source RL library Ray RLlib. The training epochs were set to 100; in this case, the trained model could maintain a good balance between the convergence speed and the strategy efficiency.

Evaluation Protocols. We used the image-level area under the ROC curve (AUC) to evaluate the prediction correctness. For the AUC classification metrics, we only evaluated the predictions of images because the output represented the safety perception of the image. To evaluate the visual perceptual pattern, we adopted cosine similarity and the Kullback–Leibler distance (KLD), with smaller KLD values being indicative of better performance.

5.2. Main Results

Comparison Methods. We compared the RL-based method for predicting the perceived safety of an urban area to several methods which performed well in the literature for perceptual safety prediction—namely, (1) the support vector machine (SVM): we predicted the perceived safety of an input image by selecting the same visual features as our method; (2) the multilayer perceptron (MLP): we trained a neural network to directly predict the safety of the input image with their visual features or image matrices, options with higher output scores being selected; and (3) the CNN: we used a 5-layer CNN to classify the input images and obtain the final classification results.

Safety Perception Predictions. Table 3 shows the AUC results of each method on the urban-safety perception prediction. Among them, the RL-based method performs best on the test set. Due to the small size of our crowdsourced dataset, the overfitting of the method using neural networks for direct image prediction is more severe, and the model generalisation is more likely to be affected.

Methods	Input Format	AUC
SVM	Vector	0.617
MLP	Vector	0.611
MLP (Layers $= 5$)	Image Matrix	0.540
CNN (Layers = 5)	Image Matrix	0.550
RL (D3QN)	Vector	0.686

Table 3. Comparison with several methods. AUC evaluates the correctness of the prediction.

The results demonstrate the limitations of supervised methods in the case of a small dataset, as well as the fact that traditional classification methods are inferior to RL-based methods, as RL-based methods can use intrinsic decision patterns and human-like considerations to derive predictions. We applied various DRL algorithms to the training process for our RL-based approach, as shown in Figure 6. Two metrics—namely, cosine similarity and KLD—are used to measure the effectiveness of each method in recovering visual perceptual patterns. The best results in each metric, obtained using D3QN, are listed in Table 4.



The mean rewards in learning process



Table 4. Quantitative results for the decision-making framework in our dataset using different reinforcement learning methods. A higher AUC value or cosine similarity value indicate a better performance. Smaller KLD values indicate better performance.

Methods	AUC	Cosine Similarity	KLD (↓)
PPO	0.619	0.231	0.302
A2C	0.612	0.215	0.295
SAC	0.684	0.369	0.237
D3QN	0.686	0.369	0.244

5.3. Interpretation of Results

Figure 7 shows the reward maps for the four different situations. In Figure 7a, the image is considered safe, with each visual feature receiving different rewards according to the functional feature. The vehicle number receives a reward of 3, whereas the Sky-FoV, Building-FoV, and Greenary-FoV all receive a reward of -1. Furthermore, the number of vehicles in the street-view image is consistent with the experts' perception of a safe image. At the same time, the sky openness, building area, and greenery area are consistent with the experts' judgements on the negative points affecting safety.



Figure 7. Visualising the expert evaluation model based on visual features for four cases. In (**a**) the label of the image is safe, and the prediction given by the model is also safe; (**b**) the label is unsafe, and the prediction is safe; (**c**) the label is safe, and the prediction is unsafe; (**d**) the label is unsafe, and the prediction is also unsafe. "r(Visual Entropy)" refers to the different rewards obtained by the feature visual entropy according to different actions.

This example reflects the experts' intrinsic patterns when judging pictures' safety, proving that the model is effective in imitating the decision making of experts. In Figure 7b, the actual label is unsafe; however, based on our model, the image's visual features satisfy the expert's perception of safety. This proves that there is noise in the labels, making it essential to seek an interpretable human-level prediction method.

6. Conclusions

Dujiangyan was the worst-hit area of the Wenchuan earthquake. Since 2008, we have been conducting a series of follow-up studies on the post-disaster recovery and reconstruction in this city. After the reconstruction, Dujiangyan not only completely retains the style of the old city before the disaster, but also has a large number of new urban landscapes. In addition to the original residents, there is also a large immigrant population. The rapid recovery and promotion of the economy and population prompts us to think about the urban construction and residents' regional identity in disaster-stricken areas. To explore this question, we proposed a RL-based decision framework to predict the perceived safety of urban street-view images, simulate human decision-making patterns, and quantify the impact of visual-semantic features on urban-safety perceptions. We used crowdsourcing and LTR algorithms to obtain the objective safety scores of street-view images. We then obtained visual and functional features using semantic segmentation methods and domain-expert knowledge, respectively, creating an image dataset with features. Subsequently, we modelled the safety evaluation process as an MDP and used RL methods to solve it and obtain a prediction policy. Using this policy, we established a complete and intelligent evaluation model for urban-safety perception and avoided the problem of weak model interpretability caused by supervised learning. We also made better use of expert knowledge and reduced the effect of the high noise rates of crowdsourced data. The experimental results of our crowdsourced dataset showed that our method achieved a satisfactory prediction performance and excellent visual interpretability. It is important to note that, as we apply this framework to other cities, expertise is

needed to fine-tune the reward function to obtain a more accurate and interpretable model. Therefore, for better results, future work should include scaling up the dataset and investigating the reward design by considering other embedding measures with expert knowledge.

Author Contributions: Conceptualisation, Qiushan Li, Yingrui Deng, Yaxuan Wang, and Zhixin Zeng; methodology, Qiushan Li, Yaxuan Wang, Zhixin Zeng, and Yingrui Deng; software, Zhixin Zeng and Yaxuan Wang; validation, Zhixin Zeng and Yaxuan Wang; formal analysis, Qiushan Li and Yingrui Deng; investigation, Qiushan Li, Yingrui Deng, Zhixin Zeng, and Yaxuan Wang; resources, Qiushan Li.; data curation, Qiushan Li, Zhixin Zeng, and Yaxuan Wang; writing—original draft preparation, Zhixin Zeng and Yaxuan Wang; writing—review and editing, Qiushan Li and Yingrui Deng; visualisation, Zhixin Zeng and Yaxuan Wang; supervision, Qiushan Li and Yingrui Deng; tunding acquisition, Qiushan Li, All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by Sichuan University, The Fundamental Research Funds for the Central Universities, grant number 2021SCU12125.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Informed consent was obtained from all subjects involved in the study

Data Availability Statement: Not applicable.

Acknowledgments: The authors would like to thank all reviewers for their helpful and valuable comments and suggestions. Thanks to the provider of street-view data, especially thank the volunteers who participated in the Street View photo safety evaluation and funding from the Sichuan University Foundation is gratefully acknowledged.

Conflicts of Interest: The authors declare no conflict of interest.

References

- 1. Kuo, F.E.; Sullivan, W.C. Environment and Crime in the Inner City: Does Vegetation Reduce Crime? *Environ. Behav.* 2001, *33*, 343–367. https://doi.org/10.1177/0013916501333002.
- Troy, A.; Morgan Grove, J.; O'Neil-Dunne, J. The relationship between tree canopy and crime rates across an urban–rural gradient in the greater Baltimore region. *Landsc. Urban Plan.* 2012, *106*, 262–270. https://doi.org/10.1016/j.landurbplan.2012.03.010.
- Arietta, S.M.; Efros, A.A.; Ramamoorthi, R.; Agrawala, M. City Forensics: Using Visual Elements to Predict Non-Visual City Attributes. *IEEE Trans. Vis. Comput. Graph.* 2014, 20, 2624–2633. https://doi.org/10.1109/TVCG.2014.2346446.
- Zhou, H.; Liu, L.; Lan, M.; Zhu, W.; Song, G.; Jing, F.; Zhong, Y.; Su, Z.; Gu, X. Using Google Street View imagery to capture micro built environment characteristics in drug places, compared with street robbery. *Comput. Environ. Urban Syst.* 2021, 88, 101631. https://doi.org/10.1016/j.compenvurbsys.2021.101631.
- Kelling, G.L.; Coles, C.M. Fixing Broken Windows: Restoring Order and Reducing Crime in Our Communities; A Touchstone Book; Simon & Schuster: New York, NY, USA, 1997; ISBN 978-0-684-83738-3.
- Quercia, D.; O'Hare, N.K.; Cramer, H. Aesthetic capital: What makes london look beautiful, quiet, and happy? In Proceedings of the 17th ACM Conference on Computer Supported Cooperative Work & Social Computing, Portland, OR, USA, 25 January– 1 February 2017; ACM: Baltimore, MD, USA, 2014; pp. 945–955.
- Naik, N.; Philipoom, J.; Raskar, R.; Hidalgo, C. Streetscore Predicting the Perceived Safety of One Million Streetscapes. In Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops, Columbus, OH, USA, 23– 28 June 2014; IEEE: Columbus, OH, USA, 2014; pp. 793–799.
- Cheng, L.; Chu, S.; Zong, W.; Li, S.; Wu, J.; Li, M. Use of Tencent Street View Imagery for Visual Perception of Streets. *ISPRS Int. J. Geo-Inf.* 2017, 6, 265. https://doi.org/10.3390/ijgi6090265.
- 9. Zhang, F.; Zhang, D.; Liu, Y.; Lin, H. Representing place locales using scene elements. *Comput. Environ. Urban Syst.* 2018, 71, 153–164. https://doi.org/10.1016/j.compenvurbsys.2018.05.005.
- 10. Lindal, P.J.; Hartig, T. Architectural variation, building height, and the restorative quality of urban residential streetscapes. *J. Environ. Psychol.* **2013**, *33*, 26–36. https://doi.org/10.1016/j.jenvp.2012.09.003.
- 11. Dubey, A.; Naik, N.; Parikh, D.; Raskar, R.; Hidalgo, C.A. Deep Learning the City : Quantifying Urban Perception At A Global Scale. *arXiv* 2016, arXiv: 160801769.
- Liu, X.; Chen, Q.; Zhu, L.; Xu, Y.; Lin, L. Place-centric Visual Urban Perception with Deep Multi-instance Regression. In Proceedings of the 25th ACM International Conference on Multimedia, Mountain View, CA, USA, 23–27 October 2017; ACM: Mountain View, CA, USA, 2017; pp. 19–27.

- Schölkopf, B.; Platt, J.; Hofmann, T. (Eds.) TrueSkill: A Bayesian Skill Rating System. In Advances in Neural Information Processing Systems 19; The MIT Press: Cambridge, MA, USA, 2007; ISBN 978-0-262-25691-9.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A.A.; Veness, J.; Bellemare, M.G.; Graves, A.; Riedmiller, M.; Fidjeland, A.K.; Ostrovski, G.; et al. Human-level control through deep reinforcement learning. *Nature* 2015, 518, 529–533. https://doi.org/10.1038/nature14236.
- 15. You, C.; Lu, J.; Filev, D.; Tsiotras, P. Advanced planning for autonomous vehicles using reinforcement learning and deep inverse reinforcement learning. *Robot. Auton. Syst.* 2019, 114, 1–18. https://doi.org/10.1016/j.robot.2019.01.003.
- Silver, D.; Huang, A.; Maddison, C.J.; Guez, A.; Sifre, L.; van den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; et al. Mastering the game of Go with deep neural networks and tree search. *Nature* 2016, 529, 484–489. https://doi.org/10.1038/nature16961.
- 17. Bellman, R. A Markovian Decision Process. Indiana Univ. Math. J. 1957, 6, 679–684. https://doi.org/10.1512/iumj.1957.6.56038.
- 18. Lynch, K. *The Image of the City*, 33th ed.; Publication of the Joint Center for Urban Studies; MIT Press: Cambridge, MA, USA, 2008; ISBN 978-0-262-62001-7.
- 19. Jacobs, J. The Death and Life of Great American Cities, Vintage Books ed.; Vintage Books: New York, NY, USA, 1992; ISBN 978-0-679-74195-4.
- 20. Jansson, M.; Fors, H.; Lindgren, T.; Wiström, B. Perceived personal safety in relation to urban woodland vegetation—A review. *Urban For. Urban Green.* **2013**, *12*, 127–133. https://doi.org/10.1016/j.ufug.2013.01.005.
- Li, F. Multilevel modelling of built environment characteristics related to neighbourhood walking activity in older adults. J. Epidemiol. Community Health 2005, 59, 558–564. https://doi.org/10.1136/jech.2004.028399.
- 22. Stafford, M.; Chandola, T.; Marmot, M. Association Between Fear of Crime and Mental Health and Physical Functioning. *Am. J. Public Health* **2007**, *97*, 2076–2081. https://doi.org/10.2105/AJPH.2006.097154.
- 23. Jackson, J.; Stafford, M. Public Health and Fear of Crime: A Prospective Cohort Study. Br. J. Criminol. 2009, 49, 832–847. https://doi.org/10.1093/bjc/azp033.
- 24. Liu, L.; Silva, E.A.; Wu, C.; Wang, H. A machine learning-based method for the large-scale evaluation of the qualities of the urban environment. *Comput. Environ. Urban Syst.* 2017, *65*, 113–125. https://doi.org/10.1016/j.compenvurbsys.2017.06.003.
- 25. He, L.; Páez, A.; Liu, D. Built environment and violent crime: An environmental audit approach using Google Street View. *Comput. Environ. Urban Syst.* 2017, *66*, 83–95. https://doi.org/10.1016/j.compenvurbsys.2017.08.001.
- Nielsen, I.; Smyth, R. Who wants safer cities? Perceptions of public safety and attitudes to migrants among China's urban population. *Int. Rev. Law Econ.* 2008, 28, 46–55. https://doi.org/10.1016/j.irle.2007.12.002.
- Yan, A.F.; Voorhees, C.C.; Clifton, K.; Burnier, C. Do you see what I see?—Correlates of multidimensional measures of neighborhood types and perceived physical activity–related neighborhood barriers and facilitators for urban youth. *Prev. Med.* 2010, 50, S18–S23. https://doi.org/10.1016/j.ypmed.2009.08.015.
- Porzi, L.; Rota Bulò, S.; Lepri, B.; Ricci, E. Predicting and Understanding Urban Perception with Convolutional Neural Networks. In Proceedings of the 23rd ACM International Conference on Multimedia, Lisboa, Portugal, 10–14 October 2022; ACM: Brisbane Australia, 2015; pp. 139–148.
- Ordonez, V.; Berg, T.L. Learning High-Level Judgments of Urban Perception. In *Computer Vision ECCV 2014*; Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T., Eds.; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2014; Volume 8694, pp. 494–510, ISBN 978-3-319-10598-7.
- 30. Acosta, S.; Camargo, J.E. Predicting city safety perception based on visual image content. arXiv 2019, arXiv: 190206871.
- Salesses, P.; Schechtner, K.; Hidalgo, C.A. The Collaborative Image of The City: Mapping the Inequality of Urban Perception. PLoS ONE 2013, 8, e68400. https://doi.org/10.1371/journal.pone.0068400.
- Kaur, T.; Gandhi, T.K. Automated Brain Image Classification Based on VGG-16 and Transfer Learning. In Proceedings of the 2019 International Conference on Information Technology (ICIT), Bhubaneswar, India, 19–21 December 2019; pp. 94–98.
- 33. Williams, R.J. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. In *Reinforcement Learning*; Sutton, R.S., Ed.; Springer: Boston, MA, USA, 1992; pp. 5–32, ISBN 978-1-4613-6608-9.
- Konda, V.; Tsitsiklis, J. Actor-Critic Algorithms. In Advances in Neural Information Processing Systems; Solla, S., Leen, T., Mül-ler, K., Eds.; MIT Press: Cambridge, MA, USA, 1999; Volume 12.
- 35. Sutton, R.S.; McAllester, D.A.; Singh, S.P.; Mansour, Y. Policy Gradient Methods for Reinforcement Learning with Function Approximation 7; AT&T Labs: New York, NY, USA, 1999.
- Sutton, R.S.; McAllester, D.; Singh, S.; Mansour, Y. Policy Gradient Methods for Reinforcement Learning with Function Approximation. In *Advances in Neural Information Processing Systems*; Solla, S., Leen, T., Müller, K., Eds.; MIT Press: Cambridge, MA, USA1999; Volume 12.
- 37. Strasburger, H. Seven Myths on Crowding and Peripheral Vision. *i-Perception* **2020**, *11*, 204166952091305. https://doi.org/10.1177/2041669520913052.
- Schölkopf, B.; Platt, J.; Hofmann, T. (Eds.) Learning to Rank with Nonsmooth Cost Functions. In Advances in Neural Information Processing Systems 19; The MIT Press: Cambridge, MA, USA, 2007; ISBN 978-0-262-25691-9.

- Nasar, J.L. Adult Viewers' Preferences in Residential Scenes: A Study of the Relationship of Environmental Attributes to Preference. *Environ. Behav.* 1983, 15, 589–614. https://doi.org/10.1177/0013916583155003.
- 40. Wohlwill, J.F. Environmental Aesthetics: The Environment as a Source of Affect. In *Human Behavior and Environment;* Altman, I., Wohlwill, J.F., Eds.; Springer US: Boston, MA, USA, 1976; pp. 37–86, ISBN 978-1-4684-2552-9.
- 41. Ewing, R.; Handy, S. Measuring the Unmeasurable: Urban Design Qualities Related to Walkability. J. Urban Des. 2009, 14, 65–84. https://doi.org/10.1080/13574800802451155.
- Zhu, D.; Zhang, F.; Wang, S.; Wang, Y.; Cheng, X.; Huang, Z.; Liu, Y. Understanding Place Characteristics in Geographic Contexts through Graph Convolutional Neural Networks. *Ann. Am. Assoc. Geogr.* 2020, 110, 408–420. https://doi.org/10.1080/24694452.2019.1694403.
- 43. Xie, E.; Wang, W.; Yu, Z.; Anandkumar, A.; Alvarez, J.M.; Luo, P. SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers. *arXiv* 2021, arXiv: 210515203.
- Yuan, Y.; Chen, X.; Wang, J. Object-Contextual Representations for Semantic Segmentation. In *Computer Vision—ECCV 2020*; Vedaldi, A., Bischof, H., Brox, T., Frahm, J.-M., Eds.; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2020; Volume 12351, pp. 173–190, ISBN 978-3-030-58538-9.
- 45. Chen, L.-C.; Zhu, Y.; Papandreou, G.; Schroff, F.; Adam, H. Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation. In *Computer Vision – ECCV 2018*; Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y., Eds.; Lecture Notes in Computer Science; Springer International Publishing: Cham, Switzerland, 2018; Volume 11211, pp. 833–851, ISBN 978-3-030-01233-5.
- 46. Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; Schiele, B. The Cityscapes Dataset for Semantic Urban Scene Understanding. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*; IEEE: Las Vegas, NV, USA, 2016; pp 3213–3223. https://doi.org/10.1109/CVPR.2016.350.
- Ren, Z.; Wang, X.; Zhang, N.; Lv, X.; Li, L.-J. Deep Reinforcement Learning-Based Image Captioning with Embedding Reward. In Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, 21–26 July 2017; pp. 1151–1159.
- Van Hasselt, H.; Guez, A.; Silver, D. Deep Reinforcement Learning with Double Q-Learning. In Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, Phoenix, AR, USA, 12–17 February 2016; Volume 30. https://doi.org/10.1609/aaai.v30i1.10295.
- Wang, Z.; Schaul, T.; Hessel, M.; van Hasselt, H.; Lanctot, M.; de Freitas, N. Dueling Network Architectures for Deep Reinforcement Learning. In Proceedings of the 33rd International Conference on Machine Learning, New York, NY, USA, 19–24 June 2016.