



Explainable, automated urban interventions to improve pedestrian and vehicle safety

C. Bustos^{a,*}, D. Rhoads^{a,*}, A. Solé-Ribalta^a, D. Masip^a, A. Arenas^c, A. Lapedriza^{a,b}, J. Borge-Holthoefer^{a,*}

^a *Internet Interdisciplinary Institute (IN3), Universitat Oberta de Catalunya, Barcelona 08860, Catalonia, Spain*

^b *Media Lab, Massachusetts Institute of Technology, 02139 Cambridge, MA, United States of America*

^c *Departament d'Enginyeria Informàtica i Matemàtiques, Universitat Rovira i Virgili, 43007 Tarragona, Spain*

ARTICLE INFO

Keywords:

Deep learning
Google Street View
Mapillary
Pedestrian
Traffic safety

ABSTRACT

At the moment, urban mobility research and governmental initiatives are mostly focused on motor-related issues, e.g. the problems of congestion and pollution. And yet, we cannot disregard the most vulnerable elements in the urban landscape: pedestrians, exposed to higher risks than other road users. Indeed, safe, accessible, and sustainable transport systems in cities are a core target of the UN's 2030 Agenda. Thus, there is an opportunity to apply advanced computational tools to the problem of traffic safety, in regards especially to pedestrians, who have been often overlooked in the past. This paper combines public data sources, large-scale street imagery and computer vision techniques to approach pedestrian and vehicle safety with an automated, relatively simple, and universally-applicable data-processing scheme. The steps involved in this pipeline include the adaptation and training of a Residual Convolutional Neural Network to determine a hazard index for each given urban scene, as well as an interpretability analysis based on image segmentation and class activation mapping on those same images. Combined, the outcome of this computational approach is a fine-grained map of hazard levels across a city, and an heuristic to identify interventions that might simultaneously improve pedestrian and vehicle safety. The proposed framework should be taken as a complement to the work of urban planners and public authorities.

1. Introduction

In the last century, the accelerated growth of urban areas has given rise to challenges at a variety of levels. Among these, mobility stands out. The ability to efficiently move people and goods is critical to a city's social and economic success (De Domenico et al., 2014; Jiang et al., 2016; Abbar et al., 2018). It is unsurprising, then, the enormous amount of economic and engineering effort that urban planners have devoted to enhance the efficiency of road networks, bus lines, and metro systems (Gakenheimer, 1999). Unlike transportation modes that operate in exclusive spaces, such as metro lines, the uncontrolled rise in urban automotive mobility has gone hand in hand with the degradation of other modes of transportation. Of all these alternative modes, walking has suffered the most, due in large part to the fact that the amount of the streetscape allotted to vehicles invades and interferes with the pedestrian space. Nevertheless, cities exhibit a growing tendency to stop and reverse this process by fostering more active, citizen-friendly transportation modes — foot, bike and personal mobility vehicles, which compete for this public space (Cervero and Duncan, 2003).

* Corresponding authors.

E-mail addresses: mbustosro@uoc.edu (C. Bustos), drhoads@uoc.edu (D. Rhoads), jborgeh@uoc.edu (J. Borge-Holthoefer).

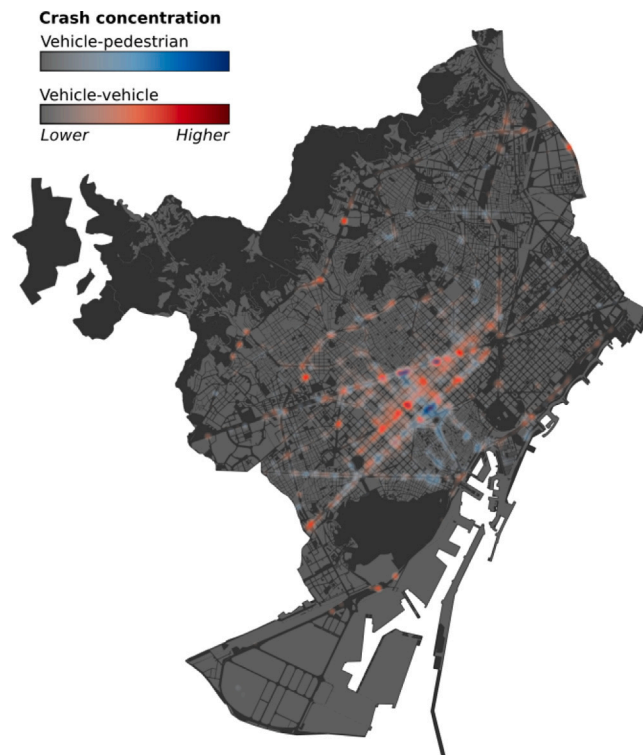


Fig. 1. Accident distribution in Barcelona. Relative concentration of accidents by type (vehicle-to-pedestrian, vehicle-to-vehicle).

One logical consequence of this paradigm shift, is the increased level of interaction between pedestrians and motor vehicles, largely due to the overlapping use of common (or adjacent) spaces such as roads, sidewalks, and zebra-crossings. Such increase gives rise to an important, negative side-effect: a growth in pedestrian injuries and fatalities. Data from the National Highway Traffic Safety Administration (NHTSA) of the United States indicate that the number of pedestrian fatalities per year is rising in the U.S. (National Highway Traffic Safety Administration, 2018). After a steady decline from the mid-1990's to a low in 2009, there has been a clear and consistent reversal until 2017 (the last year of available data), when pedestrian fatalities surpassed a previous 23-year high in 1995.

Traditionally, pedestrian safety research has focused on the impact of structural factors (e.g. road lanes Ukkusuri et al., 2012, traffic network structure Rifaat et al., 2011; Moeinaddini et al., 2014, existence of direct line-of-sight between objects Mecredy et al., 2012; Fu et al., 2019, etc.). In addition, socio-behavioral factors may be concomitant, e.g. the change of individual behavior related to the use of new, distraction-causing technologies (Nasar et al., 2008), inside and outside of vehicles, which is not likely to diminish in the future. Also, demographic variables (socio-economic status, race, gender) may play a role as well (Mukoko and Pulugurtha, 2019). Nonetheless, crashes that involve motor vehicles and pedestrians are understudied, and, at the micro level, much less so outside intersections (Hu et al., 2018).

An enlightening example, built upon real accident data, is shown in Fig. 1. Quite clear even to the naked eye, accidents involving vehicles may happen throughout a city. However, when a distinction is introduced (vehicle-to-vehicle vs. vehicle-to-pedestrian), the spatial patterns where these accidents occur are mostly non-overlapping, suggesting that the configuration of the public space – the scene where the accident happens – matters, see as well Figure S1 in the Supplementary Information (SI). All in all, the strategies for the safe coexistence of pedestrians and vehicles demand a separate and careful examination.

The combination of increasingly available street-level imagery sources and city open data portals, together with advances in the field of computer vision and larger training datasets (Zhou et al., 2014, 2017), has opened up promising new opportunities for facing challenges in urban science. Examples include the quantification of physical change and pattern identification in cities (Naik et al., 2017; Albert et al., 2017; Seiferling et al., 2017), road safety assessment (Song et al., 2018), the prediction of human-perceived features of street scenes (Naik et al., 2014; Liu et al., 2017), the automated estimation of demographic variables across the United States (Gebru et al., 2017) and Great Britain (Suel et al., 2019), or the beautification of urban images through the generation of prototypes (Kauer et al., 2018). Turning to transportation research, however, computer vision has focused mostly on traffic control and surveillance (Fadlullah et al., 2017), and automatic detection and collision prevention (Zhang et al., 2016a,?) for autonomous vehicles. Outside scene analysis, the Deep Learning paradigm has been exploited mostly on motor traffic (Polson and Sokolov, 2017; Wu et al., 2018; Zhang et al., 2018; Wang et al., 2019; Zhang et al., 2019), so far leaving aside its potential to tackle pedestrian safety.

Here, we address the complexities of vehicle-to-pedestrian interaction combining the structural (scene elements) and perceptual (scene composition) aspects of the problem. Overall, the contributions of the present work can be summarized as follows:

1. Creating a dataset of urban street-level images labeled according to accidentality, based on open data municipal accident records.
2. Developing a deep learning architecture, adapted from Deep Residual Networks (ResNet), for hazard index estimation in urban images, that works for both pedestrian and vehicle accidents, and is capable of producing city-wide hazard level landscapes at an unprecedented resolution of one value every 15–20 m.
3. Proposing a set of interpretability analyses to extract human meaning from the outputs of the classification, through customized implementations of Pyramid Scene Parsing networks (PSPNet), Gradient-weighted class activation mapping (GradCam++), radar plots, and a new measure of scene disorder.
4. Designing a greedy heuristic to propose realistic urban interventions, based on scene segmentation, class activation mapping and k-nn algorithm, which constitutes an informed guide for planners to pedestrian safety improvements.

Taken together, these points constitute a novel and comprehensive deep learning pipeline for estimating vehicle and pedestrian hazard in urban scenes, and recommending feasible physical improvements to make those same scenes safer. The building blocks of the pipeline are tailored variants of different state-of-the-art deep learning/machine learning models and techniques (Deep Residual Networks (ResNet), Pyramid Scene Parsing network (PSPNet), Gradient-weighted class activation mapping (GradCam++)).

The remainder of the paper is organized as follows: in Section 2, data (collection, processing techniques and labeling) and methods (pipeline components) are described in detail; then, in Section 3, the results on the hazard index and landscape, its connection to scene composition, and intervention heuristic are presented and discussed. Finally, Section 4 summarizes the work and discusses possible gaps and lines of development.

2. Materials and methods

In this Section we provide the details about the datasets and Deep Learning methods that are used throughout the work. For an introduction to the Deep Learning paradigm, with a focus on transportation systems, we refer to Wang et al. (2019).

2.1. Dataset collection and curation

To feed the proposed framework, we use two types of real urban data: historical accident statistics and street-level urban imagery.

In the case of Madrid and Barcelona, historical accident records for the years 2010–2018 are available from the open data portals of the respective municipal governments (Ayuntamiento de Madrid, 2019; Ajuntament de Barcelona, 2019). For San Francisco, data was available from 2015–2017 and it was filtered from the University of California, Berkeley’s Transport Injury Mapping System (TIMS) of California traffic accidents (Safe Transportation Research and Education Center, 2019). In total, the Barcelona dataset was made up of 86,414 accidents, 10,240 being pedestrian and 76,174 being vehicle accidents. The Madrid dataset had 76,026 accidents (12,533 pedestrian, 63,492 vehicle). In San Francisco, the dataset was made up of 15,492 accidents (3331 pedestrian, 12,161 vehicle). All data points are geolocated with their corresponding GPS coordinates. Besides location, due the detonating causes may be different, we distinguish between accidents where a vehicle and a pedestrian were involved (simply ‘pedestrian’, or P , onwards), from vehicle-to-vehicle accidents (simply ‘vehicle’, or V , onwards). The spatial distribution of empirical accident data for both vehicles and pedestrians can be seen in the SI Figure S1.

Street-level imagery was extracted from two data sources. The Google StreetView (GSV) (Anguelov et al., 2010) API was used for Barcelona and Madrid. In these dataset, images are, on average, 15 meters away from each other. As we wanted to capture the view of the driver, we limited our queries to images facing directly down the direction of traffic of the street. The result of this process was a comprehensive and homogeneous set of images for both cities.

For the city of San Francisco, images were provided by Mapillary (Mapillary contributors, 2019), a crowd-sourced alternative to GSV. With Mapillary, all user-uploaded images are available under the CC-BY-SA license. As images are uploaded by private individuals working with different equipment, different setup, different light conditions, different vehicles, and without central coordination, several distinct challenges were presented by this dataset. Firstly, for each point provided, usually a single image was available. Occasionally, this image did not fit our criteria of facing down the direction of traffic, and had to be discarded. Secondly, data was only available from a smaller part of the city, corresponding to the area covered by the Mapillary contributors. The part of San Francisco available in the dataset, consisting mostly of high-traffic streets, is shown in Figure S2 of the SI.

Combining data from different sources (GSV and Mapillary) allows us to test the robustness of our methods when dealing with similar, but not equally distributed, data. All the collected images, both for GSV and Mapillary, contain GPS locations in their metadata, which allows us to assign each street image a binary accident category (“safe” vs. “dangerous”). We categorize a point as “dangerous” if one or more accidents have occurred with a 50 meter radius of its location. Otherwise, the point is categorized as “safe”. More details on the creation of the image dataset can be found in Section S1 of the SI, along with a more extended discussion of the trade-offs of using a radius to assign accidents to images in Section S4.

The large collection of images tagged according to accident category was divided in 6 different datasets, resulting from the combination of the three targeted cities and two accident types (V and P). The characteristics of each dataset (number of images per dataset and category) are detailed in Table 1.

Notice that the San Francisco datasets are much smaller than Barcelona and Madrid datasets. For the 6 datasets, data was randomly split into train and test sets, containing 90% and 10% of the images respectively.

Table 1

Image dataset properties. Comparing the relative proportion of points with and without accidents across the various cities. In all 3 cities, there is a higher proportion of points with vehicle-to-vehicle accidents than vehicle-to-pedestrian accidents. Relatively less accident points in San Francisco reflects the smaller amount of accident data for that city.

City	Total	Vehicle (V)		Pedestrian (P)	
		Accident	No accident	Accident	No accident
Barcelona	177 645	61.8%	38.2%	48.1%	51.9%
Madrid	704 950	48.3%	51.7%	29.1%	70.9%
San Francisco	162 530	35.7%	64.3%	17.4%	82.6%

2.2. Hazard index estimation with Deep Learning

A variety of Deep Learning architectures have shown to be remarkably effective for many computer vision tasks (LeCun et al., 2015; Schmidhuber, 2015). In this work we use a Residual Neural Network (ResNet) (He et al., 2016a), a particular architecture of Convolutional Neural Network (CNN), to estimate the *hazard index* (H) in new, unseen images. The main characteristic of ResNets is the implementation of “shortcut connections” that skip blocks of convolutional layers, allowing the network to learn residual mappings between layers that mitigate the vanishing gradients problem. For this critical step, all of the elements used were created from scratch – training and test datasets, weight learning stage, etc. – as is detailed in the following.

We define our *hazard index* (H) as the probability that a target image is classified as ‘dangerous’ by the ResNet. For this objective, we train the ResNet to first classify images between the two defined accident categories: ‘dangerous’ and ‘safe’. For each street-level image, the classifier delivers a value H in the range of $[0, 1]$. When $H \approx 1$, the point where the image was taken is considered as dangerous. On the contrary, when $H \approx 0$, the corresponding point is considered as safe. The hazard index is defined as the output of the Softmax activation function (between 0 and 1) of the last layer of the classifier architecture:

$$H = \frac{e^{z_i}}{\sum_{j=0}^K e^{z_j}} \quad (1)$$

where z is the output logits of the last ResNet layer, i is the index of ‘dangerous’ class and K is the number of classes. H can be interpreted as the probability that the point related to a given image is hazardous.

To successfully train our ResNet architecture for the required classification task, we start with a pre-trained network that considers the Imagenet dataset (Krizhevsky et al., 2012), and then, via ‘Transfer learning’ techniques, we fine-tune the network using our data. At this stage, we remove the connections from the last layer of the pre-trained ResNet model, replace it with a new layer with two outputs (categories *dangerous* and *safe*), and randomly initialize the layer’s weights. We re-trained (fine-tuned) this last layer, leaving the rest of the CNN static. To compensate for class imbalance during training stage, class weights were adjusted in the objective cross entropy loss function according to inverse class frequency:

$$w_i = \frac{1}{\ln(c + r_i)} \quad (2)$$

with w_i as the weight assigned to each class, c is a parameter to control the range of the valid values, and r_i is the ratio of the number of samples from each class respect the total of samples, and then

$$Loss = \frac{1}{N} \sum_{i=1}^N w_i \cdot (y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)) \quad (3)$$

where N is the number of samples, and y_i and \hat{y}_i are the true label and the prediction for i class, respectively. In accordance with the defined accident types (V and P), we train our ResNet to estimate two subtypes of hazard index: H_V and H_P , corresponding to the hazard indices for vehicle-to-vehicle and vehicle-to-pedestrian accidents, respectively. Therefore, we end up training 6 models in total, two per city.

2.3. Hazard index interpretability

One of the main shortcomings of Deep Learning techniques is (the lack of) interpretability. Certainly, deep neural networks can provide a high level of discriminative power, but at the cost of introducing many model variables, which eventually hinders the interpretability of their black-box representations (Adadi and Berrada, 2018). This difficulty is especially pertinent in our case: improving pedestrian safety sometimes demands changes in the urban landscape, the question being *which* changes are pertinent. Here, we address this by using two different interpretability techniques. The first, scene disorder, is used to assess image complexity and the second, Class Activation Mapping (CAM), to assess which areas are more informative for the estimation of the hazard index. In particular, CAM methods have been recently shown to be successful for interpretability tasks in several fields (Fukui et al., 2019; Wagner et al., 2019; Desai and Ramaswamy, 2020; Patro et al., 2019), including medicine (Wang and Yang, 2017).

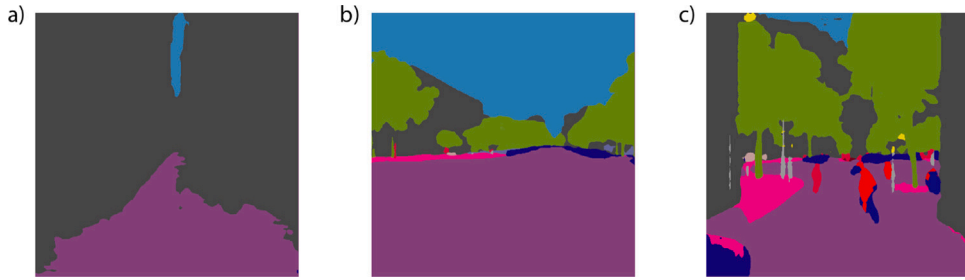


Fig. 2. Illustrating the concept of scene disorder. Segmented images with low $SD = 0.15$ scene disorder (a); mild $SD = 0.39$ scene disorder (b); and high scene disorder $SD = 0.81$ (c).

2.3.1. Urban scene segmentation and scene disorder

First, in order to identify what objects are in the scene, and where they are positioned, we use urban scene segmentation. The goal of the semantic image segmentation task is to assign a category label to each pixel of an image. Segmentation provides a comprehensive breakdown of the physical elements visible in the scene. It predicts the label, location and mask for each object. For this task, we used a high-performance method called Pyramid Scene Parsing Network (PSPNet) (Zhao et al., 2017) architecture, pre-trained with the Cityscapes dataset (Cordts et al., 2016). PSPNet is a state-of-the-art deep learning model that exploits the capability of both global and local context information aggregation through several pyramid pooling layers. It has shown outstanding performance on several semantic segmentation benchmarks. Cityscapes is a real-world, vehicle-ego-centric dataset for semantic urban scene understanding which contains 25K pixel-annotated images taken in different weather conditions. Images in Cityscapes are annotated with 30 urban object categories, but we used a subset of those (19) in our image repository segmentation — those that are common and relevant in driver-perspective scenes (e.g. “car”, “road”, “sidewalk”, “person”, “traffic light”, etc.; see right-most labels in Fig. 4).

On top of the image segmentation outcome, we propose a measure of scene disorder inspired by the gray-tone spatial-dependence matrix (Haralick et al., 1973), also known as Gray-level co-occurrence matrix (GLCM), which captures the amount of transitions between adjacent pixels labeled with different categories. It is known that complex images (related to scene disorder) may cause a division of attention (Moray, 1959; Kahneman, 1973; Alvarez and Cavanagh, 2004; Richards, 2010) and, as a consequence, reduce attention towards objects that are relevant to urban hazard.

Originally, GLCM characterizes the texture of an image by calculating how often pairs of pixels with specific values are adjacent in a specified spatial configuration. In our measure of scene disorder, the frequency of pair of pixels of different values is calculated over the segmented image, where the value of a pixel corresponds to an urban object category, instead of a gray intensity like the usual GLCM. We perform the calculation as follows:

$$SD = \sum_{i=0}^m \sum_{j=0}^n \delta [I(i, j) \neq I(i + \Delta i, j + \Delta j)] \quad (4)$$

where $\delta[x]$ is the Kronecker delta, valued 1 if the condition x is met, and 0 otherwise; and Δi and Δj represent an offset of 1, to compute the amount of pixel value transitions in two directions (right and below). With this definition, the measure SD is incremented by 1 for every pair of neighboring pixels that have differing values. Examples of scene disorder measures can be seen in Fig. 2.

2.3.2. Interpretability through Activation Mapping

Moving on to the second step of our interpretability process, Class Activation Mapping (CAM) (Zhou et al., 2016) and related techniques (e.g. gradient-weighted class activation mapping (GradCAM++) Selvaraju et al., 2017; Chattopadhyay et al., 2018) are used to interpret, visually, the patterns of images that are informative of a specific image category (Ventura et al., 2017; Adadi and Berrada, 2018), meaning, in our case, the regions that have influenced the most about the decision taken by the classifier for a certain class, in our case, classifying an image as ‘dangerous’.

GradCAM++ was used to identify the regions of the image that are dangerous. Given an input image and a our trained CNN model, GradCAM++ generates a localization map by the use of the gradient information of the specific target class ‘dangerous’ to compute the target class weights of each feature map of the last convolutional layer of the CNN before the final classification. The final localization map is synthesized from the aggregated sum of these target class weights. Generating a GradCAM++ map for the ‘dangerous’ class helps to visually identify the specific patterns and objects learned by the CNN in order to differentiate between ‘safe’ and ‘dangerous’ scenes. Since the images have been fully segmented, we can retrieve the objects that overlap with the dangerous regions. Analyzing frequencies, we can recover what object categories are more relevant to determine H_V or H_P . Fig. 4 shows one example per city in the first column and visualizations of the described techniques in the other columns. In particular, second and third column display H_P and H_V , respectively, with the corresponding Class Activation Map. Areas in red color are those that are more relevant to the hazard index, that is, areas that strongly contribute to increase the hazard indexes. Last column shows the automatic segmentation of the images.

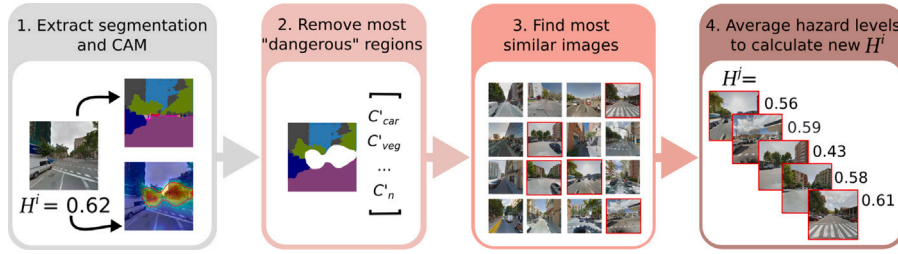


Fig. 3. Image hazard reduction flowchart. Processing pipeline to improve the most hazardous parts of a street-level image i , comparing the new image with similar partner images j , and arriving at a new H_P and H_V for the original image.

2.4. A greedy heuristic to improve H

The combination of the Class Activation Mapping and image segmentation described in the previous section gives us insight into which regions and objects of a scene contribute most to its estimated hazard level. While this information is already relevant, it provides users with no concrete recommendations for structural changes to the scene that might make it safer. Accordingly, as a final step in the pipeline, we propose a strategy to exploit the large pool of images available in order to identify, for each scene, realistic and potentially low-cost physical alterations that would diminish H_P and H_V the most.

To this end, we take advantage of the methodologies developed in the previous steps. On the one hand, the segmentation task allows us to identify which objects among C categories are present in a given scene (and to what extent). On the other, CAM provides information regarding which regions of the scene contribute most to the estimated hazard score. With this information at hand, for every image i we build a vector of characteristics $v_i \in \mathbb{R}^C$, containing information of the relative area of category C in i . For the target scene (the one for which we intend to reduce the hazard levels), we construct an additional surrogate vector of characteristics, \tilde{v}_i , in which we discard those regions that contribute most to H_P , i.e. we only consider regions of i where the class activation is mild-to-low (< 0.7), see first and second blocks in Fig. 3. Next, we deploy an exhaustive search to find the five mirror images j for \tilde{v}_i , with their respective vectors of characteristics v_j , such that their hazard index is lower:

$$\begin{aligned} & \underset{j}{\operatorname{argmin}} \|\tilde{v}_i - v_j\|_2 \\ & H_P^j < H_P^i \\ & H_V^j < H_V^i \end{aligned} \quad (5)$$

In other words, we seek the most similar locations in the city that have smaller H_P and H_V than i , see Fig. 3 for a schematic representation of the process. The search for mirror images is limited to structurally similar scenes (compared to the original one), in order to promote simple and feasible interventions. We emphasize that this strategy is designed to be used in tandem with human users, who will be able to judge which recommendations are realistic. The choice of five images allows for some diversity in the range of interventions recommended.

Finally, we remark that our approach is very similar to the regressive k -nearest neighbor (k -nn) algorithm (Harrington, 2012), as opposed to a more sophisticated, Deep Learning-based mechanism for image “safe-fication” (following the concept of “beautification” in Ref. Kauer et al., 2018). These techniques lie beyond the scope of the present work.

3. Experiments and results

3.1. Hazard index estimation

We begin the results section by assessing how well our trained ResNet performs the required classification task for the six datasets we have defined, considering the cities of Barcelona, Madrid, and San Francisco. Images belonging to the ‘dangerous’ class are defined as positive, while those belonging to the ‘safe’ class are defined as negative. In the training stage, the parameter c of the loss function was experimentally assigned as 1. For our results, we focus on the following measures: recall, precision and accuracy; and the indicators: FP (False positives), TP (True Positives), TN (True Negatives) and FN (False negatives). Recall refers to the fraction of samples detected as dangerous over the total number of dangerous samples in the dataset (TP over TP+FN). Precision is the fraction of the true dangerous points detected, over the number of points detected as dangerous by the ResNet (TP over TP+FP). Accuracy measures how good the system is at detecting dangerous points (TP+TN over all the samples).

As we can see in Table 2, the obtained accuracy is outstanding for all datasets, considering that the CNN training stage relies only on visual information, along with a binary tag indicating the occurrence (or not) of an accident within a 50 m radius (sensitivity with respect to radii is discussed in Section S4.1 and Figure S7 of the SI). As illustrated examples of hazard index estimation, see the scores in the central columns of Fig. 4.

Additionally, we compared the performance of different ResNet and other state-of-the-art architectures against the Barcelona dataset. Metrics like F1-score, area under the Precision and Recall (PR) curve, and the area under the Receiver Operating

Table 2

Results of the deep learning approach for accident prediction, considering a 50 m radius. Rows labeled as *P* and *V* correspond to pedestrian-to-vehicle and vehicle-to-vehicle accident dataset, respectively. Results for other radii can be seen on Table S1 of the SI.

	Recall	Prec.	Acc.	FP	TP	TN	FN
Barcelona <i>P</i>	0.86	0.72	0.75	17.8%	45.4%	29.8%	7%
Barcelona <i>V</i>	0.77	0.84	0.82	7.1%	37.9%	44.1%	10.9%
Madrid <i>P</i>	0.76	0.75	0.75	12.4%	37.5%	38%	12.1%
Madrid <i>V</i>	0.73	0.74	0.75	12%	35.2%	40.1%	12.7%
San Francisco <i>P</i>	0.63	0.81	0.76	6.6%	29%	47.7%	16.7%
San Francisco <i>V</i>	0.61	0.82	0.74	6.3%	30.1%	44.7%	18.9%

Table 3

Results of the deep learning approach for accident prediction, considering different classification architectures.

Model	Acc.	Prec.	Rec.l	F1-score	PR	ROC
VGG16 (Simonyan and Zisserman, 2014)	0.61	0.58	0.96	0.72	0.78	0.59
VGG19 (Simonyan and Zisserman, 2014)	0.68	0.73	0.62	0.67	0.77	0.68
Inception-V3 (Szegedy et al., 2016)	0.70	0.70	0.75	0.72	0.79	0.70
Inception-V4 (Szegedy et al., 2017)	0.57	0.80	0.24	0.37	0.72	0.59
Mobilenet (Howard et al., 2017)	0.62	0.77	0.39	0.52	0.74	0.63
ResNet-v1-50 (He et al., 2016b)	0.61	0.80	0.35	0.49	0.75	0.63
ResNet-v1-101 (He et al., 2016b)	0.59	0.56	0.99	0.71	0.78	0.57
ResNet-v1-152 (He et al., 2016b)	0.67	0.71	0.62	0.66	0.76	0.67
ResNet-v2-50 (He et al., 2016a)	0.75	0.72	0.87	0.78	0.82	0.74
ResNet-v2-101 (He et al., 2016a)	0.72	0.75	0.70	0.72	0.80	0.72
ResNet-v2-152 (He et al., 2016a)	0.72	0.74	0.72	0.73	0.80	0.72

Characteristic (ROC) curve were used for comparison as well. The F1-measure provides a balance between precision and recall in a single score:

$$F1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \quad (6)$$

Whereas the PR curve represents the balance between the measures precision and recall through different thresholds between 0 and 1. The ROC curve plots the false positive rate versus the true positive rate through different thresholds, like the PR curve. The results presented in Table 3 show that the ResNet-v2-50 offers the highest performance for this particular image classification task.

Discerning between safe and dangerous locations in a binary fashion might be limiting in several practical scenarios, such as the prioritization of urban interventions to improve pedestrian safety. To assess to what extent we can produce finer results, we have also implemented the method in Frank and Hall (2001) to learn an ordinal regressor. In this case, the Barcelona pedestrian dataset was divided in four rating classes: *no-danger*, *mild-danger*, *danger* and *high-danger*. Images tagged as ‘no-danger’, correspond those images where no accidents were observed. Images in the class ‘mild-danger’ had one accident nearby, images in class ‘danger’ have between 2 and 5 accidents nearby. Finally, images belonging to class ‘high-danger’ have more than 5 accidents in their vicinity. The dataset proportions were approximately 85k, 34k, 40k and 17k images samples, respectively. The method in Frank and Hall (2001) relies on several binary classifiers. We used our same ResNet architecture for each of those binary classifiers. After training, we obtained a balanced accuracy of 0.47 (with a dummy classifier accuracy of 0.25) which is comparable to the performance reported in Song et al. (2018) for a similar task. That is, the ResNet architecture can also provide competitive results for a finer assessment of pedestrian safety.

3.2. Urban hazard landscape

The first remarkable outcome of the described methodology (in particular, Section 2.2) is a fine-grained map of hazard indices throughout the cities under study. The Deep Learning approach, together with the short distance intervals between consecutive images, allows us to quantify the safety of all city locations at a microscopic level, i.e. every 15 meters approximately (see Figures S3 and S4 in the SI), independently of whether accidents have occurred at a given site or not.

To give a complete picture of hazard for pedestrians and vehicles, and to highlight their differences, Fig. 5 shows the spatial distribution of points that were identified as very hazardous for pedestrians ($H_p \geq 0.66$), but with low-to-moderate hazard for vehicles ($H_V < 0.66$), and vice-versa. As can be seen, in both Madrid and Barcelona, areas of high hazard for pedestrians alone are highly concentrated in the denser, older city centers. High levels of vehicle hazard tend to be distributed around arterial roads, as well as some distinct neighborhoods (e.g. Sant Martí-Poble Nou, middle right corner in Barcelona). San Francisco presents an interesting case in which the two spatial distributions are nearly homogeneous. This can likely be explained by the bias towards residential, medium-density areas in our image coverage for the city (see Materials and Methods for further discussion). Notably, we lacked image coverage in high-density downtown San Francisco, as well as peripheral low-density districts. With the inclusion of such zones, it is possible that clearer spatial patterns would emerge, although they might be distinct from those of denser European cities like Barcelona and Madrid (Louf and Barthelemy, 2014). Nevertheless, it should be noted that competitive levels of precision

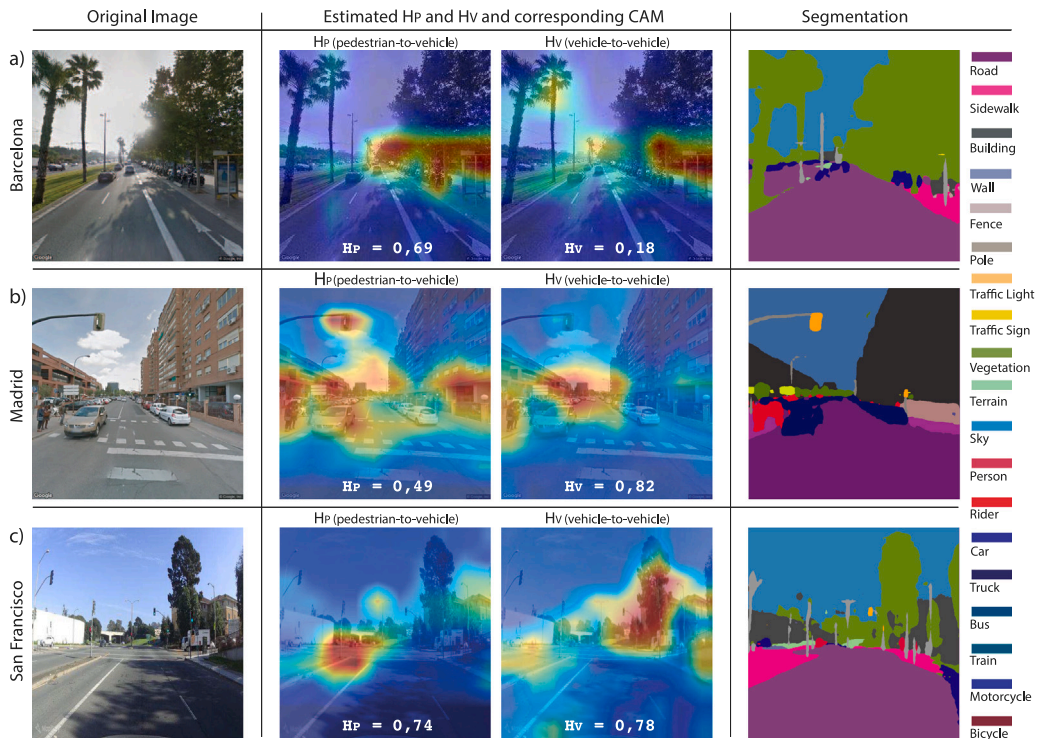


Fig. 4. Deep Learning approach: classification, segmentation and interpretability. The figures display image examples from Barcelona, San Francisco and Madrid, one location per row. First column shows the original street view image. Second and third columns correspond to the obtained CAM for pedestrian and vehicle datasets, respectively. The last column corresponds to the outcome of the segmentation task. The example in Barcelona location (top row) is classified as dangerous for vehicles (note the score in each picture), but safe for pedestrians. Finally, the third example, corresponds to a San Francisco location. Notice that, in this last case, the location is dangerous for both pedestrian and vehicle, but the CAM highlights different regions: areas increasing the hazard for pedestrians may not coincide with those increasing hazard for vehicles. Images courtesy of Google, Inc. and Mapillary. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

and accuracy were still achieved in San Francisco, indicating that our method is robust to relatively homogeneous training data. Furthermore, it shows that the classifier need not only be applied to comprehensive collections of images from an entire city, but can function well on sufficiently rich, spatially homogeneous samples of images. Separate visualizations for pedestrian and vehicle hazards are available in the SI, Figure S3.

Worth highlighting, there has been no previous attempt to associate a given street image with traffic hazard levels — unlike other urban attributes (e.g. beauty [Quercia et al., 2014](#); [Naik et al., 2017](#), or security [Naik et al., 2014](#)). Here, we do so under the assumption that street-level imagery is a good proxy for both the structural and perceptual complexity of the city landscape. Typically, traffic-related risk is either aggregated to the macro-level (neighborhoods, census tracts, even counties) ([Huang et al., 2010](#); [Ukkusuri et al., 2012](#); [Chen and Zhou, 2016](#)), or painstakingly micro-tailored to very specific settings (e.g. considering only zebra-crossings [Olszewski et al., 2016](#)). However, initiatives like Vision Zero, involving governments and organizations worldwide, demand new streams of data and methodologies that help address the street safety challenge at the finest level *and* at scale. This is achieved here combining images and accident data.

3.3. Mapping safety to scene composition

The second (segmentation) and third (Class Activation Mapping, CAM) processing steps complete the data analysis pipeline, linking hazard indices, H_P and H_V , to specific objects found in street-level images. In practice, such link is established combining the information in the central and right columns of [Fig. 4](#). Mapping each pixel label (e.g. “road”, “sidewalk”, etc.) to its corresponding activation level (heatmap in central columns of [Fig. 4](#)) provides a quantification of the contribution of that pixel to the overall hazard score of the image. Thus, at the city level, we can obtain a global perspective of the categories that most contribute to the hazard index.

[Fig. 6](#) (panels a and b) illustrates this for the central area of Barcelona. These radar plots show the level of object fixation of the CAM model for pedestrians (a) and cars (b). In both cases, the blue line represents safe scenes ($H < 0.33$), while dangerous ones ($H > 0.66$) are shown in red. Specifically, we plot the ratio between the amount of CAM fixation on a given category (in safe and dangerous scenes), with respect to the CAM fixation on that category across all the images of the dataset. Thus, values below 1 in

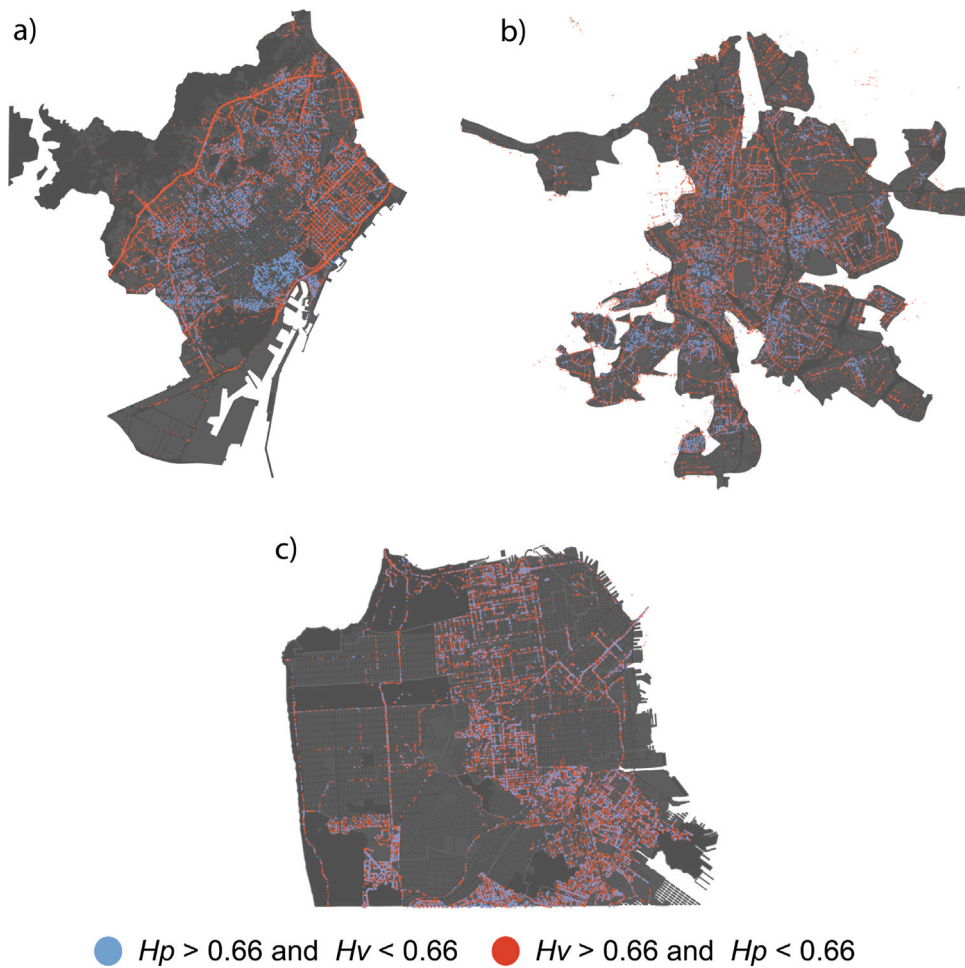


Fig. 5. Spatial distribution of hazard index. Distribution of high-hazard points for pedestrians and vehicles across all three cities of study. Points displayed are those for which hazard is high for pedestrians (vehicles) but not for vehicles (pedestrians).

the radar plots are underrepresented, while those above 1 are overrepresented. We would like to highlight that we have restricted the analysis to the city center, to avoid an exaggeration of the presence of natural elements (vegetation and sky) in low accident risk images. Remarkably, the presence of people in a scene is correlated to a dangerous classification for both vehicle-to-pedestrian and vehicle-to-vehicle predictions. Low buildings and/or wide streets (tantamount to a clear vision of the sky) correlate to safer scenes for pedestrians, whereas the presence of buildings implies a safer environment for vehicles. Also, the absence of vegetation, such as trees, could be contributing to a safe classification for vehicles.

Radar plots for Madrid (see SI, Fig. S5) show high resemblance to the Barcelona ones, while those for San Francisco (Fig. S6) show completely different patterns: for pedestrians, the presence of sidewalks – and not people – is identified as the strongest driver for high H_p . Again, the distinct layouts and walking habits of European and North American cities may be directly related to these emergent patterns.

Moving further, we can relate hazard levels to the scene complexity. While the radar plots show interesting information, they are blind to specific scene compositions in urban scenes, i.e. whether categories appear in a clustered or fragmented way. To grasp this information, we quantify scene disorder (SD) as defined in Eq. (4), see Methods above. Fig. 6c shows an hexbin scatter plot of hazard indices (H_v against H_p), with a color-coded third dimension that corresponds to scene disorder, normalized in the range $[0, 1]$. A first observation is that H_p and H_v are positively correlated. More interestingly, it is clear that more complex scenes (warmer colors) correspond to more dangerous ones. In Figure S5c of the SI, an even clearer trend is shown for Madrid. On the other hand, the level of disorder in San Francisco scenes is high when $H_p \approx H_v \approx 1$, but not clearly related to either H_p or H_v for the rest of values, see Figure S6c. All in all, the connection between image complexity and hazard (especially for vehicles) suggests that more research is needed in this direction. While certain distractions are very explicit (e.g. attending the mobile phone), the perils of scene disorder are subtle and implicit (in the sense that they are not obvious on visual inspection).

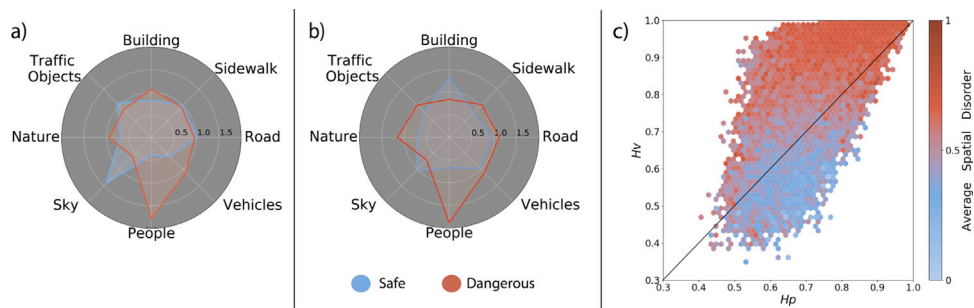


Fig. 6. Hazard level interpretability. Top: Radar plots showing the level of object fixation of the CAM model for pedestrian (a) and cars (b). For both, the blue area corresponds to images classified as safe ($H < 0.33$), while scenes classified as dangerous ($H > 0.66$) are mapped on the plot as red. To build these radars, each individual image is mapped to the radar categories (a relevant subset of those detected by the segmentation task), and the average of such mappings is shown. (c) The plot shows the triple relationship between H_p , H_v and the color-coded level of disorder (adapted from Haralick et al., 1973) –which increases towards warmer colors as the levels of hazard increase. The plot corresponds to Barcelona. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

3.4. An informed guide to pedestrian safety improvements

A precipitate analysis of Fig. 6 may render unfeasible interventions: substitution of built space with larger green areas, building height reduction, or street widening would suffice to improve pedestrian safety, but they do not represent a realistic approach. Instead, we resort on the greedy strategy developed in Section 2.4 to propose interventions conducive to scene alterations that diminish H_p and H_v most.

Fig. 7a shows the results of the application of this optimization to the set of images in Barcelona (Figure S8 in SI for Madrid and San Francisco). In some occasions the hazard index cannot be reduced (points near the (1, 1) coordinate). And yet, many locations present a potential to decrease the hazard levels, even observing, for some scenarios, extreme improvements (points near the (0, 0) coordinate). The gray intensity in Fig. 7a reflects the density of observations in that area. To provide a baseline for comparison, panel b shows alternative results considering a dummy k -nn regressor, that does not take our hazard index into account. Ratios larger than 1 indicate an increase in H_v or H_p , and ratios lower than 1 indicate a decrease. The average in both dimensions is close to zero, evidencing that, with a dummy regressor, we have no guarantee of reducing either pedestrian or vehicle hazard. Fig. 7c shows a selection of two targets and their most similar mirror image, illustrating some common interventions proposed by the heuristic (more examples, for the three cities under study, can be found in Figure S9 of the SI). Visually, all of them seem to point at simplifications of the original image — mostly removing objects on sidewalks.

Finally, Fig. 7d provides a visual overview of the most frequent interventions predicted by our optimization scheme, in the case of Barcelona. The color of the link connecting two categories expresses the source of that link. The most notable changes point – perhaps unsurprisingly – to the need to reconfigure urban scenes towards greener and wider spaces: indeed, both categories ‘road’ and ‘building’ contribute largely to ‘nature’, while the latter does the same towards ‘sky’. Madrid presents an almost identical trend, while San Francisco shows a less clear pattern (although the relevance of ‘nature’ and ‘sky’ is still clear). Both diagrams are available in the SI, Figure S10. Overall, the estimations and insights from the panels in Fig. 7 can provide initial indications to urban planners about achieving potential reductions of a local hazard score, both in terms of which items could be removed or relocated.

4. Discussion

As cities become increasingly populated, the interactions among pedestrians and motorized vehicles become permanent. This translates into a growing number of pedestrian–vehicle accidents. Complementary to the efforts by urban planners, public authorities and sensor technology designers, we present here an automated scheme that exploits a wide range of Computer Vision methods (classification, segmentation and interpretability techniques) to reduce traffic-related fatalities. The proposed processing pipeline, conveniently fed with rich sources of open data, renders an holistic characterization of a city’s hazard landscape, capturing the physical (scene structure) and perceptual (scene complexity) characteristics from a car driver’s point of view. Beyond its informative value, the hazard landscape provides actionable insights to planners.

The main strength of our proposal lies in its simplicity, and its potentially universal applicability out of a comprehensive street image collection and a rich accident dataset. Even crowd-sourced imagery, which is unavoidably diverse and often sparse, provides a solid starting point to quantify safety at a below-segment level. A global, automated, data-driven endeavor towards improving pedestrian safety is not out of reach, considering the advances in cities’ public data portals, and the wide coverage of proprietary services like Google Street View or open initiatives like Mapillary.

Our approach opens a promising line of development. The hazard landscape is defined at an unprecedented, sub-segment resolution level – roughly a hazard score every 15 meters – through an automated and scalable classification process. This is well beyond macroscale approaches (e.g. crash hotspots), and extends the emphasis on intersections (Hu et al., 2018). Such fine-grained map adds a valuable geoinformation layer to those already in use — traffic and pollution levels (Xu et al., 2019), land and

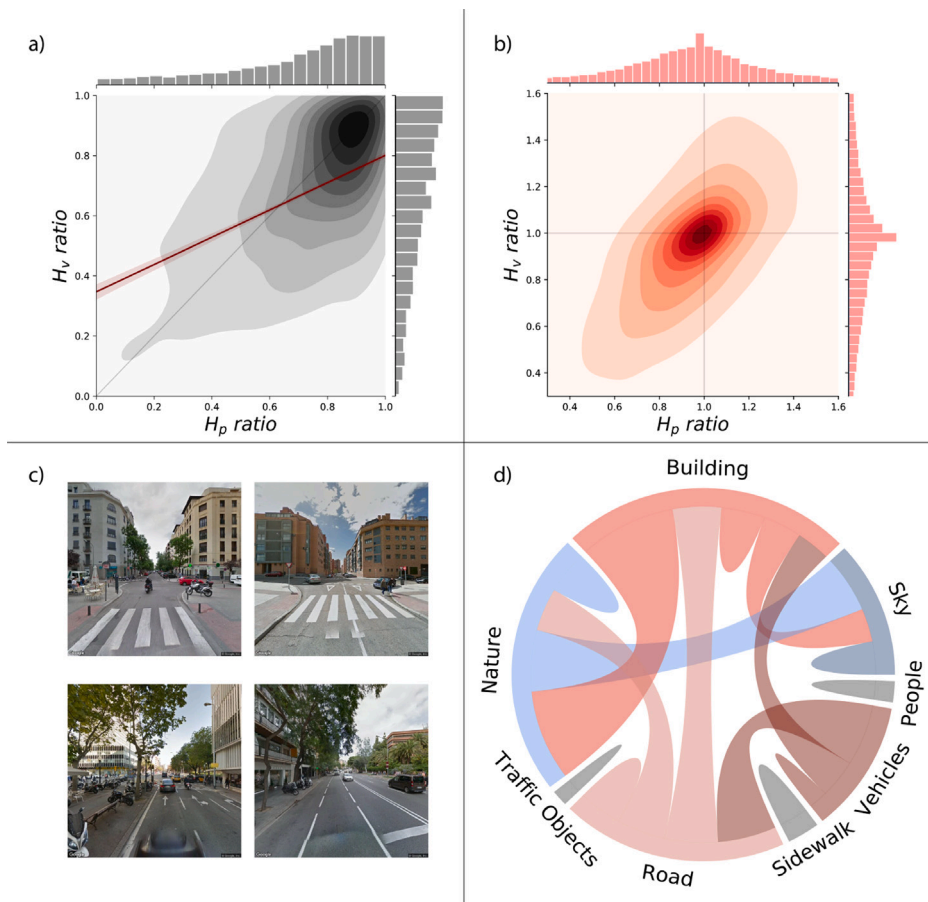


Fig. 7. Hazard reduction: results. (a) Expected improvement for pedestrian and vehicle hazards, with respect to their original values. The horizontal axis corresponds to the ratio between the improved and the original pedestrian hazard index, \tilde{H}_p/H_p ; while the vertical axis represents the equivalent ratio for vehicles, \tilde{H}_v/H_v . Gray intensity represents the density of observations in a given area of the plot. (b) Expected improvement of a dummy k -nn algorithm that only considers similarity between images. This can be regarded as a baseline for results in panel (a) (c) Examples of original and mirror images in Barcelona and Madrid. (d) Chord diagram representing an aggregate overview of proposed interventions in Barcelona. The most notable outcome from the diagram is the propensity to reduce the space allotted to roads and buildings, exchanging it emptier, greener scenes.

underground transportation systems, crime, etc.– enabling better route design: safe paths, along with clean, beautiful, or shortest ones.

Additionally, segmentation and interpretability methods unveil the relationship between potential danger and specific objects in urban scenes. What is more, the disposition of those objects is related to hazard indices, adding a perceptual–attentional link to other possible concomitant variables that affect vehicle and pedestrian safety. Along this line, our work can be used in conjunction with other similar pipelines, such as Song et al. (2018), which automates road safety assessment in terms of infrastructure and estimates road attributes, or may contribute to more focused analysis, relating what a person pays attention to while driving (Palazzi et al., 2018). Additionally, further information such as temporal accident data, or factors known to influence accident rate (e.g. weather, lighting condition, distraction, asphalt conditions, road signaling) could be included by using, for instance, a multi-branch convolutional neural network, to obtain a richer prediction model.

On the other hand, the step from descriptive (hazard landscape) to actionable insights paves the way to automatized, computer-aided prioritization of urban interventions. The proposed heuristic towards safety improvements can serve as a novel tool for planners and policy makers, and might trigger the development of more sophisticated approaches such as the use of Generative Adversarial Networks to produce virtual, plausible alternatives to target scenes (seeking for instance “safe-fication”, instead of “beautification” Kauer et al., 2018). These techniques could be complemented with intervention cost quantification, considering as well cost-safety gain trade-offs.

CRediT authorship contribution statement

C. Bustos: Designed the research, Collected data, Implemented the methods, Performed analyses, Wrote the article, Discussed the results, Reviewed and edited the manuscript. **D. Rhoads:** Designed the research, Collected data, Implemented the methods,

Performed analyses, Wrote the article, Discussed the results, Reviewed and edited the manuscript. **A. Solé-Ribalta**: Conceived the research, Designed the research, Wrote the article, Discussed the results, Reviewed and edited the manuscript. **D. Masip**: Designed the research, Discussed the results, Reviewed and edited the manuscript. **A. Arenas**: Conceived the research, Designed the research, Discussed the results, Reviewed and edited the manuscript. **A. Lapedriza**: Designed the research, Wrote the article, Discussed the results, Reviewed and edited the manuscript. **J. Borge-Holthoefer**: Conceived the research, Designed the research, Wrote the article, Discussed the results, Reviewed and edited the manuscript.

Acknowledgments

All authors acknowledge financial support from the Dirección General de Tráfico (Spain), Project No. SPIP2017-02263, as well as TIN2015-66951-C2-2-R and RTI2018-095232-B-C22 grants from the Spanish Ministry of Science, Innovation and Universities (FEDER funds). CB and DR acknowledge as well the support of a doctoral grant from the Universitat Oberta de Catalunya (UOC). CB, DM and AL acknowledge the NVIDIA Hardware grant program. Street network data copyrighted OpenStreetMap contributors and available from <https://www.openstreetmap.org>.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.trc.2021.103018>.

References

- Abbar, S., Zanoua, T., Borge-Holthoefer, J., 2018. Structural robustness and service reachability in urban settings. *Data Min. Knowl. Discov.* 32, 830–847.
- Adadi, A., Berrada, M., 2018. Peeking inside the black-box: A survey on explainable artificial intelligence (xai). *IEEE Access* 6, 52138–52160.
- Ajuntament de Barcelona, 2019. Open data bcn. <https://opendata-ajuntament.barcelona.cat/en/>. Accessed: 2019-04-20.
- Albert, A., Kaur, J., Gonzalez, M.C., 2017. Using convolutional networks and satellite imagery to identify patterns in urban environments at a large scale. In: Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining, ACM, pp. 1357–1366.
- Alvarez, G.A., Cavanagh, P., 2004. The capacity of visual short-term memory is set both by visual information load and by number of objects. *Psychol. Sci.* 15, 106–111.
- Anguelov, D., Dulong, C., Filip, D., Frueh, C., Lafon, S., Lyon, R., Ogale, A., Vincent, L., Weaver, J., 2010. Google street view: Capturing the world at street level. *Computer* 43, 32–38.
- Ayuntamiento de Madrid, 2019. Portal de datos abiertos del ayuntamiento de madrid. <https://datos.madrid.es/portal/site/egob/>. Accessed: 2019-04-20.
- Cervero, R., Duncan, M., 2003. Walking, bicycling, and urban landscapes: Evidence from the san francisco bay area. *Am J Public Health* 93, 1478–1483, PMID: 12948966.
- Chattopadhyay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N., 2018. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV), IEEE, pp. 839–847.
- Chen, P., Zhou, J., 2016. Effects of the built environment on automobile-involved pedestrian crash frequency and risk. *J. Trans. Health* 3, 448–456.
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B., 2016. The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3213–3223.
- De Domenico, M., Solé-Ribalta, A., Gómez, S., Arenas, A., 2014. Navigability of interconnected networks under random failures. *Proc. Natl. Acad. Sci.* 111, 8351–8356.
- Desai, S., Ramaswamy, H.G., 2020. Ablation-cam: Visual explanations for deep convolutional network via gradient-free localization. In: 2020 IEEE Winter Conference on Applications of Computer Vision (WACV), IEEE, pp. 972–980.
- Fadlullah, Z.M., Tang, F., Mao, B., Kato, N., Akashi, O., Inoue, T., Mizutani, K., 2017. State-of-the-art deep learning: Evolving machine intelligence toward tomorrow's intelligent network traffic control systems. *IEEE Commun. Surv. Tutor.* 19, 2432–2455.
- Frank, E., Hall, M., 2001. A simple approach to ordinal classification. In: European Conference on Machine Learning, Springer, pp. 145–156.
- Fu, T., Hu, W., Miranda-Moreno, L., Saunier, N., 2019. Investigating secondary pedestrian-vehicle interactions at non-signalized intersections using vision-based trajectory data. *Transp. Res. C* 105, 222–240.
- Fukui, H., Hirakawa, T., Yamashita, T., Fujiyoshi, H., 2019. Attention branch network: Learning of attention mechanism for visual explanation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 10705–10714.
- Gakenheimer, R., 1999. Urban mobility in the developing world. *Trans. Res. Part A* 33, 671–689.
- Gebru, T., Krause, J., Wang, Y., Chen, D., Deng, J., Aiden, E.L., Fei-Fei, L., 2017. Using deep learning and google street view to estimate the demographic makeup of neighborhoods across the united states. *Proc. Natl. Acad. Sci.* 114, 13108–13113.
- Haralick, R.M., Shanmugam, K., Dinstein, I.H., 1973. Textural features for image classification. *IEEE Trans. Syst. Man Cybern.* 610–621.
- Harrington, P., 2012. Machine Learning in Action. Manning Publications Co..
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Identity mappings in deep residual networks. In: European conference on computer vision, Springer, pp. 630–645.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778.
- Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H., 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*.
- Hu, Y., Zhang, Y., Shelton, K.S., 2018. Where are the dangerous intersections for pedestrians and cyclists: A colocation-based approach. *Transp. Res. C* 95, 431–441.
- Huang, H., Abdel-Aty, M.A., Darwiche, A.L., 2010. County-level crash risk analysis in florida: Bayesian spatial modeling. *Transp. Res. Rec.* 2148, 27–37.
- Jiang, S., Yang, Y., Gupta, S., Veneziano, D., Athavale, S., González, M.C., 2016. The timegeo modeling framework for urban mobility without travel surveys. *Proc. Natl. Acad. Sci.* 113, E5370–E5378.
- Kahneman, D., 1973. Attention and Effort, Volume 1063. Citeseer.
- Kauer, T., Joglekar, S., Redi, M., Aiello, L.M., Quercia, D., 2018. Mapping and visualizing deep-learning urban beautification. *IEEE Comput. Graph. Appl.* 38, 70–83.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp. 1097–1105.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 521, 436.

- Liu, L., Silva, E.A., Wu, C., Wang, H., 2017. A machine learning-based method for the large-scale evaluation of the qualities of the urban environment. *Comput. Environ. Urban Syst.* 65, 113–125.
- Louf, R., Barthelemy, M., 2014. A typology of street patterns. *J. R. Soc. Interface* 11, 20140924.
- Mapillary contributors, 2019. Mapillary - street-level imagery, powered by collaboration and computer vision. <https://www.mapillary.com/app>.
- Mecredy, G., Janssen, I., Pickett, W., 2012. Neighbourhood street connectivity and injury in youth: a national study of built environments in Canada. *Injury Prev.* 18, 81–87.
- Moainnadin, M., Asadi-Shekari, Z., Shah, M.Z., 2014. The relationship between urban street networks and the number of transport fatalities at the city level. *Saf. Sci.* 62, 114–120.
- Moray, N., 1959. Attention in dichotic listening: Affective cues and the influence of instructions. *Q. J. Exp. Psychol.* 11, 56–60.
- Mukoko, K.K., Pulugurtha, S.S., 2019. Examining the influence of network, land use, and demographic characteristics to estimate the number of bicycle-vehicle crashes on urban roads. *IATSS Res.*
- Naik, N., Kominers, S.D., Raskar, R., Glaeser, E.L., Hidalgo, C.A., 2017. Computer vision uncovers predictors of physical urban change. *Proc. Natl. Acad. Sci.* 114, 7571–7576.
- Naik, N., Philipoom, J., Raskar, R., Hidalgo, C., 2014. Streetscore-predicting the perceived safety of one million streetscapes. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 779–785.
- Nasar, J., Hecht, P., Wener, R., 2008. Mobile telephones, distracted attention, and pedestrian safety. *Accid. Anal. Prev.* 40, 69–75.
- National Highway Traffic Safety Administration, 2018. Fatality analysis reporting system (fars) encyclopedia. <https://www.fars.nhtsa.dot.gov/Main/index.aspx>, Accessed: 2019-06-27.
- Olszewski, P., Buttler, I., Czajewski, W., Dabkowski, P., Kraśkiewicz, C., Szagała, P., Zielińska, A., 2016. Pedestrian safety assessment with video analysis. *Trans. Res. Proc.* 14, 2044–2053.
- Palazzi, A., Abati, D., Solera, F., Cucchiara, R., et al., 2018. Predicting the driver's focus of attention: the dr (eye) ve project. *IEEE Trans. Pattern Anal. Mach. Intell.* 41, 1720–1733.
- Patro, B.N., Lunayach, M., Patel, S., Nambodiri, V.P., 2019. U-cam: Visual explanation using uncertainty based class activation maps. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 7444–7453.
- Polson, N.G., Sokolov, V.O., 2017. Deep learning for short-term traffic flow prediction. *Transp. Res. C* 79, 1–17.
- Quercia, D., Schifanella, R., Aiello, L.M., 2014. The shortest path to happiness: Recommending beautiful, quiet, and happy routes in the city. In: *Proceedings of the 25th ACM conference on Hypertext and social media*, ACM, pp. 116–125.
- Richards, J.E., 2010. The development of attention to simple and complex visual stimuli in infants: Behavioral and psychophysiological measures. *Dev. Rev.* 30, 203–219.
- Rifaat, S.M., Tay, R., De Barros, A., 2011. Effect of street pattern on the severity of crashes involving vulnerable road users. *Accid. Anal. Prev.* 43, 276–283.
- Safe Transportation Research and Education Center, 2019. Transportation Injury Mapping System (Tims). University of California, Berkeley, Accessed: 2019-06-27.
- Schmidhuber, J., 2015. Deep learning in neural networks: An overview. *Neural Netw.* 61, 85–117.
- Seiferling, I., Naik, N., Ratti, C., Proulx, R., 2017. Green streets- quantifying and mapping urban trees with street-level imagery and computer vision. *Landsc. Urban Plan.* 165, 93–101.
- Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D., 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 618–626.
- Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. arXiv preprint [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
- Song, W., Workman, S., Hadzic, A., Zhang, X., Green, E., Chen, M., Souleyrette, R., Jacobs, N., 2018. Farsa: Fully automated roadway safety assessment. In: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, pp. 521–529.
- Suel, E., Polak, J.W., Bennett, J.E., Ezzati, M., 2019. Measuring social, environmental and health inequalities using deep learning and street imagery. *Sci. Rep.* 9, 6229.
- Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A.A., 2017. Inception-v4, inception-resnet and the impact of residual connections on learning. In: *Thirty-first AAAI conference on artificial intelligence*.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z., 2016. Rethinking the inception architecture for computer vision. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826.
- Ukkusuri, S., Miranda-Moreno, L.F., Ramadurai, G., Isa-Tavarez, J., 2012. The role of built environment on pedestrian crash frequency. *Saf. Sci.* 50, 1141–1151.
- Ventura, C., Masip, D., Lapedriza, A., 2017. Interpreting cnn models for apparent personality trait regression. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 55–63.
- Wagner, J., Kohler, J.M., Gindele, T., Hetzel, L., Wiedemer, J.T., Behnke, S., 2019. Interpretable and fine-grained visual explanations for convolutional neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9097–9107.
- Wang, Z., Yang, J., 2017. Diabetic retinopathy detection via deep convolutional networks for discriminative localization and visual explanation. arXiv preprint [arXiv:1703.10757](https://arxiv.org/abs/1703.10757).
- Wang, Y., Zhang, D., Liu, Y., Dai, B., Lee, L.H., 2019. Enhancing transportation systems via deep learning: A survey. *Transp. Res. C* 99, 144–163.
- Wu, Y., Tan, H., Qin, L., Ran, B., Jiang, Z., 2018. A hybrid deep learning based traffic flow prediction method and its understanding. *Transp. Res. C* 90, 166–180.
- Xu, Y., Jiang, S., Li, R., Zhang, J., Zhao, J., Abbar, S., González, M.C., 2019. Unraveling environmental justice in ambient pm2.5 exposure in Beijing: A big data approach. *Comput. Environ. Urban Syst.* 75, 12–21.
- Zhang, S., Benenson, R., Omran, M., Hosang, J., Schiele, B., 2016. How far are we from solving pedestrian detection? In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1259–1267.
- Zhang, Z., He, Q., Gao, J., Ni, M., 2018. A deep learning approach for detecting traffic accidents from social media data. *Transp. Res. C* 86, 580–596.
- Zhang, Z., Li, M., Lin, X., Wang, Y., He, F., 2019. Multistep speed prediction on traffic networks: A deep learning approach considering spatio-temporal dependencies. *Transp. Res. C* 105, 297–322.
- Zhang, L., Lin, L., Liang, X., He, K., 2016a. Is faster r-cnn doing well for pedestrian detection? In: *European Conference on Computer Vision*, Springer, pp. 443–457.
- Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J., 2017. Pyramid scene parsing network. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2881–2890.
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A., 2016. Learning deep features for discriminative localization. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2921–2929.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., Torralba, A., 2017. Places: A 10 million image database for scene recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 40, 1452–1464.
- Zhou, B., Lapedriza, A., Xiao, J., Torralba, A., Oliva, A., 2014. Learning deep features for scene recognition using places database. In: *Advances in neural information processing systems*, pp. 487–495.