

# Explainable, automated urban interventions to improve pedestrian and vehicle safety

C. Bustos<sup>a,\*</sup>, D. Rhoads<sup>a,\*</sup>, A. Solé-Ribalta<sup>a</sup>, D. Masip<sup>a</sup>, A. Arenas<sup>c</sup>, A. Lapedriza<sup>a,b</sup>, J. Borge-Holthoefer<sup>a,\*</sup>

<sup>a</sup>*Internet Interdisciplinary Institute (IN3), Universitat Oberta de Catalunya, Barcelona 08860, Catalonia, Spain*

<sup>b</sup>*Media Lab, Massachusetts Institute of Technology, 02139 Cambridge, MA*

<sup>c</sup>*Departament d'Enginyeria Informàtica i Matemàtiques, Universitat Rovira i Virgili, 43007 Tarragona, Spain*

---

## 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 can not 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 frame-

---

\*Corresponding author

*Email addresses:* [mbustosro@uoc.edu](mailto:mbustosro@uoc.edu) (C. Bustos), [drhoads@uoc.edu](mailto:drhoads@uoc.edu) (D. Rhoads), [jborgeh@uoc.edu](mailto:jborgeh@uoc.edu) (J. Borge-Holthoefer)

work should be taken as a complement to the work of urban planners and public authorities.

*Keywords:* Deep Learning, Google Street View, Mapillary, Pedestrian, Traffic safety

---

## 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 [1–3]. 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 [4]. 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 [5].

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. [6]. 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 [7], traffic network structure [8, 9], existence of direct line-of-sight between objects [10, 11], 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 [12], in-

32 side and outside of vehicles, which is not likely to diminish in the future.  
33 Also, demographic variables (socio-economic status, race, gender) may play  
34 a role as well [13]. Nonetheless, crashes that involve motor vehicles and  
35 pedestrians are understudied, and, at the micro level, much less so outside  
36 intersections [14].

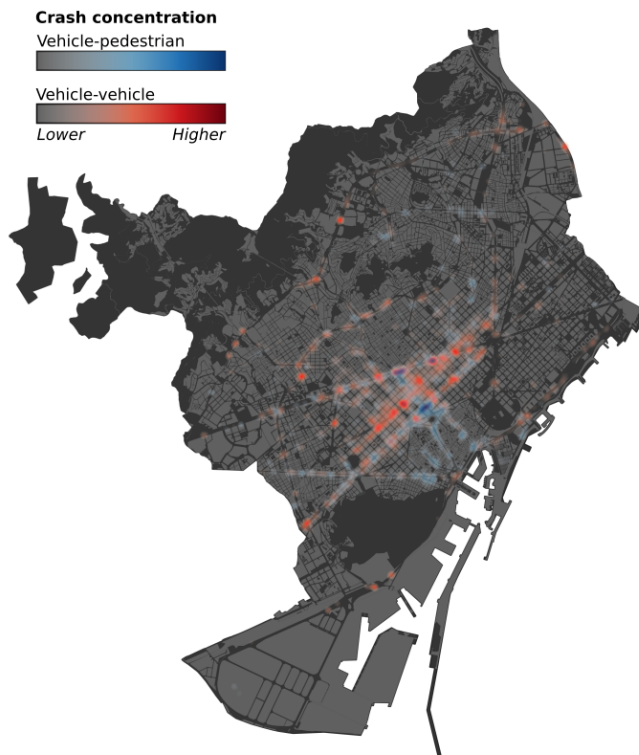


Figure 1: **Accident distribution in Barcelona.** Relative concentration of accidents by type (vehicle-to-pedestrian, vehicle-to-vehicle).

37 An enlightening example, built upon real accident data, is shown in Fig-  
38 ure 1. Quite clear even to the naked eye, accidents involving vehicles may  
39 happen throughout a city. However, when a distinction is introduced (vehicle-  
40 to-vehicle *vs.* vehicle-to-pedestrian), the spatial patterns where these acci-  
41 dents occur are mostly non-overlapping, suggesting that the configuration of  
42 the public space –the scene where the accident happens– matters, see as well  
43 Figure S1 in the Supplementary Information (SI). All in all, the strategies

44 for the safe coexistence of pedestrians and vehicles demand a separate and  
45 careful examination.

46 The combination of increasingly available street-level imagery sources and  
47 city open data portals, together with advances in the field of computer vision  
48 and larger training datasets [15, 16], has opened up promising new oppor-  
49 tunities for facing challenges in urban science. Examples include the quan-  
50 tification of physical change and pattern identification in cities [17–19], road  
51 safety assessment [20], the prediction of human-perceived features of street  
52 scenes [21, 22], the automated estimation of demographic variables across  
53 the United States [23] and Great Britain [24], or the beautification of urban  
54 images through the generation of prototypes [25]. Turning to transportation  
55 research, however, computer vision has focused mostly on traffic control and  
56 surveillance [26], and automatic detection and collision prevention [27, 28] for  
57 autonomous vehicles. Outside scene analysis, the Deep Learning paradigm  
58 has been exploited mostly on motor traffic [29–33], so far leaving aside its  
59 potential to tackle pedestrian safety.

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

- 64 1. Creating a dataset of urban street-level images labelled according to  
65 accidentality, based on open data municipal accident records.
- 66 2. Developing a deep learning architecture, adapted from Deep Residual  
67 Networks (ResNet), for hazard index estimation in urban images, that  
68 works for both pedestrian and vehicle accidents, and is capable of pro-  
69 ducing city-wide hazard level landscapes at an unprecedented resolution  
70 of one value every 15-20 meters.
- 71 3. Proposing a set of interpretability analyses to extract human meaning  
72 from the outputs of the classification, through customized implementa-  
73 tions of Pyramid Scene Parsing networks (PSPNet), Gradient-weighted  
74 class activation mapping (GradCam++), radar plots, and a new mea-  
75 sure of scene disorder.
- 76 4. Designing a greedy heuristic to propose realistic urban interventions,  
77 based on scene segmentation, class activation mapping and k-nn algo-  
78 rithm, which constitutes an informed guide for planners to pedestrian  
79 safety improvements.

80 Taken together, these points constitute a novel and comprehensive deep



81 learning pipeline for estimating vehicle and pedestrian hazard in urban scenes,  
82 and recommending feasible physical improvements to make those same scenes  
83 safer. The building blocks of the pipeline are tailored variants of differ-  
84 ent state-of-the-art deep learning/machine learning models and techniques  
85 (Deep Residual Networks (ResNet), Pyramid Scene Parsing network (PSP-  
86 Net), Gradient-weighted class activation mapping (GradCam++)).

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

## 93 2. Materials and Methods

94 In this Section we provide the details about the datasets and Deep Learn-  
95 ing methods that are used throughout the work. For an introduction to the  
96 Deep Learning paradigm, with a focus on transportation systems, we refer  
97 to Wang *et. al.* [32].

### 98 2.1. Dataset collection and curation

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

101 In the case of Madrid and Barcelona, historical accident records for the  
102 years 2010-2018 are available from the open data portals of the respective  
103 municipal governments [34, 35]. For San Francisco, data was available from  
104 2015-2017 and it was filtered from the University of California, Berkeley’s  
105 Transport Injury Mapping System (TIMS) of California traffic accidents [36].  
106 In total, the Barcelona dataset was made up of 86,414 accidents, 10,240 be-  
107 ing pedestrian and 76,174 being vehicle accidents. The Madrid dataset had  
108 76,026 accidents (12,533 pedestrian, 63,492 vehicle). In San Francisco, the  
109 dataset was made up of 15,492 accidents (3331 pedestrian, 12,161 vehicle).  
110 All data points are geolocated with their corresponding GPS coordinates.  
111 Besides location, due the detonating causes may be different, we distinguish  
112 between accidents where a vehicle and a pedestrian were involved (simply  
113 ‘pedestrian’, or  $P$ , onwards), from vehicle-to-vehicle accidents (simply ‘vehi-  
114 cle’, or  $V$ , onwards). The spatial distribution of empirical accident data for  
115 both vehicles and pedestrians can be seen in the SI Figure S1.

116 Street-level imagery was extracted from two data sources. The Google  
117 StreetView (GSV) [37] API was used for Barcelona and Madrid. In these  
118 dataset, images are, on average, 15 meters away from each other. As we  
119 wanted to capture the view of the driver, we limited our queries to images  
120 facing directly down the direction of traffic of the street. The result of this  
121 process was a comprehensive and homogeneous set of images for both cities.

122 For the city of San Francisco, images were provided by Mapillary [38], a  
123 crowd-sourced alternative to GSV. With Mapillary, all user-uploaded images  
124 are available under the CC-BY-SA license. As images are uploaded by private  
125 individuals working with different equipment, different setup, different  
126 light conditions, different vehicles, and without central coordination, several  
127 distinct challenges were presented by this dataset. Firstly, for each point  
128 provided, usually a single image was available. Occasionally, this image did  
129 not fit our criteria of facing down the direction of traffic, and had to be dis-  
130 carded. Secondly, data was only available from a smaller part of the city,  
131 corresponding to the area covered by the Mapillary contributors. The part  
132 of San Francisco available in the dataset, consisting mostly of high-traffic  
133 streets, is shown in Figure S2 of the SI.

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

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

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

## 151 2.2. Hazard index estimation with Deep Learning

City	Total	Vehicle ( $V$ )		Pedestrian ( $P$ )	
		Accident	No accident	Accident	No accident
Barcelona	177645	61.8%	38.2%	48.1%	51.9%
Madrid	704950	48.3%	51.7%	29.1%	70.9%
San Francisco	162530	35.7%	64.3%	17.4%	82.6%

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.

152 A variety of Deep Learning architectures have shown to be remarkably  
153 effective for many computer vision tasks [39, 40]. In this work we use a  
154 Residual Neural Network (ResNet) [41], a particular architecture of Convo-  
155 lutional Neural Network (CNN), to estimate the *hazard index* ( $H$ ) in new,  
156 unseen images. The main characteristic of ResNets is the implementation of  
157 “shortcut connections” that skip blocks of convolutional layers, allowing the  
158 network to learn residual mappings between layers that mitigate the vanish-  
159 ing gradients problem. For this critical step, all of the elements used were  
160 created from scratch – training and test datasets, weight learning stage, etc.  
161 – as is detailed in the following.

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

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

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

174 To successfully train our ResNet architecture for the required classifica-

175 tion task, we start with a pre-trained network that considers the Imagenet  
 176 dataset [42], and then, via 'Transfer learning' techniques, we fine-tune the  
 177 network using our data. At this stage, we remove the connections from  
 178 the last layer of the pre-trained ResNet model, replace it with a new layer  
 179 with two outputs (categories *dangerous* and *safe*), and randomly initialize  
 180 the layer's weights. We re-trained (fine-tuned) this last layer, leaving the  
 181 rest of the CNN static. To compensate for class imbalance during training  
 182 stage, class weights were adjusted in the objective cross entropy loss function  
 183 according to inverse class frequency:

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

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

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

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

### 193 2.3. Hazard index interpretability

194 One of the main shortcomings of Deep Learning techniques is (the lack  
 195 of) interpretability. Certainly, deep neural networks can provide a high level  
 196 of discriminative power, but at the cost of introducing many model vari-  
 197 ables, which eventually hinders the interpretability of their black-box repre-  
 198 sentations [43]. This difficulty is especially pertinent in our case: improving  
 199 pedestrian safety sometimes demands changes in the urban landscape, the  
 200 question being *which* changes are pertinent. Here, we address this by using  
 201 two different interpretability techniques. The first, scene disorder, is used to  
 202 assess image complexity and the second, Class Activation Mapping (CAM),  
 203 to assess which areas are more informative for the estimation of the hazard

204 index. In particular, CAM methods have been recently shown to be suc-  
205 cessful for interpretability tasks in several fields [44–47], including medicine  
206 [48].

### 207 *2.3.1. Urban scene segmentation and scene disorder*

208 First, in order to identify what objects are in the scene, and where they  
209 are positioned, we use urban scene segmentation. The goal of the semantic  
210 image segmentation task is to assign a category label to each pixel of an  
211 image. Segmentation provides a comprehensive breakdown of the physical  
212 elements visible in the scene. It predicts the label, location and mask for  
213 each object. For this task, we used a high-performance method called Pyra-  
214 mid Scene Parsing Network (PSPNet) [49] architecture, pre-trained with the  
215 Cityscapes dataset [50]. PSPNet is a state-of-the-art deep learning model  
216 that exploits the capability of both global and local context information ag-  
217 gregation through several pyramid pooling layers. It has shown outstanding  
218 performance on several semantic segmentation benchmarks. Cityscapes is a  
219 real-world, vehicle-egocentric dataset for semantic urban scene understanding  
220 which contains 25K pixel-annotated images taken in different weather condi-  
221 tions. Images in Cityscapes are annotated with 30 urban object categories,  
222 but we used a subset of those (19) in our image repository segmentation –  
223 those that are common and relevant in driver-perspective scenes (e.g. “car”,  
224 “road”, “sidewalk”, “person”, “traffic light”, etc.; see right-most labels in  
225 Figure 4).

226 On top of the image segmentation outcome, we propose a measure of scene  
227 disorder inspired by the gray-tone spatial-dependence matrix [51], also known  
228 as Gray-level co-occurrence matrix (GLCM), which captures the amount of  
229 transitions between adjacent pixels labelled with different categories. It is  
230 known that complex images (related to scene disorder) may cause a division  
231 of attention [52–55] and, as a consequence, reduce attention towards objects  
232 that are relevant to urban hazard.

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

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

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

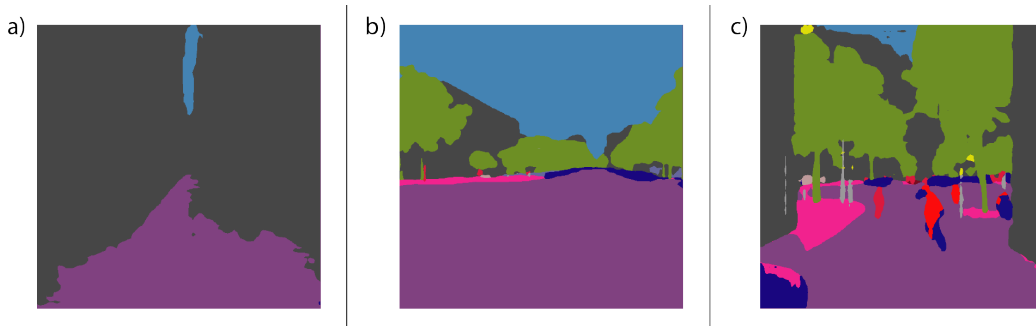


Figure 2: **Illustrating the concept of scene disorder.** Segmented images with low  $SD = 0.15$  scene disorder (a); mild  $SD = 0.39$  scene disorder (b); and high scene disorder  $SD = 0.81$  (c).

### 245 2.3.2. Interpretability through Activation Mapping

246 Moving on to the second step of our interpretability process, Class Acti-  
 247 vation Mapping (CAM) [56] and related techniques (e.g. gradient-weighted  
 248 class activation mapping (GradCAM++) [57, 58]) are used to interpret, visu-  
 249 ally, the patterns of images that are informative of a specific image category  
 250 [43, 59], meaning, in our case, the regions that have influenced the most  
 251 about the decision taken by the classifier for a certain class, in our case,  
 252 classifying an image as 'dangerous'.

253 GradCAM++ was used to identify the regions of the image that are dan-  
 254 gerous. Given an input image and a our trained CNN model, GradCAM++  
 255 generates a localization map by the use of the gradient information of the  
 256 specific target class 'dangerous' to compute the target class weights of each

257 feature map of the last convolutional layer of the CNN before the final clas-  
 258 sification. The final localization map is synthesized from the aggregated  
 259 sum of these target class weights. Generating a GradCAM++ map for the  
 260 'dangerous' class helps to visually identify the specific patterns and objects  
 261 learned by the CNN in order to differentiate between 'safe' and 'dangerous'  
 262 scenes. Since the images have been fully segmented, we can retrieve the objects  
 263 that overlap with the dangerous regions. Analyzing frequencies, we can  
 264 recover what object categories are more relevant to determine  $H_V$  or  $H_P$ .  
 265 Figure 4 shows one example per city in the first column and visualizations  
 266 of the described techniques in the other columns. In particular, second and  
 267 third column display  $H_P$  and  $H_V$ , respectively, with the corresponding Class  
 268 Activation Map. Areas in red color are those that are more relevant to the  
 269 hazard index, that is, areas that strongly contribute to increase the hazard  
 270 indexes. Last column shows the automatic segmentation of the images.

271 *2.4. A greedy heuristic to improve  $H$*

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

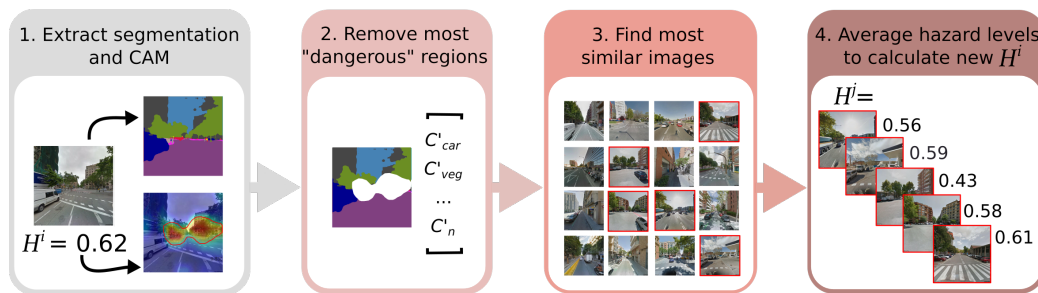


Figure 3: **Image hazard reduction flowchart.** Processing pipeline to improve the most hazardous parts of a street-level image  $i$ , comparing the new image with similar partner images  $j$ , and arriving at a new  $H_P$  and  $H_V$  for the original image.

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

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

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

303 Finally, we remark that our approach is very similar to the regressive  $k$ -  
 304 nearest neighbor ( $k$ -nn) algorithm [60], as opposed to a more sophisticated,  
 305 Deep Learning-based mechanism for image “safe-fication” (following the con-  
 306 cept of “beautification” in ref. [25]). These techniques lie beyond the scope  
 307 of the present work.

### 308 3. Experiments and Results

#### 309 3.1. Hazard index Estimation

310 We begin the results section by assessing how well our trained ResNet  
 311 performs the required classification task for the six datasets we have defined,



312 considering the cities of Barcelona, Madrid, and San Francisco. Images be-  
 313 longing to the ‘dangerous’ class are defined as positive, while those belonging  
 314 to the ‘safe’ class are defined as negative. In the training stage, the param-  
 315 eter  $c$  of the loss function was experimentally assigned as 1. For our results,  
 316 we focus on the following measures: recall, precision and accuracy; and the  
 317 indicators: FP (False positives), TP (True Positives), TN (True Negatives)  
 318 and FN (False negatives). Recall refers to the fraction of samples detected  
 319 as dangerous over the total number of dangerous samples in the dataset  
 320 (TP over TP+FN). Precision is the fraction of the true dangerous points  
 321 detected, over the number of points detected as dangerous by the ResNet  
 322 (TP over TP+FP). Accuracy measures how good the system is at detecting  
 323 dangerous points (TP+TN over all the samples).

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

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

Table 2: Results of the Deep Learning approach for accident prediction, considering a 50 meters radius. Rows labelled 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.

330 Additionally, we compared the performance of different ResNet and other  
 331 state-of-the-art architectures against the Barcelona dataset. Metrics like F1-  
 332 score, area under the Precision and Recall (PR) curve, and the area under the  
 333 Receiver Operating Characteristic (ROC) curve were used for comparison as  
 334 well. The F1-measure provides a balance between precision and recall in a

335 single score:

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

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

342 Discerning between safe and dangerous locations in a binary fashion might  
343 be limiting in several practical scenarios, such as the prioritization of urban  
344 interventions to improve pedestrian safety. To assess to what extent we  
345 can produce finer results, we have also implemented the method in [61] to  
346 learn an ordinal regressor. In this case, the Barcelona pedestrian dataset  
347 was divided in four rating classes: *no-danger*, *mild-danger*, *danger* and *high-*  
348 *danger*. Images tagged as ‘no-danger’, correspond those images where no  
349 accidents were observed. Images in the class ‘mild-danger’ had one accident  
350 nearby, images in class ‘danger’ have between 2 and 5 accidents nearby.  
351 Finally, images belonging to class ‘high-danger’ have more than 5 accidents in  
352 their vicinity. The dataset proportions were approximately 85k, 34k, 40k and  
353 17k images samples, respectively. The method in [61] relies on several binary  
354 classifiers. We used our same ResNet architecture for each of those binary  
355 classifiers. After training, we obtained a balanced accuracy of 0.47 (with a the  
356 dummy classifier accuracy of 0.25) which is comparable to the performance  
357 reported in [20] for a similar task. That is, the ResNet architecture can also  
358 provide competitive results for a finer assessment of pedestrian safety.

### 359 3.2. *Urban hazard landscape*

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

367 To give a complete picture of hazard for pedestrians and vehicles, and to  
368 highlight their differences, Figure 5 shows the spatial distribution of points  
369 that were identified as very hazardous for pedestrians ( $H_P \geq 0.66$ ), but with

Model	Acc.	Prec.	Rec.1	F1-Score	PR	ROC
VGG16 [62]	0.61	0.58	0.96	0.72	0.78	0.59
VGG19 [62]	0.68	0.73	0.62	0.67	0.77	0.68
Inception-V3 [63]	0.70	0.70	0.75	0.72	0.79	0.70
Inception-V4 [64]	0.57	0.80	0.24	0.37	0.72	0.59
Mobilenet [65]	0.62	0.77	0.39	0.52	0.74	0.63
ResNet-v1-50 [66]	0.61	0.80	0.35	0.49	0.75	0.63
ResNet-v1-101 [66]	0.59	0.56	0.99	0.71	0.78	0.57
ResNet-v1-152 [66]	0.67	0.71	0.62	0.66	0.76	0.67
ResNet-v2-50 [41]	<b>0.75</b>	0.72	0.87	<b>0.78</b>	<b>0.82</b>	<b>0.74</b>
ResNet-v2-101 [41]	0.72	0.75	0.70	0.72	0.80	0.72
ResNet-v2-152 [41]	0.72	0.74	0.72	0.73	0.80	0.72

Table 3: Results of the Deep Learning approach for accident prediction, considering different classification architectures.

370 low-to-moderate hazard for vehicles ( $H_V < 0.66$ ), and vice-versa. As can be  
371 seen, in both Madrid and Barcelona, areas of high hazard for pedestrians  
372 alone are highly concentrated in the denser, older city centers. High levels of  
373 vehicle hazard tend to be distributed around arterial roads, as well as some  
374 distinct neighborhoods (e.g. Sant Martí-Poble Nou, middle right corner in  
375 Barcelona). San Francisco presents an interesting case in which the two  
376 spatial distributions are nearly homogeneous. This can likely be explained  
377 by the bias towards residential, medium-density areas in our image coverage  
378 for the city (see Materials and Methods for further discussion). Notably,  
379 we lacked image coverage in high-density downtown San Francisco, as well  
380 as peripheral low-density districts. With the inclusion of such zones, it is  
381 possible that clearer spatial patterns would emerge, although they might be  
382 distinct from those of denser European cities like Barcelona and Madrid [67].  
383 Nevertheless, it should be noted that competitive levels of precision and  
384 accuracy were still achieved in San Francisco, indicating that our method  
385 is robust to relatively homogeneous training data. Furthermore, it shows  
386 that the classifier need not only be applied to comprehensive collections of  
387 images from an entire city, but can function well on sufficiently rich, spatially  
388 homogeneous samples of images. Separate visualizations for pedestrian and  
389 vehicle hazards are available in the SI, Figure S3.

390 Worth highlighting, there has been no previous attempt to associate a  
391 given street image with traffic hazard levels –unlike other urban attributes

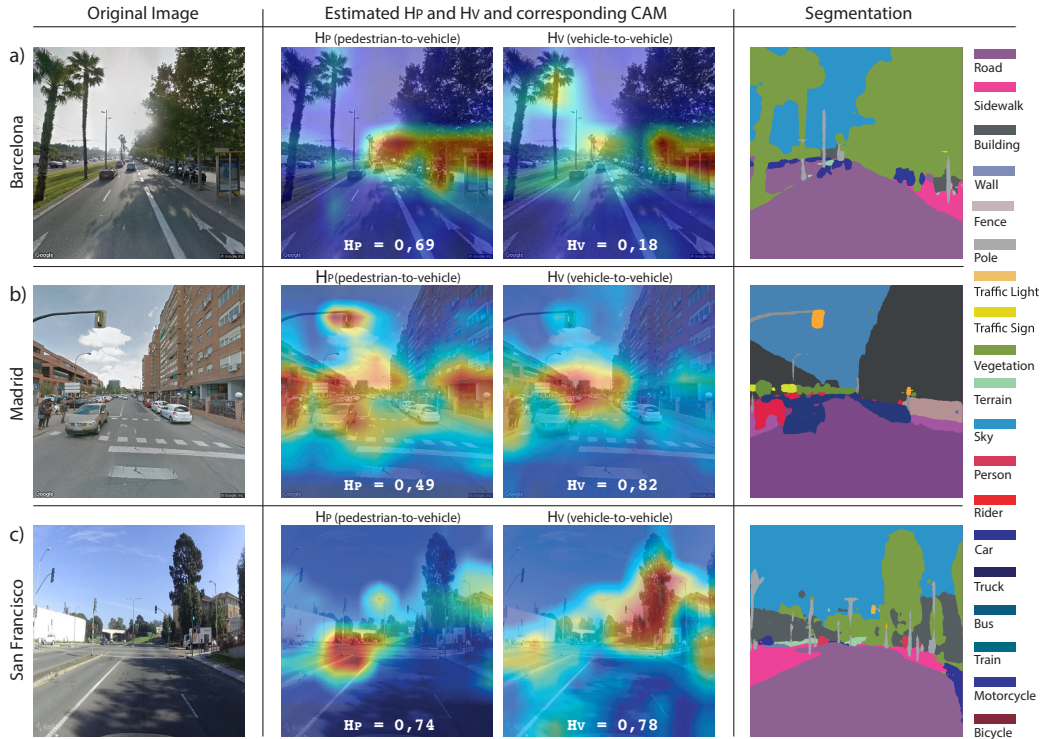


Figure 4: **Deep Learning approach: classification, segmentation and interpretability.** The figures display image examples from Barcelona, San Francisco and Madrid, one location per row. First column shows the original street view image. Second and third columns correspond to the obtained CAM for pedestrian 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 pedestrians (note the score in each picture), but safe for vehicles. The second example, corresponding to a Madrid location, is classified as dangerous for vehicles, 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.

392 (e.g. beauty [17, 68], or security [21]). Here, we do so under the assumption  
 393 that street-level imagery is a good proxy for both the structural and  
 394 perceptual complexity of the city landscape. Typically, traffic-related risk  
 395 is either aggregated to the macro-level (neighborhoods, census tracts, even  
 396 counties)[7, 69, 70], or painstakingly micro-tailored to very specific settings

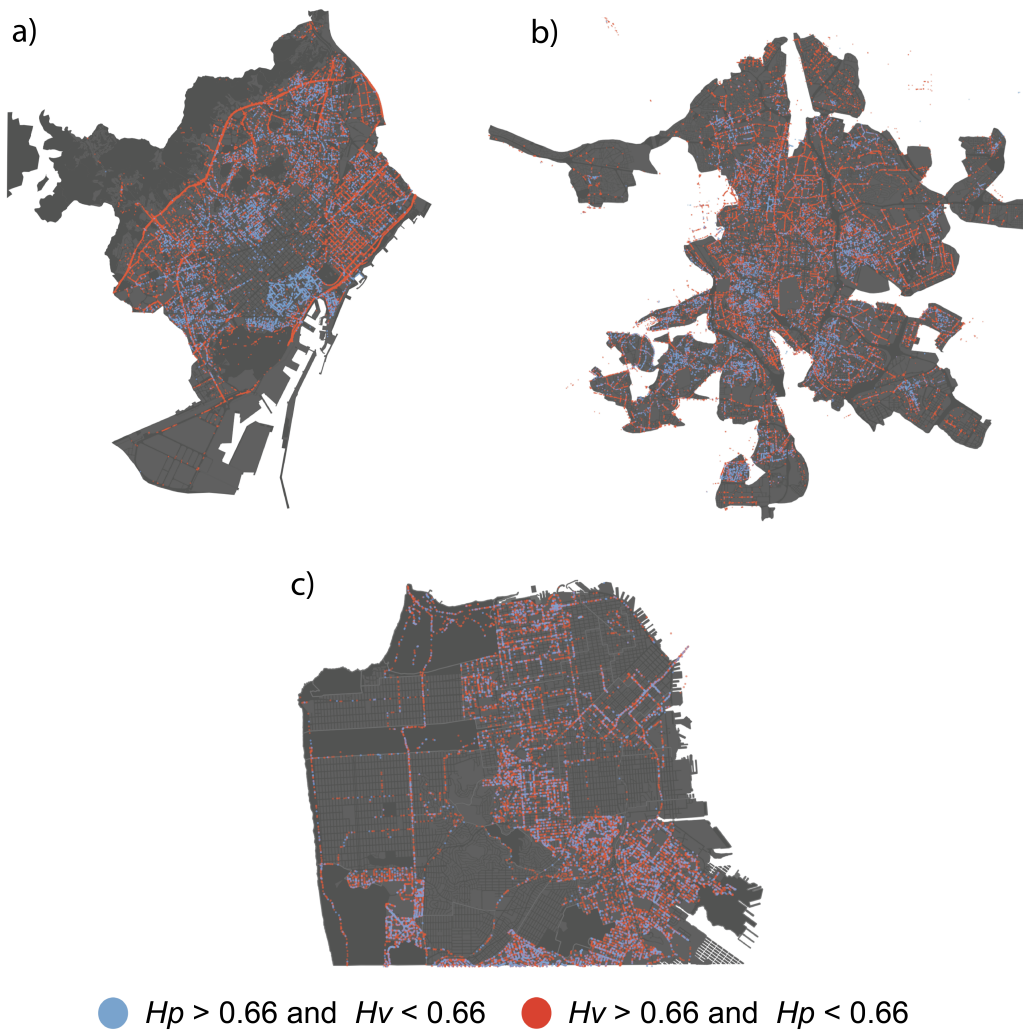


Figure 5: **Spatial distribution of hazard index.** Distribution of high-hazard points for pedestrians and vehicles across all three cities of study. Points displayed are those for which hazard is high for pedestrians (vehicles) but not for vehicles (pedestrians).

397 (e.g. considering only zebra-crossings [71]). However, initiatives like Vi-  
398 sion Zero, involving governments and organizations worldwide, demand new  
399 streams of data and methodologies that help address the street safety chal-  
400 lenge at the finest level *and* at scale. This is achieved here combining images  
401 and accident data.

### 402 3.3. Mapping safety to scene composition

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

412 Figure 6 (panels a and b) illustrates this for the central area of Barcelona.  
413 These radar plots show the level of object fixation of the CAM model for  
414 pedestrians (a) and cars (b). In both cases, the blue line represents safe  
415 scenes ( $H < 0.33$ ), while dangerous ones ( $H > 0.66$ ) are shown in red.  
416 Specifically, we plot the ratio between the amount of CAM fixation on a  
417 given category (in safe and dangerous scenes), with respect to the CAM fix-  
418 ation on that category across all the images of the dataset. Thus, values  
419 below 1 in the radar plots are underrepresented, while those above 1 are  
420 overrepresented. We would like to highlight that we have restricted the anal-  
421 ysis to the city center, to avoid an exaggeration of the presence of natural  
422 elements (vegetation and sky) in low accident risk images. Remarkably, the  
423 presence of people in a scene is correlated to a dangerous classification for  
424 both vehicle-to-pedestrian and vehicle-to-vehicle predictions. Low buildings  
425 and/or wide streets (tantamount to a clear vision of the sky) correlate to  
426 safer scenes for pedestrians, whereas the presence of buildings implies a safer  
427 environment for vehicles. Also, the absence of vegetation, such as trees, could  
428 be contributing to a safe classification for vehicles.

429 Radar plots for Madrid (see SI, Fig. S5) show high resemblance to the  
430 Barcelona ones, while those for San Francisco (Fig. S6) show completely  
431 different patterns: for pedestrians, the presence of sidewalks –and not people–  
432 is identified as the strongest driver for high  $H_P$ . Again, the distinct layouts

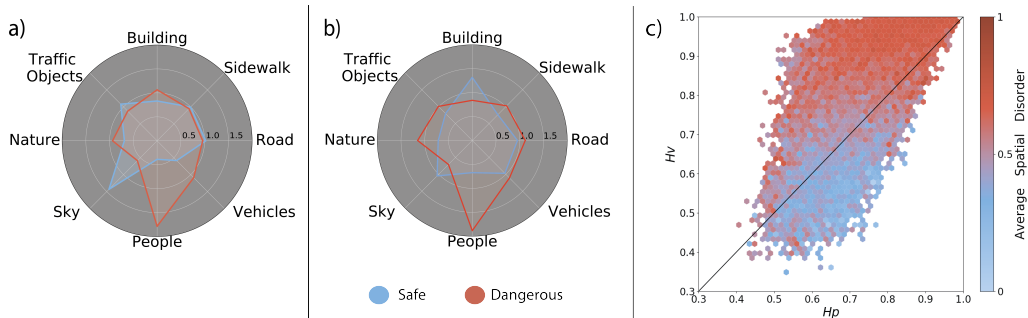


Figure 6: **Hazard level interpretability.** **Top:** Radar plots showing the level of object fixation of the CAM model for pedestrian (a) and cars (b). For both, the blue area corresponds to images classified as safe ( $H < 0.33$ ), while scenes classified as dangerous ( $H > 0.66$ ) are mapped on the 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 [51]) –which increases towards warmer colors as the levels of hazard increase. The plot corresponds to Barcelona.

433 and walking habits of European and North American cities may be directly  
 434 related to these emergent patterns.

435 Moving further, we can relate hazard levels to the scene complexity. While  
 436 the radar plots show interesting information, they are blind to specific scene  
 437 compositions in urban scenes, i.e. whether categories appear in a clustered  
 438 or fragmented way. To grasp this information, we quantify scene disorder  
 439 ( $SD$ ) as defined in Equation 4, see Methods above. Figure 6c shows an  
 440 hexbin scatter plot of hazard indices ( $H_V$  against  $H_P$ ), with a color-coded  
 441 third dimension that corresponds to scene disorder, normalized in the range  
 442  $[0, 1]$ . A first observation is that  $H_P$  and  $H_V$  are positively correlated. More  
 443 interestingly, it is clear that more complex scenes (warmer colors) correspond  
 444 to more dangerous ones. In Figure S5c of the SI, an even clearer trend is  
 445 shown for Madrid. On the other hand, the level of disorder in San Francisco  
 446 scenes is high when  $H_P \approx H_V \approx 1$ , but not clearly related to either  $H_P$   
 447 or  $H_V$  for the rest of values, see Figure S6c. All in all, the connection  
 448 between image complexity and hazard (especially for vehicles) suggests that  
 449 more research is needed in this direction. While certain distractions are  
 450 very explicit (e.g. attending the mobile phone), the perils of scene disorder  
 451 are subtle and implicit (in the sense that they are not obvious on visual

452 inspection).

### 453 3.4. An informed guide to pedestrian safety improvements

454 A precipitate analysis of Figure 6 may render unfeasible interventions:  
455 substitution of built space with larger green areas, building height reduction,  
456 or street widening would suffice to improve pedestrian safety, but they do not  
457 represent a realistic approach. Instead, we resort on the greedy strategy de-  
458 veloped in Section 2.4 to propose interventions conducive to scene alterations  
459 that diminish  $H_P$  and  $H_V$  most.

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

477 Finally, Figure 7d provides a visual overview of the most frequent in-  
478 terventions predicted by our optimization scheme, in the case of Barcelona.  
479 The color of the link connecting two categories expresses the source of that  
480 link. The most notable changes point –perhaps unsurprisingly– to the need  
481 to reconfigure urban scenes towards greener and wider spaces: indeed, both  
482 categories 'road' and 'building' contribute largely to 'nature', while the lat-  
483 ter does the same towards 'sky'. Madrid presents an almost identical trend,  
484 while San Francisco shows a less clear pattern (although the relevance of  
485 'nature' and 'sky' is still clear). Both diagrams are available in the SI, Fig-  
486 ure S10. Overall, the estimations and insights from the panels in Fig. 7 can  
487 provide initial indications to urban planners about achieving potential reduc-  
488 tions of a local hazard score, both in terms of which items could be removed



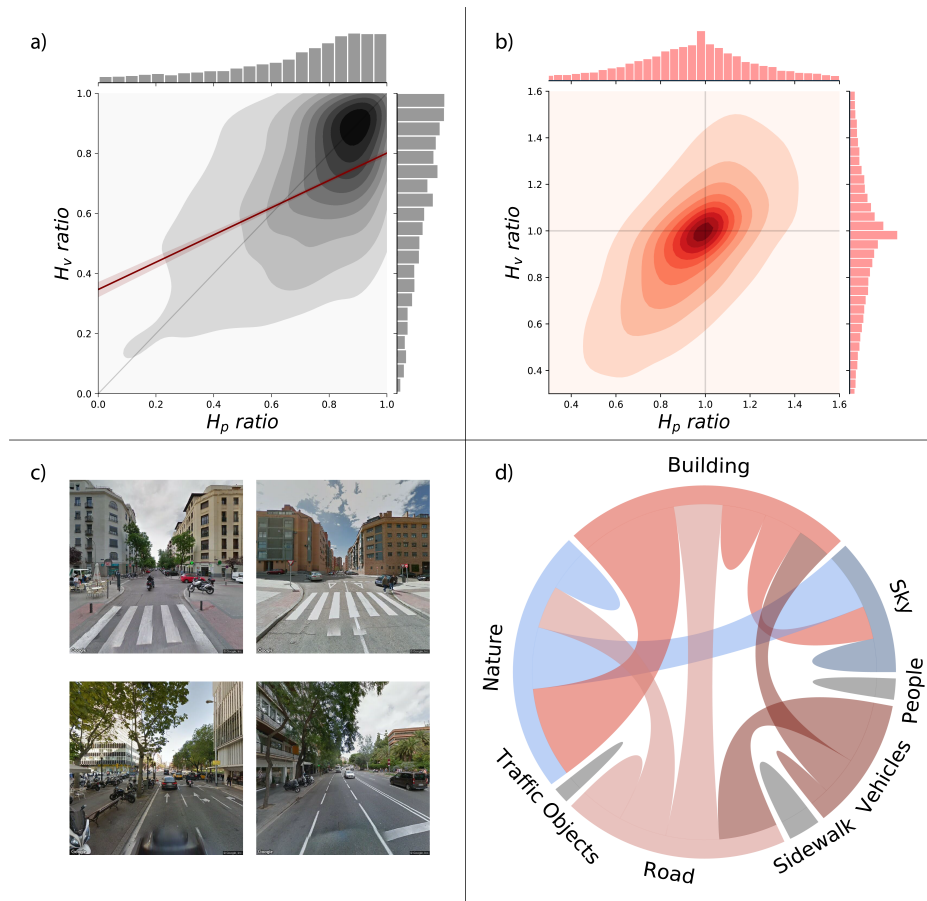


Figure 7: **Hazard reduction: results.** (a) Expected improvement for pedestrian and vehicle hazards, with respect to their original values. The horizontal axis corresponds to the ratio between the improved and the original pedestrian hazard index,  $\tilde{H}_P/H_P$ ; while the vertical axis represents the equivalent ratio for vehicles,  $\tilde{H}_V/H_V$ . Grey 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.

489 or relocated.

#### 490 4. Discussion

491 As cities become increasingly populated, the interactions among pedestri-  
492 ans and motorized vehicles become permanent. This translates into a growing  
493 number of pedestrian-vehicle accidents. Complementary to the efforts by ur-  
494 ban planners, public authorities and sensor technology designers, we present  
495 here an automated scheme that exploits a wide range of Computer Vision  
496 methods (classification, segmentation and interpretability techniques) to re-  
497 duce traffic-related fatalities. The proposed processing pipeline, conveniently  
498 fed with rich sources of open data, renders an holistic characterization of a  
499 city’s hazard landscape, capturing the physical (scene structure) and per-  
500 ceptual (scene complexity) characteristics from a car driver’s point of view.  
501 Beyond its informative value, the hazard landscape provides actionable in-  
502 sights to planners.

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

511 Our approach opens a promising line of development. The hazard land-  
512 scape is defined at an unprecedented, sub-segment resolution level –roughly  
513 a hazard score every 15 meters– through an automated and scalable clas-  
514 sification process. This is well beyond macroscale approaches (e.g. crash  
515 hotspots), and extends the emphasis on intersections [14]. Such fine-grained  
516 map adds a valuable geoinformation layer to those already in use –traffic and  
517 pollution levels [72], land and underground transportation systems, crime,  
518 etc.– enabling better route design: safe paths, along with clean, beautiful, or  
519 shortest ones.

520 Additionally, segmentation and interpretability methods unveil the re-  
521 lationship between potential danger and specific objects in urban scenes.  
522 What’s more, the disposition of those objects is related to hazard indices,  
523 adding a perceptual-attentional link to other possible concomitant variables  
524 that affect vehicle and pedestrian safety. Along this line, our work can be  
525 used in conjunction with other similar pipelines, such as [20], which auto-  
526 mates road safety assessment in terms of infrastructure and estimates road

527 attributes, or may contribute to more focused analysis, relating what a per-  
528 son pays attention to while driving [73]. Additionally, further information  
529 such as temporal accident data, or factors known to influence accident rate  
530 (e.g. weather, lighting condition, distraction, asphalt conditions, road signal-  
531 ing) could be included by using, for instance, a multi-branch convolutional  
532 neural network, to obtain a richer prediction model.

533 On the other hand, the step from descriptive (hazard landscape) to ac-  
534 tionable insights paves the way to automatized, computer-aided prioritization  
535 of urban interventions. The proposed heuristic towards safety improvements  
536 can serve as a novel tool for planners and policy makers, and might trig-  
537 ger the development of more sophisticated approaches such as the use of  
538 Generative Adversarial Networks to produce virtual, plausible alternatives  
539 to target scenes (seeking for instance “safe-fication”, instead of “beautifica-  
540 tion” [25]). These techniques could be complemented with intervention cost  
541 quantification, considering as well cost-safety gain trade-offs.

## 542 Acknowledgements

543 All authors acknowledge financial support from the Dirección General de  
544 Tráfico (Spain), Project No. SPIP2017-02263, as well as TIN2015-66951-C2-  
545 2-R and RTI2018-095232-B- C22 grants from the Spanish Ministry of Science,  
546 Innovation and Universities (FEDER funds). CB and DR acknowledge as  
547 well the support of a doctoral grant from the Universitat Oberta de Catalunya  
548 (UOC). CB, DM and AL acknowledge the NVIDIA Hardware grant program.  
549 Street network data copyrighted OpenStreetMap contributors and available  
550 from <https://www.openstreetmap.org>.

## 551 References

- 552 [1] M. De Domenico, A. Solé-Ribalta, S. Gómez, A. Arenas, Navigability  
553 of interconnected networks under random failures, Proceedings of the  
554 National Academy of Sciences 111 (2014) 8351–8356.
- 555 [2] S. Jiang, Y. Yang, S. Gupta, D. Veneziano, S. Athavale, M. C. González,  
556 The timegeo modeling framework for urban mobility without travel sur-  
557 veys, Proceedings of the National Academy of Sciences 113 (2016)  
558 E5370–E5378.

- 559 [3] S. Abbar, T. Zanouda, J. Borge-Holthoefer, Structural robustness and  
560 service reachability in urban settings, *Data Mining and Knowledge Dis-*  
561 *covery* 32 (2018) 830–847.
- 562 [4] R. Gakenheimer, Urban mobility in the developing world, *Transporta-*  
563 *tion Research Part A: Policy and Practice* 33 (1999) 671 – 689.
- 564 [5] R. Cervero, M. Duncan, Walking, bicycling, and urban landscapes:  
565 Evidence from the san francisco bay area, *American Journal of Public*  
566 *Health* 93 (2003) 1478–1483. PMID: 12948966.
- 567 [6] National Highway Traffic Safety Administration, Fatality analysis re-  
568 porting system (fars) encyclopedia, [https://www-fars.nhtsa.dot.](https://www-fars.nhtsa.dot.gov/Main/index.aspx)  
569 [gov/Main/index.aspx](https://www-fars.nhtsa.dot.gov/Main/index.aspx), 2018. Accessed: 2019-06-27.
- 570 [7] S. Ukkusuri, L. F. Miranda-Moreno, G. Ramadurai, J. Isa-Tavarez, The  
571 role of built environment on pedestrian crash frequency, *Safety Science*  
572 50 (2012) 1141–1151.
- 573 [8] S. M. Rifaat, R. Tay, A. De Barros, Effect of street pattern on the  
574 severity of crashes involving vulnerable road users, *Accident Analysis &*  
575 *Prevention* 43 (2011) 276–283.
- 576 [9] M. Moeinaddini, Z. Asadi-Shekari, M. Z. Shah, The relationship between  
577 urban street networks and the number of transport fatalities at the city  
578 level, *Safety Science* 62 (2014) 114–120.
- 579 [10] G. Mecredy, I. Janssen, W. Pickett, Neighbourhood street connectivity  
580 and injury in youth: a national study of built environments in canada,  
581 *Injury Prevention* 18 (2012) 81–87.
- 582 [11] T. Fu, W. Hu, L. Miranda-Moreno, N. Saunier, Investigating sec-  
583 ondary pedestrian-vehicle interactions at non-signalized intersections  
584 using vision-based trajectory data, *Transportation Research Part C:*  
585 *Emerging Technologies* 105 (2019) 222–240.
- 586 [12] J. Nasar, P. Hecht, R. Wener, Mobile telephones, distracted attention,  
587 and pedestrian safety, *Accident analysis & prevention* 40 (2008) 69–75.
- 588 [13] K. K. Mukoko, S. S. Pulugurtha, Examining the influence of network,  
589 land use, and demographic characteristics to estimate the number of  
590 bicycle-vehicle crashes on urban roads, *IATSS Research* (2019).

- 591 [14] Y. Hu, Y. Zhang, K. S. Shelton, Where are the dangerous intersections  
592 for pedestrians and cyclists: A colocation-based approach, *Transporta-  
593 tion Research Part C: Emerging Technologies* 95 (2018) 431–441.
- 594 [15] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, A. Oliva, Learning deep  
595 features for scene recognition using places database, in: *Advances in  
596 neural information processing systems*, pp. 487–495.
- 597 [16] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, A. Torralba, Places: A  
598 10 million image database for scene recognition, *IEEE transactions on  
599 pattern analysis and machine intelligence* 40 (2017) 1452–1464.
- 600 [17] N. Naik, S. D. Kominers, R. Raskar, E. L. Glaeser, C. A. Hidalgo, Com-  
601 puter vision uncovers predictors of physical urban change, *Proceedings  
602 of the National Academy of Sciences* 114 (2017) 7571–7576.
- 603 [18] A. Albert, J. Kaur, M. C. Gonzalez, Using convolutional networks and  
604 satellite imagery to identify patterns in urban environments at a large  
605 scale, in: *Proceedings of the 23rd ACM SIGKDD international confer-  
606 ence on knowledge discovery and data mining*, ACM, pp. 1357–1366.
- 607 [19] I. Seiferling, N. Naik, C. Ratti, R. Proulx, Green streets- quantifying  
608 and mapping urban trees with street-level imagery and computer vision,  
609 *Landscape and Urban Planning* 165 (2017) 93–101.
- 610 [20] W. Song, S. Workman, A. Hadzic, X. Zhang, E. Green, M. Chen,  
611 R. Souleyrette, N. Jacobs, Farsa: Fully automated roadway safety as-  
612 sessment, in: *2018 IEEE Winter Conference on Applications of Com-  
613 puter Vision (WACV)*, IEEE, pp. 521–529.
- 614 [21] N. Naik, J. Philipoom, R. Raskar, C. Hidalgo, Streetscore-predicting the  
615 perceived safety of one million streetscapes, in: *Proceedings of the IEEE  
616 Conference on Computer Vision and Pattern Recognition Workshops*,  
617 pp. 779–785.
- 618 [22] L. Liu, E. A. Silva, C. Wu, H. Wang, A machine learning-based method  
619 for the large-scale evaluation of the qualities of the urban environment,  
620 *Computers, Environment and Urban Systems* 65 (2017) 113–125.
- 621 [23] T. Gebru, J. Krause, Y. Wang, D. Chen, J. Deng, E. L. Aiden, L. Fei-Fei,  
622 Using deep learning and google street view to estimate the demographic

- 623       makeup of neighborhoods across the united states, Proceedings of the  
624       National Academy of Sciences 114 (2017) 13108–13113.
- 625 [24] E. Suel, J. W. Polak, J. E. Bennett, M. Ezzati, Measuring social, en-  
626       vironmental and health inequalities using deep learning and street im-  
627       agery, Scientific Reports 9 (2019) 6229.
- 628 [25] T. Kauer, S. Joglekar, M. Redi, L. M. Aiello, D. Quercia, Mapping and  
629       visualizing deep-learning urban beautification, IEEE Computer Graph-  
630       ics and Applications 38 (2018) 70–83.
- 631 [26] Z. M. Fadlullah, F. Tang, B. Mao, N. Kato, O. Akashi, T. Inoue,  
632       K. Mizutani, State-of-the-art deep learning: Evolving machine intel-  
633       ligence toward tomorrow’s intelligent network traffic control systems,  
634       IEEE Communications Surveys & Tutorials 19 (2017) 2432–2455.
- 635 [27] L. Zhang, L. Lin, X. Liang, K. He, Is faster r-cnn doing well for  
636       pedestrian detection?, in: European Conference on Computer Vision,  
637       Springer, pp. 443–457.
- 638 [28] S. Zhang, R. Benenson, M. Omran, J. Hosang, B. Schiele, How far are  
639       we from solving pedestrian detection?, in: Proceedings of the IEEE  
640       Conference on Computer Vision and Pattern Recognition, pp. 1259–  
641       1267.
- 642 [29] N. G. Polson, V. O. Sokolov, Deep learning for short-term traffic flow  
643       prediction, Transportation Research Part C: Emerging Technologies 79  
644       (2017) 1–17.
- 645 [30] Y. Wu, H. Tan, L. Qin, B. Ran, Z. Jiang, A hybrid deep learning based  
646       traffic flow prediction method and its understanding, Transportation  
647       Research Part C: Emerging Technologies 90 (2018) 166–180.
- 648 [31] Z. Zhang, Q. He, J. Gao, M. Ni, A deep learning approach for detecting  
649       traffic accidents from social media data, Transportation Research Part  
650       C: Emerging Technologies 86 (2018) 580–596.
- 651 [32] Y. Wang, D. Zhang, Y. Liu, B. Dai, L. H. Lee, Enhancing transportation  
652       systems via deep learning: A survey, Transportation Research Part C:  
653       Emerging Technologies 99 (2019) 144–163.

- 654 [33] Z. Zhang, M. Li, X. Lin, Y. Wang, F. He, Multistep speed prediction on  
655 traffic networks: A deep learning approach considering spatio-temporal  
656 dependencies, *Transportation Research Part C: Emerging Technologies*  
657 105 (2019) 297–322.
- 658 [34] Ayuntamiento de Madrid, Portal de datos abiertos del ayuntamiento  
659 de madrid, <https://datos.madrid.es/portal/site/egob/>, 2019. Ac-  
660 cessed: 2019-04-20.
- 661 [35] Ajuntament de Barcelona, Open data bcn, [https://  
662 opendata-ajuntament.barcelona.cat/en/](https://opendata-ajuntament.barcelona.cat/en/), 2019. Accessed: 2019-04-  
663 20.
- 664 [36] Safe Transportation Research and Education Center, University of Cal-  
665 ifornia, Berkeley, Transportation injury mapping system (tims), 2019.  
666 Accessed: 2019-06-27.
- 667 [37] D. Anguelov, C. Dulong, D. Filip, C. Frueh, S. Lafon, R. Lyon, A. Ogale,  
668 L. Vincent, J. Weaver, Google street view: Capturing the world at street  
669 level, *Computer* 43 (2010) 32–38.
- 670 [38] Mapillary contributors, Mapillary - Street-level imagery, powered by  
671 collaboration and computer vision , <https://www.mapillary.com/app>,  
672 2019.
- 673 [39] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, *Nature* 521 (2015) 436.
- 674 [40] J. Schmidhuber, Deep learning in neural networks: An overview, *Neural  
675 networks* 61 (2015) 85–117.
- 676 [41] K. He, X. Zhang, S. Ren, J. Sun, Identity mappings in deep residual  
677 networks, in: *European conference on computer vision*, Springer, pp.  
678 630–645.
- 679 [42] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with  
680 deep convolutional neural networks, in: *Advances in neural information  
681 processing systems*, pp. 1097–1105.
- 682 [43] A. Adadi, M. Berrada, Peeking inside the black-box: A survey on ex-  
683 plainable artificial intelligence (xai), *IEEE Access* 6 (2018) 52138–52160.

- 684 [44] H. Fukui, T. Hiraakawa, T. Yamashita, H. Fujiyoshi, Attention branch  
685 network: Learning of attention mechanism for visual explanation, in:  
686 Proceedings of the IEEE Conference on Computer Vision and Pattern  
687 Recognition, pp. 10705–10714.
- 688 [45] J. Wagner, J. M. Kohler, T. Gindele, L. Hetzel, J. T. Wiedemer,  
689 S. Behnke, Interpretable and fine-grained visual explanations for con-  
690 volutional neural networks, in: Proceedings of the IEEE Conference on  
691 Computer Vision and Pattern Recognition, pp. 9097–9107.
- 692 [46] S. Desai, H. G. Ramaswamy, Ablation-cam: Visual explanations for  
693 deep convolutional network via gradient-free localization, in: 2020 IEEE  
694 Winter Conference on Applications of Computer Vision (WACV), IEEE,  
695 pp. 972–980.
- 696 [47] B. N. Patro, M. Lunayach, S. Patel, V. P. Namboodiri, U-cam: Visual  
697 explanation using uncertainty based class activation maps, in: Pro-  
698 ceedings of the IEEE International Conference on Computer Vision, pp.  
699 7444–7453.
- 700 [48] Z. Wang, J. Yang, Diabetic retinopathy detection via deep convolutional  
701 networks for discriminative localization and visual explanation, arXiv  
702 preprint arXiv:1703.10757 (2017).
- 703 [49] H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid scene parsing net-  
704 work, in: Proceedings of the IEEE Conference on Computer Vision and  
705 Pattern Recognition, pp. 2881–2890.
- 706 [50] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benen-  
707 son, U. Franke, S. Roth, B. Schiele, The cityscapes dataset for semantic  
708 urban scene understanding, in: Proceedings of the IEEE conference on  
709 computer vision and pattern recognition, pp. 3213–3223.
- 710 [51] R. M. Haralick, K. Shanmugam, I. H. Dinstein, Textural features for  
711 image classification, IEEE Transactions on Systems, Man, and Cyber-  
712 netics (1973) 610–621.
- 713 [52] N. Moray, Attention in dichotic listening: Affective cues and the influ-  
714 ence of instructions, Quarterly journal of experimental psychology 11  
715 (1959) 56–60.



- 716 [53] D. Kahneman, *Attention and effort*, volume 1063, Citeseer, 1973.
- 717 [54] G. A. Alvarez, P. Cavanagh, The capacity of visual short-term mem-  
718 ory is set both by visual information load and by number of objects,  
719 *Psychological science* 15 (2004) 106–111.
- 720 [55] J. E. Richards, The development of attention to simple and complex  
721 visual stimuli in infants: Behavioral and psychophysiological measures,  
722 *Developmental Review* 30 (2010) 203–219.
- 723 [56] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Learning deep  
724 features for discriminative localization, in: *Proceedings of the IEEE*  
725 *conference on computer vision and pattern recognition*, pp. 2921–2929.
- 726 [57] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra,  
727 Grad-cam: Visual explanations from deep networks via gradient-based  
728 localization, in: *Proceedings of the IEEE International Conference on*  
729 *Computer Vision*, pp. 618–626.
- 730 [58] A. Chattopadhyay, A. Sarkar, P. Howlader, V. N. Balasubramanian,  
731 Grad-cam++: Generalized gradient-based visual explanations for deep  
732 convolutional networks, in: *2018 IEEE Winter Conference on Applica-*  
733 *tions of Computer Vision (WACV)*, IEEE, pp. 839–847.
- 734 [59] C. Ventura, D. Masip, A. Lapedriza, Interpreting cnn models for appar-  
735 ent personality trait regression, in: *Proceedings of the IEEE Conference*  
736 *on Computer Vision and Pattern Recognition Workshops*, pp. 55–63.
- 737 [60] P. Harrington, *Machine learning in action*, Manning Publications Co.,  
738 2012.
- 739 [61] E. Frank, M. Hall, A simple approach to ordinal classification, in:  
740 *European Conference on Machine Learning*, Springer, pp. 145–156.
- 741 [62] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-  
742 scale image recognition, *arXiv preprint arXiv:1409.1556* (2014).
- 743 [63] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the  
744 inception architecture for computer vision, in: *Proceedings of the IEEE*  
745 *conference on computer vision and pattern recognition*, pp. 2818–2826.

- 746 [64] C. Szegedy, S. Ioffe, V. Vanhoucke, A. A. Alemi, Inception-v4, inception-  
747 resnet and the impact of residual connections on learning, in: Thirty-  
748 first AAAI conference on artificial intelligence.
- 749 [65] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang,  
750 T. Weyand, M. Andreetto, H. Adam, Mobilenets: Efficient convolu-  
751 tional neural networks for mobile vision applications, arXiv preprint  
752 arXiv:1704.04861 (2017).
- 753 [66] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recog-  
754 nition, in: Proceedings of the IEEE Conference on Computer Vision and  
755 Pattern Recognition, pp. 770–778.
- 756 [67] R. Louf, M. Barthelemy, A typology of street patterns, *Journal of The*  
757 *Royal Society Interface* 11 (2014) 20140924.
- 758 [68] D. Quercia, R. Schifanella, L. M. Aiello, The shortest path to happiness:  
759 Recommending beautiful, quiet, and happy routes in the city, in: Pro-  
760 ceedings of the 25th ACM conference on Hypertext and social media,  
761 ACM, pp. 116–125.
- 762 [69] H. Huang, M. A. Abdel-Aty, A. L. Darwiche, County-level crash risk  
763 analysis in florida: Bayesian spatial modeling, *Transportation Research*  
764 *Record* 2148 (2010) 27–37.
- 765 [70] P. Chen, J. Zhou, Effects of the built environment on automobile-  
766 involved pedestrian crash frequency and risk, *Journal of Transport &*  
767 *Health* 3 (2016) 448–456.
- 768 [71] P. Olszewski, I. Buttler, W. Czajewski, P. Dabkowski, C. Kraśkiewicz,  
769 P. Szagała, A. Zielińska, Pedestrian safety assessment with video anal-  
770 ysis, *Transportation Research Procedia* 14 (2016) 2044–2053.
- 771 [72] Y. Xu, S. Jiang, R. Li, J. Zhang, J. Zhao, S. Abbar, M. C. González,  
772 Unraveling environmental justice in ambient pm<sub>2.5</sub> exposure in beijing:  
773 A big data approach, *Computers, Environment and Urban Systems* 75  
774 (2019) 12–21.
- 775 [73] A. Palazzi, D. Abati, F. Solera, R. Cucchiara, et al., Predicting the  
776 driver’s focus of attention: the dr (eye) ve project, *IEEE Transactions*  
777 *on Pattern Analysis and Machine Intelligence* 41 (2018) 1720–1733.

- 778 [74] OpenStreetMap contributors, Planet dump retrieved from  
779 <https://planet.osm.org> , <https://www.openstreetmap.org>, 2017.
- 780 [75] G. Boeing, Osmnx: New methods for acquiring, constructing, analyzing,  
781 and visualizing complex street networks, *Computers, Environment and*  
782 *Urban Systems* 65 (2017) 126–139.