

UNDERSTANDING CITIES WITH MACHINE EYES:

A review of deep computer vision in urban analytics

Mohamed R Ibrahim - Manuscript – AUG 2019

Abstract—Modelling urban systems has interested planners and modellers for decades. Different models have been achieved relying on mathematics, cellular automation, complexity, and scaling. While most of these models tend to be a simplification of reality, today within the paradigm shifts of artificial intelligence across the different fields of science, the applications of computer vision show promising potential in understanding the realistic dynamics of cities. While cities are complex by nature, computer vision shows progress in tackling a variety of complex physical and non-physical visual tasks. In this article, we review the tasks and algorithms of computer vision and their applications in understanding cities. We attempt to subdivide computer vision algorithms into tasks, and cities into layers to show evidence of where computer vision is intensively applied and where further research is needed. We focus on highlighting the potential role of computer vision in understanding urban systems related to the built environment, natural environment, human interaction, transportation, and infrastructure. After showing the diversity of computer vision algorithms and applications, the challenges that remain in understanding the integration between these different layers of cities and their interactions with one another relying on deep learning and computer vision. We also show recommendations for practice and policy-making towards reaching AI-generated urban policies.

Keywords—Cities, computer vision, deep learning, Convolutional Neural Networks (CNN), urban studies

1. INTRODUCTION

Cities are complex entities by nature and modelling urban systems has interested planners for decades (Batty, 2008; Bettencourt, 2013; Isalgue, Coch, & Serra, 2007). A range of approaches have been used to model urban processes, examples of which include cellular automata (Batty, 1997; Batty, Couclelis, & Eichen, 1997; de Almeida et al., 2003), fractals (Batty and Longley, 1994; Batty and Xie, 1996; Frankhauser, 1998; Murcio et al., 2015) and multi-agent models (Batty, 2005; Heppenstall, Crooks, See, & Batty, 2012). These models aim to understand cities by modelling their underlying components and exploring their systems, ultimately intending to inform decision making and policy (Batty, 2009; Calder et al., 2018). Due to the complexity and nonlinearity of cities, these models tend to explore or predict urban systems in a sectoral fashion. For example, transport models are used to simulate the potential impact of policy and infrastructure investment. Such models may fail to represent complex events in cities, in which multiple systems interact.

The success of deep learning and computer vision in pattern recognition over the past decade (LeCun, Bengio, & Hinton, 2015) has created opportunities to understand cities through images (Reichstein et al., 2019). So far, the diversity of the algorithms of computer vision has enabled researchers to tackle and predict a wide spectrum of issues in more accurate and precise fashion (Goodfellow, Bengio, & Courville, 2017; LeCun et al., 2015; Reichstein et al., 2019).

In this paper, we review the algorithms and applications of computer vision related to urban analytics. Urban analytics

can be defined as urban research that exploits new data resources that are captured, for example, from sensors (e.g. imagery, the internet of things), crowdsources and social media (Batty, 2019). Deep learning and computer vision technologies have tremendous potential in this area for dealing with heterogeneous data types, many of which are image-based. In the review, we identify the areas that have been intensively modelled using computer vision while also revealing the areas in which further research is needed. This is achieved by categorising the application areas of urban analytics into five layers of the city (the built environment, human interaction, transportation and traffic, the natural environment, and infrastructure). In doing so, we demonstrate that, while many urban processes are a result of interactions across these layers, the current approach is to tackle these layers differently and separately. Here, we note the potential of extracting data of different disciplines using a unified input (images/videos) that relies on computer vision methods to cover a wide spectrum of urban and transport research.

This review aims to provide a resource for urban planners and practitioners by: 1) reviewing the main methodologies of computer vision, and their applicability to various tasks of urban analytics, 2) illustrating the variation and nuances of deep learning and computer vision algorithms and their limitations in understanding cities, 3) giving a descriptive understanding of the algorithms of computer vision for policy-makers and planners, and how they are used in cities, 4) paving the way for developing AI-generated urban policies by highlighting the key enabling technologies and research directions. The remainder of this review is structured as follows: In section 2, the methodology of the review is described. In section 3, the key tasks of computer vision are described, along with the main algorithms. The applications of computer vision in urban analytics are reviewed in section 4. Section 5 summarises what remains missing in current research, before section 6 shows how we can move from prediction to decision making and policy recommendation. Finally, some conclusions are given in section 6.

2. REVIEW METHODOLOGY

The methodology of this review is divided into two parts: 1) manuscripts are collected that summarise the progress in deep learning methods and algorithms that are applicable to computer vision tasks, 2) manuscripts are collected that reflect the application of deep learning and computer vision in understanding cities in the last decade (since 2010). For the first part, we present only the major methodological approaches. Papers that vary or improve on these main approaches are excluded. Most of these studies are presented in premier computer science conferences, including, but not limited to CVPR, ICCV, ECCV and NeurIPS, or in ArXiv. For the second part, we extend the search to peer-reviewed journals and conference proceedings listed in Scopus, Web of Science, Google Scholar and Science Direct, that can be

accessed via a combination of keywords such as: deep learning, cities, computer vision, land-use modelling, urban perception, prediction, detection, street-level images, aerial or satellite images. This is because the applied computer vision literature is often found in domain specific journals, rather than computer science conferences.

In total, 641 manuscripts were collected to cover the two parts of the methodology. For the second part, the collected manuscripts were filtered to include only those related to computer vision of street-level or aerial images, which use deep learning or hybrid models that include a convolutional structure. Studies that involve deep learning of other data types such as 2D/3D LIDAR data are excluded. Studies that use classical machine learning or computer vision algorithms without involving deep learning are also excluded, except where they are required to draw a baseline to emphasise advancement or contrast. The algorithms are presented at a descriptive level and readers are referred to the relevant literature for further details.

3. THE BASICS AND TASKS OF COMPUTER VISION

Before exploring the domains where computer vision is applied in cities, it is worth identifying first what computer vision is and what its algorithms are capable of achieving from a generic perspective. Computer vision can be narrowed to the task of learning the qualitative representation of visual elements in their raw form in order to quantify them (LeCun et al., 2015). Similar to human eyes, the computer sees visual objects and creates a cognitive understanding of a scene based on a sequential sample of the presented images or frames of images in a task-specific manner. While computer vision is not new (i.e. Viola & Jones, 2001), deep learning, most specifically Convolutional Neural Networks (CNN), has made it possible for computer vision to tackle various issues and process images more precisely and efficiently (K. He, Zhang, Ren, & Sun, 2015; LeCun et al., 2015). These deep models, computation capabilities, and the availability of large datasets have made it possible for computer vision to permeate a wide range of applications in realistic settings (Cordts et al., 2016; T.-Y. Lin et al., 2014; Russakovsky et al., 2015). Generally, the logic of computer vision, relying on these deep models, can be summarized as the construction of multiple hidden layers that are capable of accomplishing a range of vision tasks by extracting digital features that may or may not be recognisable to human eyes (Y. Guo et al., 2016; Kuo, 2016; LeCun et al., 2015). The most commonly used are convolutional, pooling, flatten, and fully-connected layers. The general functions of these layers can be summarised as follows:

- Convolutional layers are responsible for extracting features coupled with activation functions, such as Rectified linear units (ReLU), to add nonlinearity to the model,
- Pooling layers are responsible for reducing the dimensionality of the data,
- Flatten layers are responsible for converting the features of the model into neurons to be fed forward to the fully-connected layers
- Fully-connected layers aim to adjust the weights and predict the output for a given task.

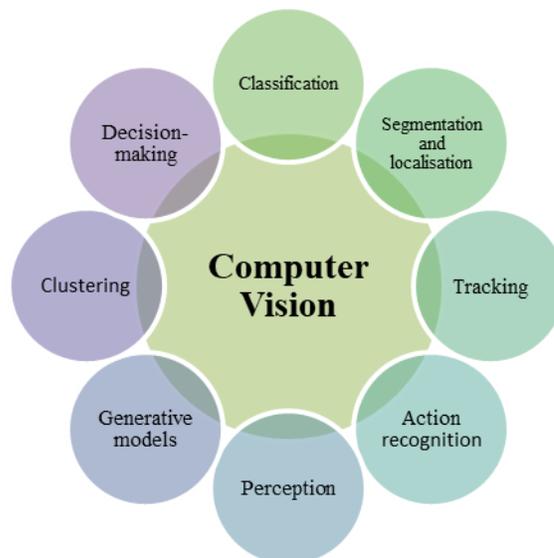


FIG. 1
COMPUTER VISION TASKS

The types, numbers, and orders of these layers are responsible for determining functionality and the optimisation of both accuracy and time needed for the training and the inference of the model. The structure of the model and the fine-tuning of the various hyperparameters represents the innovation and the advancements of the state-of-the-art for pattern recognition for a given task (LeCun et al., 2015).

Depending on the type of visual task, deep models can be trained differently with different layers and different sets of algorithms (Y. Guo et al., 2016). As shown in Fig. 1, these algorithms of computer vision can be subdivided based on eight fundamental tasks, upon which other tasks can be framed and built. These are; image classification, segmentation and localisation, tracking, action-recognition, perception, generative models, clustering, and decision-making. Table 1 shows the literature related to different computer vision tasks. It expands on the methods related to each task and their subcategories.

TABLE 1: METHODS RELATED TO THE TASKS OF COMPUTER VISION

VISION TASK	SUB-CATEGORY	METHOD	
CLASSIFICATION		ALEXNET	(krizhevsky, sutskever, & hinton, 2012)
		VGGNET	(simonyan & zisserman, 2014)
		GOOGLNET	(szegedy, liu, jia, sermanet, & reed, 2015)
		RESNET	(he et al., 2015)
		DENSNET	(huang, liu, weinberger, & van der maaten, 2017)
SEGMENTATION AND LOCALISATION	OBJECT-BASED DETECTION	R-CNN	(Girshick, Donahue, Darrell, & Malik, 2014)
		FAST R-CNN	(Ren, He, Girshick, & Sun, 2016)

		YOLO SSD YOLOV2 YOLOV3 RETINANET	(Redmon, Divvala, Girshick, & Farhadi, 2016) (Liu et al., 2016) (Redmon & Farhadi, 2017) (Redmon & Farhadi, 2018) (Lin, Goyal, Girshick, He, & Dollár, 2018)
	SEMANTIC SEGMENTATION	DEEPLAB U-NET SEGNET - - - - REFINENET - FOVEANET LINKNET -	(Chen, Papandreou, Kokkinos, Murphy, & Yuille, 2016) (Ronneberger, Fischer, & Brox, 2015) (Badrinarayanan, Kendall, & Cipolla, 2016) (Long, Shelhamer, & Darrell, 2015) (Peng, Zhang, Yu, Luo, & Sun, 2017) (Chen, Papandreou, Kokkinos, Murphy, & Yuille, 2016) (Zhao, Shi, Qi, Wang, & Jia, 2017) (Yu & Koltun, 2015) (Lin, Milan, Shen, & Reid, 2017) (Chen, Papandreou, Schroff, & Adam, 2017) (Jégou, Drozdal, Vazquez, Romero, & Bengio, 2016) (Li, Wang, et al., 2017) (Chaurasia & Culurciello, 2017) (Yang, Yu, Zhang, Li, & Yang, 2018)
TRACKING OBJECTS		MOTS - - - - - ECO CNNTRACKER ARTTRACK - - PATHTRACK	(Voigtlaender, Krause, Sekar, Geiger, & Leibe, 2019) (Jiang, Xiao, Xie, Tillo, & Huang, 2018) (Kang, Ouyang, Li, & Wang, 2016) (Girdhar et al., 2017) (Danelljan, Hager, Khan, & Felsberg, 2015) (Held, Thrun, & Savarese, 2016) (Danelljan, Bhat, Khan, & Felsberg, 2016) (Y. Chen et al., 2016) (Insafutdinov, Andriluka, et al., 2016) (Wu, Lu, Gao, Zhao, & Liu, 2016) (Chu et al., 2017) (Manen, Gygli, Dai, & Gool, 2017)
ACTION RECOGNITION	HUMAN POSE ESTIMATION	DENSEPOSE MULTIPOSENET - RMPE DEEPCUT - -	(Guler, Neverova, & Kokkinos, 2018) (Kocabas, Karagoz, & Akbas, 2018) (Papandreou et al., 2017) (Fang, Xie, Tai, & Lu, 2016) (Insafutdinov, Pishchulin, Andres, Andriluka, & Schiele, 2016) (Cao, Simon, Wei, & Sheikh, 2016) (Pflister, Charles, & Zisserman, 2015)
	ACTION CLASSIFICATION	- - - -	(Girdhar & Ramanan, 2017) (Bilen, Fernando, Gavves, Vedaldi, & Gould, 2016) (Zhu, Lan, Newsam, & Hauptmann, 2017) (Guo et al., 2018) (Zhang, Wang, Wang, Qiao, & Wang, 2016)
	TEMPORAL ACTION DETECTION	DAPS - - - - - - - - -	(Escorcia, Caba Heilbron, Niebles, & Ghanem, 2016) (Diba et al., 2017) (Gemert, Jain, Gati, & Snoek, 2015) (Shou, Chan, Zareian, Miyazawa, & Chang, 2017) (Escorcia et al., 2016) (Li et al., 2016) (Xu, Das, & Saenko, 2017) (Chao et al., 2018) (Buch, Escorcia, Shen, Ghanem, & Niebles, 2017) (Zhao et al., 2017)
	SPATIO-TEMPORAL ACTION DETECTION	- - - - - - - - -	(Chen & Corso, 2015) (Becattini, Uricchio, Seidenari, Del Bimbo, & Ballan, 2017) (Saha, Singh, & Cuzzolin, 2017) (Gemert et al., 2015) (Zhu, Vial, & Lu, 2017) (El-Nouby & Taylor, 2018) (Saha, Singh, Sapienza, Torr, & Cuzzolin, 2016) (Singh, Saha, Sapienza, Torr, & Cuzzolin, 2016) (Mettes, van Gemert, & Snoek, 2016)

		-	(Weinzaepfel, Harchaoui, & Schmid, 2015)
PERCEPTION	UNDERSTANDING SCENES		(Eslami et al., 2018)
	ESTIMATING DEPTH	-	(Cao, Wu, & Shen, 2017) (He, Wang, & Hu, 2018)
GENERATIVE MODELS	GANS	-	(Goodfellow et al., 2014)
		-	(Radford, Metz, & Chintala, 2015)
		-	(Reed et al., 2016)
		STACKGAN	(Zhang et al., 2016)
		-	(Isola, Zhu, Zhou, & Efros, 2016)
		BIGGAN	(Brock, Donahue, & Simonyan, 2018)
CLUSTERING		-	(Caron, Bojanowski, Joulin, & Douze, 2018)
		-	(Xie, Girshick, & Farhadi, 2016)
		DEEPCUSTER	(Tian, Zhou, & Guan, 2017)
MAKING DECISIONS	DEEP Q-LEARNING	-	(Mnih et al., 2013)
		-	(Hester et al., 2017)
	DOUBLE DEEP Q-LEARNING	-	(van Hasselt, Guez, & Silver, 2015)
	DUEL DEEP Q-LEARNING	-	(Wang et al., 2015)
	A3C	-	(Mnih et al., 2016)

3.1 Classification

Deep learning models, most specifically Convolutional Neural Networks (CNN), have shown substantial progress in classifying images of a wide spectrum of classes (LeCun et al., 2015). Various deep CNN models with different architectures and hyper-parameters have been computed to recognize visual objects in large repositories of images, such as the ImageNET dataset that contains 15 million images that belong to 22,000 different classes (Russakovsky et al., 2015, 2015). Starting with AlexNet (Krizhevsky et al., 2012), VGGNet (Simonyan & Zisserman, 2014), GoogLeNet (Szegedy et al., 2015), ResNet (K. He et al., 2015) and most recently, DenseNet (Huang et al., 2017), these CNN models are able to accurately recognize and classify a wide range of images. For instance, ResNet-152 achieved 4.49% top-5 error score on the validation set of ImageNET (K. He et al., 2015).

3.2 Segmentation and localisation

Segmentation and localisation are the processes of identifying multiple objects in a single image. These models use a single deep model in an end-to-end fashion, in which the first part of the model is an image classifier followed by different types of layers to localise different objects with a given confidence. Notable examples include the Region-based CNN model (R-CNN) (Girshick et al., 2014), Fast R-CNN (Ren et al., 2016), You Only Look Once (YOLO) (Redmon & Farhadi, 2017, 2018) and the MultiBox Detectors for fast image segmentation, or so-called; Single Shot Multi-Box Detector (SSD) technique (W. Liu et al., 2016). CNN models have shown significant progress in recognising and detecting objects in images with a minimal inference time and high overall validation accuracy. YOLOv3 achieves 93.8% top-5 score on the COCO dataset (Redmon & Farhadi, 2018).

For further explanation related to localisation and object detection, see (Zou, Shi, Guo, & Ye, 2019).

3.3 Tracking objects

After building a system of object detection, computer vision can be used for tracking multiple objects in a complex scene by adding features that correlate a pair of consecutive frames.

This tracker system is capable of identifying a candidate box at each frame-level jointly with their time deformations (Girdhar, Gkioxari, Torresani, Paluri, & Tran, 2017). While different tracker systems can be built based on correlation filtering and online learning techniques between consecutive frames (X. Zhang, Xia, Lu, Shen, & Zhang, 2018), the state-of-the-art research in object tracking uses an end-to-end CNN model to tackle both detection and tracking, which can add more advanced features (i.e. dealing with occlusion issues) for tracking various elements (Girdhar et al., 2017; Hou, Chen, & Shah, 2017; K. Kang, Ouyang, Li, & Wang, 2016). For further explanation related to deep visual tracking, see P. Li, Wang, Wang, & Lu (2018).

3.4 Action recognition

Computer vision coupled with deep CNN models is not only capable of tracking the motion of an object in a complex scene, but also classifying its multiple actions while tracking (Bilen et al., 2016; Limin Wang, Qiao, & Tang, 2015; B. Zhang et al., 2016). Various computer vision algorithms have been developed to tackle humans poses and their interaction with an external object in a complex scene (El-Nouby & Taylor, 2018; Saha et al., 2016; Soomro & Shah, 2017; Weinzaepfel, Martin, & Schmid, 2016). 2D or 3D convolution layers (with or without the spatiotemporal dimensions) can identify the action of the object from its pose in relation to another target object. For instance, from the pose of a person sitting on a bike, the algorithms of computer vision can identify cycling as an action. This concept of the triplet inputs (object, verb, target) has been seminal for tackling real-world events and behaviours, from a simple still image to multi-frame images (Girdhar et al., 2017).

3.5 Perception

Perception tasks can be seen as classification or regression tasks that predict information that is not necessarily embedded directly in the image but can be inferred from the overall structure of the image. Perceiving a neighbourhood as safe or unsafe for example can be seen as a perception task, in which the machine extracts features from the structure of an image to classify the safety of the image. Even though understanding

the overall gist of a scene is seminal for understanding more than an object in an image (Oliva & Torralba, 2006), few works have been done in this domain. The complexity of tackling this subject lies in sensing the class of an image by sensing the overall profound features of the image, rather than identifying an object in the image. For instance, identifying and sensing the planning status of a region from the image (Ibrahim, Haworth, & Cheng, 2019).

Moreover, seeing what is far and what is close just by looking at a still image is another advantage of computer vision relying on deep CNN models. Cao, Wu, & Shen (2017) trained deep CNN models to estimate the depth in a single image by labelling the different depths on the image and dealing with training the model as a classification task. In contrast, He, Wang, & Hu (2018) trained a deep CNN model to estimate the depth of a monocular image relying on the information of focal length that has proven to outperform the other state-of-the-art depth estimation algorithms based on deep learning models.

3.6 Generative models

Generative models refer to the ones that tend to output synthesized data by learning the representation of their input data in an unsupervised fashion, conditionally or unconditionally.

There is a range of algorithms that are classified as generative models, such as Restricted Boltzmann Machine (RBM), deep belief networks, Autoencoders, and Generative adversarial Networks (GANs) (Goodfellow et al., 2017). This section refers only to GANs, which generate synthetic graphical data in an unsupervised training fashion relying on images as input. Unlike other tasks related to computer vision, the deep models of GANs, introduced in 2014, enable machines to generate new information that is similar to what the model has been trained to identify (Goodfellow et al., 2014). In other words, if the model is trained on images of trees, by using GANs the model can generate a new image of a tree that preserves the fundamental features of a tree, but with a new visual identity. This progress of deep learning enables the creation of unique objects or scenes by understanding the underlying features of the trained images or videos.

GANs are trained differently from the abovementioned deep models, not only in term of layers but rather, instead of the single end-to-end model, two parallel deep models are trained that compete with one another (Goodfellow, 2016; Goodfellow et al., 2014; Radford et al., 2015). The first one, the Generator model, generates new images to deceive the second model that holds the ground truth data, while the second model, the Discriminator model, blocks this new image until the generator model becomes advanced enough to generate new images that are similar enough to the ground truth that the discriminator model can no longer refuse them. This computationally intensive training, in an unsupervised manner, opens the door for computer-based creativity without the prior supervision of humans.

GANs have been utilised in various applications. Isola, Zhu, Zhou, & Efros (2016) used conditional GANs to translate from one form of an image to another. For instance, by giving the model a satellite image of a location, the model can give the semantic segmentation of the location or vice versa. Zhang

et al. (2016) created stackGAN model to transform a text description of an image into a photo-realistic synthesis. Moreover, Reed et al. (2016) have pushed the algorithms of GANs further. The machines can learn to draw not only from text distributions but also by telling the machine what and where to draw on the canvas. Apart from the daily-life applications, GANs have been used in the simulation of 3D energy particle showers and physics-related applications (Paganini, de Oliveira, & Nachman, 2018).

3.7 Clustering

Clustering is a form of unsupervised learning, in which the machines are able to cluster different still images or multi-frame images based on their content or embedded objects without prior human supervision (Caron et al., 2018; Tian et al., 2017; Xie et al., 2016). So far, different computer vision algorithms have been developed to tackle this task and eliminate the need for a long process of manual labelling from still images. Recently, Eslami et al. (2018) introduced the Generative Query Network (GQN) for scene representation without human supervision. The GQN takes images from a different perspective as an input and generates a visual representation of the scene from an unobserved perspective. This process of coupling generative models with clustering introduces a new form of machine intelligence to understand scene representation without human supervision.

3.8 Decision-making

By looking at the edge of computer vision and coupling its deep models with reinforcement learning, or so-called Deep Reinforcement Learning (DRL), machines can be trained to explore and compute the outcomes of different scenarios in order to make an real-time decisions based on visual aspects of the environment (Hester et al., 2017; Mnih et al., 2016). This level of cognitive ability of machines by applying one or more of the abovementioned tasks can enable an agent to grasp information and interact with an environment to optimize target resources without human supervision.

Due to the complexity of the algorithms related to this subject, most examples are in virtual or gaming environments (Mnih et al., 2013). However, most significantly, Mirowski et al. (2018) utilised DRL to enable a machine to navigate through the unstructured environment of the street network relying on street-level images. In this work, the machine learns to navigate by understanding landmarks from images and to determine its location and its target destination.

4. RECOGNISING THE URBAN WORLD

Understanding the dynamics of cities remains a complex issue. Data collection, for instance, is one of the crucial domains where automation is highly desirable, in which computer vision has been successfully applied in capturing and analysing various objects depicted in urban scenes. Specifically, scene parsing and semantic segmentation represent crucial tasks of computer vision for a better understanding of the elements of an urban scene. From images, computer vision can localize multiple objects in cities, or simply segment the entire scene based on a group of themes, such as sky, ground, road, building, vegetation, etc

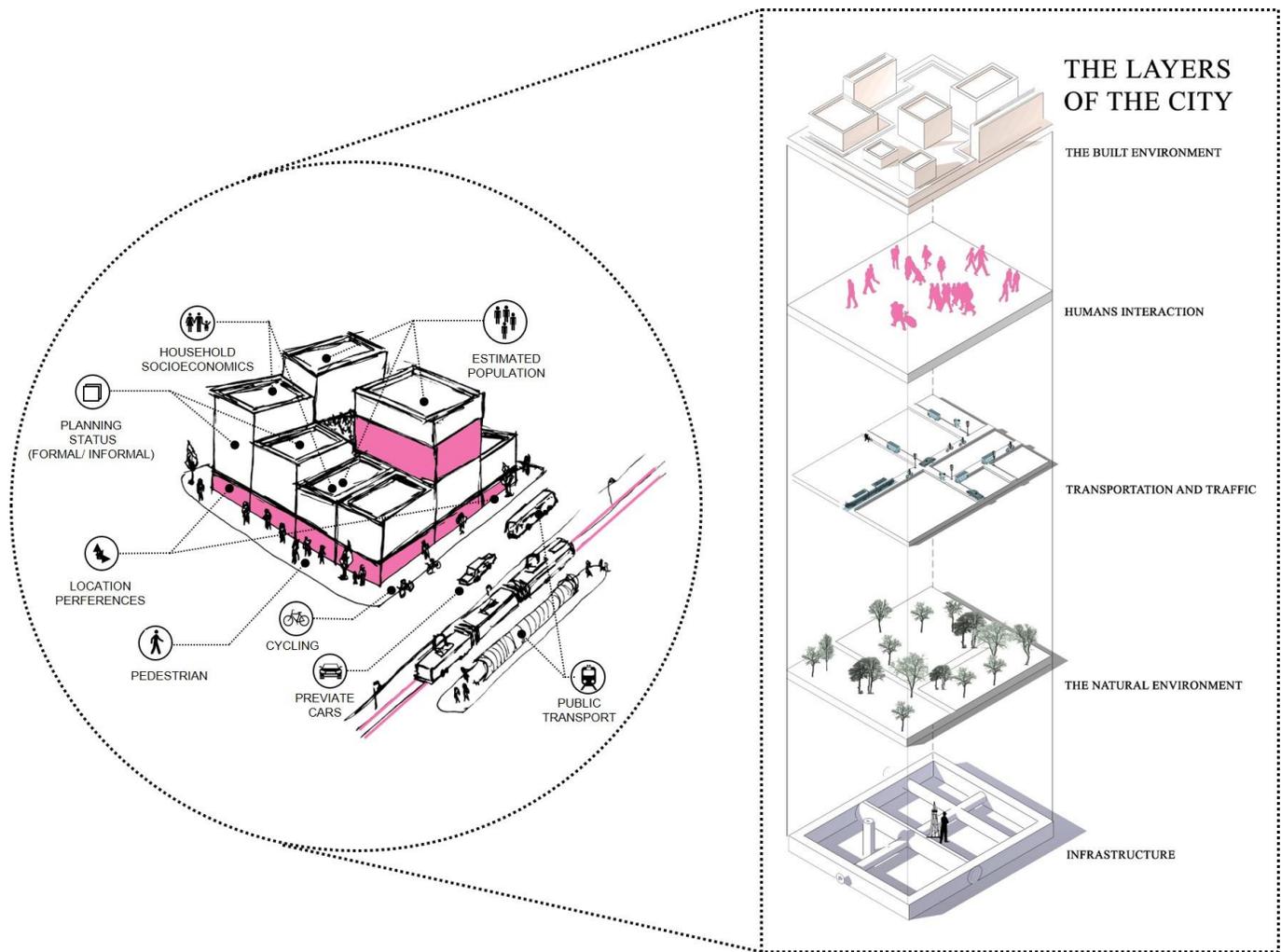


FIG. 2
THE LAYERS OF THE CITIES WHERE COMPUTER VISION HAS BEEN APPLIED
CREATED BY THE AUTHORS

(Chaurasia & Culurciello, 2017; Zhou et al., 2017). Putting all the above-mentioned tasks together, computer vision shows good potential in urban analytics for analysing the multi-layers of cities. For the purposes of this review, we define these layers as; the built environment, the natural environment, humans and their physical interactions, transport modes and traffic-related issues, and infrastructure. The main reason for breaking-down cities in these layers is to be able to tackle the applications of computer vision in each individual field of science related to urban analytics, in which the methods, scope, language used, and the nature of work may vary depending on the discipline. For instance, research that has been done in understanding the built environment may vary in nature from that done to understand transportation, even though the methods of deep learning and computer vision may be similar.

Fig. 2 shows examples of computer vision applications in cities to detect multidisciplinary tasks that belong to the five layers of cities, whereas table 2 shows the applications of

computer vision to these layers. Each layer is broken down into further subcategories as appropriate.

4.1 The built environment

This section addresses cities from an architectural and urban design perspective, for example, understanding cities from a land-use perspective, the level of the physical appearance of the street-level that may indicate or measure housing prices, or even the level of safety with a certain neighbourhood.

When it comes to understanding the built environment, there are different challenges that face urban planners and policy-makers. For example, modelling the physical appearance of complex urban areas is a multi-faceted issue that is vital for planners and policy-makers for making decisions for improving living conditions in cities. The collection of data that reflects the current status of the built environment is a critical issue for urban analytics. So far, the applications of computer vision have merged not only to detect various urban components but also to understand the appearance and the safety factors of an urban scene. While there is a wide range of applications of computer vision in cities, these applications

can be divided into two approaches that either analyse cities from street-level images or remote sensing data such as satellite images.

4.1.1 Seeing cities from above

Analysing cities from above relying on remote sensing and geographical information systems (GIS), perhaps, is the most common approach for planners (J. Chen et al., 2016). Applications of computer vision jointly with these systems are

TABLE 2: COMPUTER VISION ALGORITHMS THAT TACKLE URBAN-RELATED ISSUES

CITY LAYER	CATEGORY	METHOD	
THE BUILT ENVIRONMENT	URBAN COMPONENTS	SEMANTIC SEGMENTATION	(Zhou et al., 2017) (Chaurasia & Culurciello, 2017) (Chen et al., 2016) (H. He, Yang, Wang, Wang, & Li, 2019) (Helbich et al., 2019) (Amirkolaei & Arefi, 2019) (Wurm, Stark, Zhu, Weigand, & Taubenböck, 2019) (Cordts et al., 2016)
		OBJECT-BASED DETECTION	(D. Yang, Liu, He, & Li, 2019) (R. Chew et al., 2018)
	LAND USE CLASSIFICATION	CLASSIFICATION AND SEMANTIC SEGMENTATION	(Demir et al., 2018) (Sharma, Liu, Yang, & Shi, 2017) (Audebert, Le Saux, & Lefèvre, 2018)
		CLASSIFICATION	(Wang, Xu, Dong, Gui, & Pu, 2018) (Srivastava, Vargas-Muñoz, & Tuia, 2019) (R. F. Chew et al., 2018)
	URBAN PERCEPTION	CLASSIFICATION AND PERCEPTION	(Ibrahim et al., 2019) (J. Zhao, Liu, Kuang, Chen, & Yang, 2018) (Law, Seresinhe, Shen, & Gutierrez-Roig, 2018) (F. Zhang, Wu, Zhu, & Liu, 2019) (Seresinhe, Preis, & Moat, 2017) (Oliva & Torralba, 2006) (W. Wang et al., 2018) (Salesses, Schechtner, & Hidalgo, 2013) (Dubey, Naik, Parikh, Raskar, & Hidalgo, 2016) (Naik, Raskar, & Hidalgo, 2016) (Quercia, O'Hare, & Cramer, 2014)
			URBAN SAFETY
	HUMAN INTERACTION		OBJECT-BASED DETECTION
TRANSPORTATION AND TRAFFIC	TRAFFIC SURVEILLANCE	CLASSIFICATION AND OBJECT-BASED DETECTION	(Bottino, Garbo, Loiacono, & Quer, 2016)
		ACTION RECOGNITION	(H. Yu, Wu, Wang, Wang, & Ma, 2017)
		OBJECT-BASED DETECTION	(Z. Yang & Pun-Cheng, 2018)
	SAFETY/ ACCIDENTS	CLASSIFICATION AND OBJECT-BASED DETECTION	(Sayed, Zaki, & Autey, 2013) (Zaki, Sayed, Tageldin, & Hussein, 2013)
THE NATURAL ENVIRONMENT	FLORA AND FAUNA	OBJECT-BASED DETECTION	(Cai, Li, Seiferling, & Ratti, 2018) (Hong, Han, Kim, Lee, & Kim, 2019)
		SEMANTIC SEGMENTATION	(Krause, Sugita, Baek, & Lim, 2018) (Williams et al., 2017)
		CLASSIFICATION	(Mohanty, Hughes, & Salathé, 2016) (Sun, Liu, Wang, & Zhang, 2017)
	ENVIRONMENTAL AND WEATHER CONDITIONS	CLASSIFICATION AND PERCEPTION	(C. Liu, Tsow, Zou, & Tao, 2016) (W. Liu, Yang, Wei, & School of Automation, China University of Geosciences, 2017) (Guerra, Khanam, Ehsan, Stolkin, & McDonald-Maier, 2018) (Elhoseiny, Huang, & Elgammal, 2015) (Sirirattanapol, Nagai, Witayangkurn, Pravinvongvuth, & Ekpanyapong, 2019)
INFRASTRUCTURE	CONCRETE CONDITION	OBJECT-BASED DETECTION	(Cha, Choi, & Büyükoztürk, 2017) (B. Wang, Zhao, Gao, Zhang, & Wang, 2018)
	PAVEMENT/ ROAD CONDITION	OBJECT-BASED DETECTION	(Maeda, Sekimoto, Seto, Kashiya, & Omata, 2018)
	BRIDGE COMPONENT RECOGNITION	SEMANTIC SEGMENTATION	(Narazaki, Hoskere, Hoang, & Jr, 2017)

capable of automating urban tasks such as mapping and zoning. Most recently, the notion of DeepGlobe (Demir et al., 2018) aimed to describe the earth from satellite images. DeepGlobe can extract streets, buildings and the different types of land-cover. Similarly, (Wang, Xu, Dong, Gui, & Pu, 2018) used a CNN model to segment satellite images into multi-classes at the pixel level. Marcos, Volpi, Kellenberger, & Tuia (2018) used the CNN model for land cover mapping, solving the issue of rotation of objects. Vanhoey et al. (2017) introduced VarCity as an approach of automating the construction of a city-scale 3D model based on semantic segmentation and machine processing of urban components (buildings, built environment, vegetation, roads, etc).

Furthermore, relying on deep learning, Amirkolaei & Arefi (2019) estimated heights from single aerial images, Wang et al. (2018) used deep CNN models for remote sensing image registration. Wurm, Stark, Zhu, Weigand, & Taubenböck (2019) relied on semantic segmentation to classify slum areas from aerial images.

These presented methods may differ from one another in terms of accuracies or purposes. However, the main limitation remains in how these models can be generalised to fit for multiple locations beyond the context where the models are trained and tested.

4.1.2 Seeing cities from a street-level

While it is vital to understand the overall urban systems of cities from an aerial view, seeing cities from the street-level adds more layers of information. These images can capture rapid urban changes in day-to-day life and offer more opportunities to model urban dynamics. However, capturing these rapid urban changes is a more complex task. Street-level images, taken by individuals or represented in Google's Street View API, have been used to identify a wide range of urban components from buildings to small objects such as street signs. For instance, Nguyen et al. (2018) used a CNN model to detect building types, crosswalks, and street greenness as a way to automatically quantify neighbourhood qualities.

Similarly, a range of applications based on classifying, segmenting and localising pixels from street-level images was a common approach for understanding the components of an urban scene (Chaurasia & Culurciello, 2017; Li, Jie, et al., 2017; Yang, Yu, Zhang, Li, & Yang, 2018; Zhou et al., 2017). Scene parsing relying on semantic segmentation is a continual success of CNN models for understanding and classifying the different components of the built environment at a pixel-level (Badrinarayanan et al., 2016; L.-C. Chen et al., 2016a, 2017; G. Lin et al., 2017; Long et al., 2015; Peng et al., 2017; F. Yu & Koltun, 2015; H. Zhao et al., 2017). Relying on both street-level images and satellite images, Kang, Körner, Wang, Taubenböck, & Zhu (2018) used a deep CNN model to classify land use in satellite images by learning from building blocks of similar functions.

Quantifying the physical and non-physical appearance of cities is another area that has been intensively researched. Naik et al. (2016) quantified the physical appearance of neighbourhoods based on individuals' ranking perceptions of the urban spaces using a framework of two CNN models that are concatenated and fused to predict a score for paired street-level images, known as Streetscore-CNN. Similarly, Zhang et al. (2018) quantified urban spaces of street-level images

labelled into six categories (Depressing, Boring, Beautiful, Safe, Lively, Wealthy) based on a crowdsourced dataset (MIT places pulse). By applying a supervised deep CNN model, they are able to predict the class for a given street view image. Liu, Silva, Wu, & Wang (2017) evaluated the urban visual appearance based on two indicators of the quality of street façade and the continuity of the street walls relying on the expert ranking that is evaluated with a public survey. Moreover, Naik, Kominers, Raskar, Glaeser, & Hidalgo (2017) have used computer vision to measure the dynamics of neighbourhood characteristics from time series street view images adjoined with socioeconomic data in five US cities. As a different approach, Law, Paige, & Russell (2018) used street view images to identify housing prices from urban perception relying on computer vision.

While seeing cities at street-level adds more information and gives an opportunity to understand the rapid changes that occur in an everyday urban scene in cities, the images used from Google street-view images only represent urban areas at a single weather condition, commonly clear weather, neglecting other visual and weather conditions that impact the appearance of cities. Furthermore, more research is needed on how to make best use of street level images coming from various sources, such as CCTV, dashcams or crowd sources, within and across domains.

4.2 Human interaction

Deep learning and computer vision have shown substantial progress in understanding a wide range of applications not only related to human detection but also understanding their activities and interaction with other objects (Kale & Patil, 2016; Mohamed & Ali, 2013; Zhang et al., 2017). Such approaches can assist planners and policy-makers to better understand tasks related to wellbeing and human behaviour in cities. For instance, Priya, Paul, & Singh (2015) used deep learning and computer vision to classify human actions, such as walking, running, sitting or dancing for multi-frame images. Guler, Neverova, & Kokkinos (2018) used a region-based CNN model (RCCN) to estimate the various human poses from a single image to better understand human interactions. Gkioxari, Girshick, Dollar, & He (2017) used computer vision to predict human actions over a specific target object from every day still images. This novel approach provides substantial progress in understanding human interaction with different objects. Furthermore, adjoining human pose detection with tracking, (Girdhar et al., 2017) used computer vision to detect and track key human body points from videos. This could enable, for example, tackling various issues related to human safety and wellbeing in cities such as detecting when a person falls, or detecting abnormal behaviour such as crime-related actions. Indeed, a knowledge gap appears in this field of study in scaling-up deep computer vision algorithms for monitoring and detecting irregular behaviours at a city level in real-time.

4.3 Transportation and traffic

Transportation and traffic is a crucial and complex layer that merges and interacts with other layers of the city. There is a wide range of computer vision applications that aim to tackle transport modes and their common issues, such as road safety

and optimisation of traffic (N. Buch, Velastin, & Orwell, 2011; Priya et al., 2015). Subjectively, traffic surveillance and intelligent transportation systems hold the largest share of computer vision related applications in cities. Typical tasks include vehicle detection, counting, overtake detection, and traffic incident detection (Mahmud, Ferreira, Hoque, & Tavassoli, 2017; Yang & Pun-Cheng, 2018). A full review of the literature on vehicle detection is beyond the scope of this article, for a comprehensive review consult Yang and Pun-Cheng (2018).

Understanding the different traffic scenarios and interactions of the different transport modes by computer vision is crucial. Bottino, Garbo, Loiacono, & Quer (2016) introduced 'Street Viewer' as a system to tackle and analyse the different scenarios of traffic behaviour from street view images. Sayed, Zaki, & Autey (2013) used computer vision to evaluate the safety measures of vehicle-bicycle conflicts. Zaki, Sayed, Tageldin, & Hussein (2013) used computer vision to analyse the conflicts among pedestrians and vehicles at a signalized intersection. Zaki & Sayed (2013) introduced a framework relying on computer vision to classify the different types of road-users.

Building on the aforementioned artificial intelligence approaches for traffic-related issues, computer vision is a core element when it comes to smart mobility and autonomous vehicles. Different applications relying on computer vision are being used to make transport modes aware of the surrounding environments either for safety indications or moving towards a self-navigation system. However, the technology of autonomous vehicles is not the focus of this research but rather the interactions of transport modes with the aforementioned layers in cities (Faisal, Yigitcanlar, Kamruzzaman, & Currie, 2019).

4.4 The natural environment

The natural environment (i.e. green space, landscape, climate conditions, etc.) is a crucial layer when it comes to understanding cities. It influences our perception of the visual appearance of the built environment and also affects mobility and human interaction in cities. Different aspects related to this natural layer of cities have been tackled by computer vision. These applications vary from mapping vegetation and greenery in cities, or so-called 'Treepedia' (Cai et al., 2018), identifying plant types (Krause et al., 2018; Sun et al., 2017), to deeper understanding of the natural environment and wildlife such as detecting plant-related diseases (Mohanty et al., 2016) and understanding the patterns of social interaction among animals (Robie, Seagraves, Egnor, & Branson, 2017).

Deep learning and computer vision have also been used to infer the weather, climatic and air conditions in cities. Liu et al. (2016) used the CNN model to identify extreme weather conditions from aerial images of climate simulations and reanalysis products. Liu, Tsow, Zou, & Tao, (2016) used images to analyse particle pollution for Beijing, Shanghai and Phoenix relying on region of interest selection, feature extraction and regression models. Z. Li et al. (2019) developed a model to detect clouds from high-resolution aerial view images relying on CNNs, named multi-scale convolutional feature fusion.

While there is noticeable progress in term of methods development and accuracy enhancement among the presented papers, the common limitation remains in the lack of a single model or a framework that fuses various models to infer the different weather and environmental conditions.

4.5 Infrastructure

Cities comprise a range of infrastructure systems that represent a large portion of their economy. Inspecting these systems and detecting their deficiencies is a crucial aspect for engineers and planners in cities. The focus of this section differs from the built environment section by analysing materials and the civil engineering related issues that are not covered in the aforementioned sections.

So far, the applications of computer vision have been seen in a wide range of domains related to infrastructure and civil engineering (Gopalakrishnan, 2018; Griffiths & Boehm, 2018), most importantly in analysing defects (Feng, Liu, Kao, & Lee, 2017). For instance, B. Wang, Zhao, Gao, Zhang, & Wang (2018) used computer vision to detect concrete crack damage. Similarly, Cha, Choi, & Büyüköztürk (2017) applied computer vision relying on deep CNN model to detect crack damage of concrete. On the other hand, Maeda, Sekimoto, Seto, Kashiya, & Omata (2018) used computer vision to detect road damage from images that are taken from mobile devices.

5. WHAT REMAINS MISSING?

Section 3 of this paper presented the different types of computer vision algorithms that are available to researchers, and the sectors in which they have been applied were presented in section 4. Typically, these models have been applied in a sectoral fashion to a specific problem. Comparatively little attention has been placed on how to understand the interconnections between the different layers of the city. These interconnections will eventually lead to increased capabilities of computer vision and AI to aid decision making and policy. In this section, we outline 2 under-researched areas in which computer vision has enormous potential.

5.1 Integrated models of the layers of the city

A significant challenge remains in modelling the interconnectedness and dependencies of the different layers of the city that were introduced in figure 2. A first step in this regard is the integration of models that have been developed for each layer in isolation. For example, there is still a knowledge gap in how to use computer vision coupled with deep learning to understand the interaction between people in cities and transport modes, or the influence of one mode on the others in terms of accessibility and safety. While the technology is there, the challenge remains in combining different models in a framework that enables them to tackle complex, multi-layered issues using the same data source, rather than just combining or fusing outputs from different data sources. On the other hand, even if the knowledge of the models is transferable among the different layers of cities, the challenges remain in finding comprehensive image data sources that cover a wide scope of tasks and functions in cities.

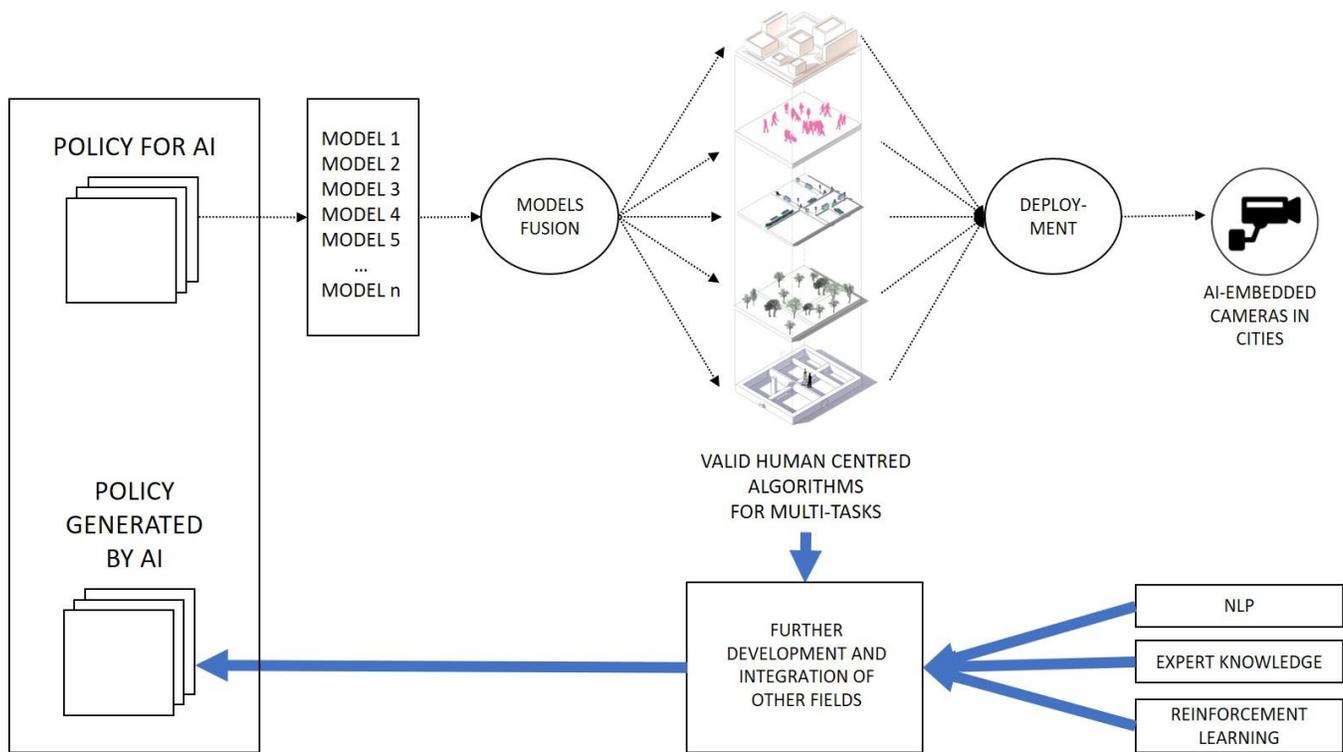


FIG. 3
AI-RELATED URBAN POLICY

5.2 The scale of applying computer vision in cities

Understanding cities requires both local and global perspectives, in which scale plays a crucial role in tackling urban issues. There are different algorithms that have been used to understand, for example, individuals' actions and activities. Challenges remain in applying and scaling up such algorithms to the city level. Although there are different models, as discussed in the literature, that extract information at the city scale, the nature of the developed algorithms is still limited towards the analysis of certain area or city. The reason for this is either because of a lack of computational resources or the inability of trained models to generalise to a larger dataset at a city-level. Models often require further training and optimisation to be deployed in real-life applications. It is well known that computer vision algorithms require large sets of labelled data, which must often be manually labelled. Labels can be crowd sourced but there is often a cost involved and accuracy is difficult to guarantee. Semi- or weakly supervised learning methods are promising approaches in this regard (S. Guo et al., 2018).

6. MOVING FROM PREDICTION TO DECISION-MAKING TO POLICY

After addressing the limitations of the stated models, the superior performance of modern computer vision algorithms is in little doubt. However, the extent to which model outputs can be used for automated and optimised policy and decision making remains an important research frontier. Big data, of which image data is a subset, is increasingly having an impact on decision and policy making, whether explicitly or not. Government authorities rely on algorithmic outputs to inform their decisions on a daily basis. The practical, ethical and

societal implications of this are still unclear and (Duarte & Álvarez, 2019) note the lack of synchronicity between the potential societal impact of AI technologies and our cultural discussions around them. An option that shows promise in this direction is the concept of living labs and policy labs. These provide testbeds within which to test data driven policies, which use ICT to realise the benefits of new data sources and support collaboration with relevant stakeholders and citizens (van Veenstra & Kotterink, 2017).

Alongside other sources of big data, images and video play a particularly important role in this effort because they capture the action and interaction of humans within their environment. This provides the opportunity to understand a range of issues, such as how the structure of the built environment affects pedestrian safety, or how street lighting influences crime. These issues are inextricably linked, and urban planning and policy making must take a holistic view of them to avoid disadvantaging certain groups.

6.1 Enabling Technologies

There are two enabling technologies that will be important in this area. Firstly, multi-agent reinforcement learning (MARL) will enable more realistic human agents to be simulated in more realistic urban environments. The behaviour of these agents can be learned and validated using image and video data. Such models could support or supersede traditional land use and transport planning approaches, as well as optimise the performance of urban systems such as transportation.

The second technology is GANs. It is not inconceivable that GANs, fed with images of a city, whether street view or aerial, could eventually be trained to design effective urban

environments. In the same way that GANs can generate synthetic human faces that are indistinguishable from real faces (Karras, Laine, & Aila, 2018), they could be used to plan new cities or neighbourhoods that perform like existing cities. This is certainly a long way off, but advancements in AI will enable predictions that are beyond what humans of social groups may achieve, or even conceive of (Duarte & Álvarez, 2019).

6.2 AI-embedded cameras in cities for real-time insights

The implementation of computer vision model pipelines in (near) real-time is a crucial issue for urban analytics and Internet of Things (IoT) systems. This deployment at the edge in urban contexts can show a direct impact of the current research for developing urban theories and policies. For example, AI-embedded cameras may alert police or transport control rooms of incidents, which they can verify and respond to. This type of system should be managed in a coordinated fashion so that the needs of various authorities can be met, which requires integration of the different layers of the city. However, while this approach will enable fast decision making and response, it falls short of being a fully intelligent and automated system able to implement or generate policy.

6.3 Policy for AI and by AI

After the deployment of AI in cities based on accepted norms and ethics, their deployment in cities will also lead to the generation of adaptive urban policies by AI. AI has the potential to generate dynamic and place-based policies. However, challenges remain in the innovation and fusion of different domains of knowledge to reach this critical step where the machine not only predicts and makes decisions but generates short- and long-term plans. Most importantly, it is a mixture of the tackled deep learning and computer vision research in urban settings with Natural Language Processing (NLP) research and reinforcement learning. By merging these different knowledge domains and integrating models that are capable of addressing multiple tasks in cities, theories and more flexible place-oriented policies can be generated for cities. Nevertheless, knowledge can be transferred from one city to another.

6.4 Conceptual framework towards AI generated policy and decision making

Fig. 3 shows a conceptual framework and a recommended process for achieving the two crucial steps outlined in sections 6.2 and 6.3, and how they can be reached from the current perspective of the deep computer vision research that is highlighted in this review. It shows the overall system for policy-makers and developers showing the important aspect of this process and the domains that are still under-developed and require further integration with urban analytics research.

Currently we are at the stage where policy for AI is being developed to mitigate the risks of reliance on the technology to make decisions. However, we envisage a future where the integration of the layers of the city through AI enables understanding of urban processes that is not possible by viewing them in isolation, leading to AI generated policies, as stated in section 6.3.

7. CONCLUSIONS

Understanding cities has been a profound interest for many scholars across a wide range of disciplines. Modelling the different urban systems of cities is a longevitous purpose for many urban and transport planners. While cities are complex by nature and classical urban modelling may not capture the actual complexities of urban systems, computer vision shows progress in tackling a variety of complex physical and non-physical visual tasks. In this article, we provide a review of deep learning and computer vision and its application so far in understanding cities. The article highlights the different types of algorithms of computer vision and their application to cities and their multifaceted issues. It aimed to show the nuances of the variations of these algorithms within the same task. It also aimed to show what has been done so far to understand cities by machine vision and what remains missing for future research work within this domain.

We attempt to highlight the potential role of computer vision in understanding the interactions between the built environment, people and transportation in order to tackle the complexity and nonlinearity of many urban and transport issues for better policy-making and planning safer cities. We also highlight the current limitations that require further work to reach an integrated computer vision-based urban models that capable of making automatic decisions.

REFERENCES

- [1] Amirkolaei, H. A., & Arefi, H. (2019). Height estimation from single aerial images using a deep convolutional encoder-decoder network. *ISPRS Journal of Photogrammetry and Remote Sensing*, 149, 50–66. <https://doi.org/10.1016/j.isprsjprs.2019.01.013>
- [2] Audebert, N., Le Saux, B., & Lefèvre, S. (2018). Beyond RGB: Very high resolution urban remote sensing with multimodal deep networks. *ISPRS Journal of Photogrammetry and Remote Sensing*, 140, 20–32. <https://doi.org/10.1016/j.isprsjprs.2017.11.011>
- [3] Badrinarayanan, V., Kendall, A., & Cipolla, R. (2016). *SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation*. arXiv:1511.00561v3 [cs.CV].
- [4] Batty, M. (2008). The Size, Scale, and Shape of Cities. *Science*, 319(5864), 769–771. <https://doi.org/10.1126/science.1151419>
- [5] Batty, M. (2009). Urban Modeling. In *International Encyclopedia of Human Geography* (pp. 51–58). Oxford, UK: Elsevier.
- [6] Batty, M., & Longley, P. (1994). *Fractal Cities: A Geometry of Form and Function*. New York: Academic Press.
- [7] Batty, M., & Xie, Y. (1996). Preliminary Evidence for a Theory of the Fractal City. *Environment and Planning A*, 28(10), 1745–1762. <https://doi.org/10.1068/a281745>
- [8] Batty, Michael. (2005). Agents, Cells, and Cities: New Representational Models for Simulating Multiscale Urban Dynamics. *Environment and Planning A*, 37(8), 1373–1394. <https://doi.org/10.1068/a3784>
- [9] Batty, Michael. (2019). Urban analytics defined. *Environment and Planning B: Urban Analytics and City Science*, 46(3), 403–405. <https://doi.org/10.1177/2399808319839494>
- [10] Becattini, F., Uricchio, T., Seidenari, L., Del Bimbo, A., & Ballan, L. (2017). Am I Done? Predicting Action Progress in Videos. *ArXiv:1705.01781 [Cs]*. Retrieved from <http://arxiv.org/abs/1705.01781>
- [11] Bettencourt, L. (2013). The origins of scaling in cities. *Science*, 340(6139), 1438–1441.
- [12] Bilen, H., Fernando, B., Gavves, E., Vedaldi, A., & Gould, S. (2016). Dynamic Image Networks for Action Recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3034–3042. <https://doi.org/10.1109/CVPR.2016.331>
- [13] Bottino, A., Garbo, A., Loiacono, C., & Quer, S. (2016). Street Viewer: An Autonomous Vision Based Traffic Tracking System. *Sensors*, 16(6), 813. <https://doi.org/10.3390/s16060813>

- [14] Brock, A., Donahue, J., & Simonyan, K. (2018). Large Scale GAN Training for High Fidelity Natural Image Synthesis. *ArXiv:1809.11096 [Cs, Stat]*. Retrieved from <http://arxiv.org/abs/1809.11096>
- [15] Buch, N., Velastin, S. A., & Orwell, J. (2011). A Review of Computer Vision Techniques for the Analysis of Urban Traffic. *IEEE Transactions on Intelligent Transportation Systems*, 12(3), 920–939. <https://doi.org/10.1109/TITS.2011.2119372>
- [16] Buch, S., Escorcia, V., Shen, C., Ghanem, B., & Niebles, J. C. (2017). SST: Single-Stream Temporal Action Proposals. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 6373–6382. <https://doi.org/10.1109/CVPR.2017.675>
- [17] Cai, B. Y., Li, X., Seiferling, I., & Ratti, C. (2018). Treepedia 2.0: Applying Deep Learning for Large-scale Quantification of Urban Tree Cover. *ArXiv:1808.04754 [Cs]*. Retrieved from <http://arxiv.org/abs/1808.04754>
- [18] Calder, M., Craig, C., Culley, D., de Cani, R., Donnelly, C. A., Douglas, R., ... Wilson, A. (2018). Computational modelling for decision-making: Where, why, what, who and how. *Royal Society Open Science*, 5(6), 172096. <https://doi.org/10.1098/rsos.172096>
- [19] Cao, Y., Wu, Z., & Shen, C. (2017). Estimating Depth from Monocular Images as Classification Using Deep Fully Convolutional Residual Networks. *ArXiv:1605.02305 [Cs]*. Retrieved from <http://arxiv.org/abs/1605.02305>
- [20] Cao, Z., Simon, T., Wei, S.-E., & Sheikh, Y. (2016). Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. *ArXiv:1611.08050 [Cs]*. Retrieved from <http://arxiv.org/abs/1611.08050>
- [21] Caron, M., Bojanowski, P., Joulin, A., & Douze, M. (2018). *Deep Clustering for Unsupervised Learning of Visual Features*. 29.
- [22] Cha, Y.-J., Choi, W., & Büyükoztürk, O. (2017). Deep Learning-Based Crack Damage Detection Using Convolutional Neural Networks: Deep learning-based crack damage detection using CNNs. *Computer-Aided Civil and Infrastructure Engineering*, 32(5), 361–378. <https://doi.org/10.1111/mice.12263>
- [23] Chao, Y.-W., Vijayanarasimhan, S., Seybold, B., Ross, D. A., Deng, J., & Sukthankar, R. (2018). Rethinking the Faster R-CNN Architecture for Temporal Action Localization. *ArXiv:1804.07667 [Cs]*. Retrieved from <http://arxiv.org/abs/1804.07667>
- [24] Chaurasia, A., & Culurciello, E. (2017). LinkNet: Exploiting Encoder Representations for Efficient Semantic Segmentation. *2017 IEEE Visual Communications and Image Processing (VCIP)*, 1–4. <https://doi.org/10.1109/VCIP.2017.8305148>
- [25] Chen, J., Dowman, I., Li, S., Li, Z., Madden, M., Mills, J., ... Heipke, C. (2016a). Information from imagery: ISPRS scientific vision and research agenda. *ISPRS Journal of Photogrammetry and Remote Sensing*, 115, 3–21. <https://doi.org/10.1016/j.isprsjprs.2015.09.008>
- [26] Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. (2016a). *SEMANTIC IMAGE SEGMENTATION WITH DEEP CONVOLUTIONAL NETS AND FULLY CONNECTED CRFS*.
- [27] Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2016b). DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. *ArXiv:1606.00915 [Cs]*. Retrieved from <http://arxiv.org/abs/1606.00915>
- [28] Chen, L.-C., Papandreou, G., Schroff, F., & Adam, H. (2017). Rethinking atrous convolution for semantic image segmentation. *ArXiv Preprint ArXiv:1706.05587*.
- [29] Chen, W., & Corso, J. J. (2015). Action Detection by Implicit Intentional Motion Clustering. *2015 IEEE International Conference on Computer Vision (ICCV)*, 3298–3306. <https://doi.org/10.1109/ICCV.2015.377>
- [30] Chen, Y., Yang, X., Zhong, B., Pan, S., Chen, D., & Zhang, H. (2016). CNNTracker: Online discriminative object tracking via deep convolutional neural network. *Applied Soft Computing*, 38, 1088–1098. <https://doi.org/10.1016/j.asoc.2015.06.048>
- [31] Chew, R. F., Amer, S., Jones, K., Unangst, J., Cajka, J., Allpress, J., & Bruhn, M. (2018). Residential scene classification for gridded population sampling in developing countries using deep convolutional neural networks on satellite imagery. *International Journal of Health Geographics*, 17(1). <https://doi.org/10.1186/s12942-018-0132-1>
- [32] Chew, R., Jones, K., Unangst, J., Cajka, J., Allpress, J., Amer, S., & Krotki, K. (2018). Toward Model-Generated Household Listing in Low- and Middle-Income Countries Using Deep Learning. *ISPRS International Journal of Geo-Information*, 7(11), 448. <https://doi.org/10.3390/ijgi7110448>
- [33] Chu, Q., Ouyang, W., Li, H., Wang, X., Liu, B., & Yu, N. (2017). Online Multi-Object Tracking Using CNN-based Single Object Tracker with Spatial-Temporal Attention Mechanism. *ArXiv:1708.02843 [Cs]*. Retrieved from <http://arxiv.org/abs/1708.02843>
- [34] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., ... Schiele, B. (2016). The cityscapes dataset for semantic urban scene understanding. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3213–3223.
- [35] Danelljan, M., Bhat, G., Khan, F. S., & Felsberg, M. (2016). ECO: Efficient Convolution Operators for Tracking. *ArXiv:1611.09224 [Cs]*. Retrieved from <http://arxiv.org/abs/1611.09224>
- [36] Danelljan, M., Hager, G., Khan, F. S., & Felsberg, M. (2015). Convolutional Features for Correlation Filter Based Visual Tracking. *2015 IEEE International Conference on Computer Vision Workshop (ICCVW)*, 621–629. <https://doi.org/10.1109/ICCVW.2015.84>
- [37] De Nadai, M., Vieriu, R. L., Zen, G., Dragicevic, S., Naik, N., Caraviello, M., ... Lepri, B. (2016). Are Safer Looking Neighborhoods More Lively?: A Multimodal Investigation into Urban Life. *Proceedings of the 2016 ACM on Multimedia Conference - MM '16*, 1127–1135. <https://doi.org/10.1145/2964284.2964312>
- [38] Demir, I., Koperski, K., Lindenbaum, D., Pang, G., Huang, J., Basu, S., ... Raskar, R. (2018). *DeepGlobe 2018: A Challenge to Parse the Earth through Satellite Images*. 10.
- [39] Diba, A., Fayyaz, M., Sharma, V., Karami, A. H., Arzani, M. M., Yousefzadeh, R., & Gool, L. V. (2017). *Temporal 3D ConvNets Using Temporal Transition Layer*. 5.
- [40] Duarte, F., & Álvarez, R. (2019). The data politics of the urban age. *Palgrave Communications*, 5(1), 1–7. <https://doi.org/10.1057/s41599-019-0264-3>
- [41] Dubey, A., Naik, N., Parikh, D., Raskar, R., & Hidalgo, C. A. (2016). Deep learning the city: Quantifying urban perception at a global scale. *European Conference on Computer Vision*, 196–212. Springer.
- [42] Elhoseiny, M., Huang, S., & Elgammal, A. (2015). Weather classification with deep convolutional neural networks. *2015 IEEE International Conference on Image Processing (ICIP)*, 3349–3353. <https://doi.org/10.1109/ICIP.2015.7351424>
- [43] El-Nouby, A., & Taylor, G. W. (2018). Real-Time End-to-End Action Detection with Two-Stream Networks. *ArXiv:1802.08362 [Cs]*. Retrieved from <http://arxiv.org/abs/1802.08362>
- [44] Escorcia, V., Caba Heilbron, F., Niebles, J. C., & Ghanem, B. (2016). DAPs: Deep Action Proposals for Action Understanding. In B. Leibe, J. Matas, N. Sebe, & M. Welling (Eds.), *Computer Vision – ECCV 2016* (Vol. 9907, pp. 768–784). https://doi.org/10.1007/978-3-319-46487-9_47
- [45] Eslami, S. M. A., Jimenez Rezende, D., Besse, F., Viola, F., Morcos, A. S., Garnelo, M., ... Hassabis, D. (2018). Neural scene representation and rendering. *Science*, 360(6394), 1204–1210. <https://doi.org/10.1126/science.aar6170>
- [46] Faisal, A., Yigitcanlar, T., Kamruzzaman, Md., & Currie, G. (2019). Understanding autonomous vehicles: A systematic literature review on capability, impact, planning and policy. *Journal of Transport and Land Use*, 12(1). <https://doi.org/10.5198/jtlu.2019.1405>
- [47] Fang, H.-S., Xie, S., Tai, Y.-W., & Lu, C. (2016). RMPE: Regional Multi-person Pose Estimation. *ArXiv:1612.00137 [Cs]*. Retrieved from <http://arxiv.org/abs/1612.00137>
- [48] Feng, C., Liu, M.-Y., Kao, C.-C., & Lee, T.-Y. (2017). Deep Active Learning for Civil Infrastructure Defect Detection and Classification. *Computing in Civil Engineering 2017*, 298–306. <https://doi.org/10.1061/9780784480823.036>
- [49] Frankhauser, P. (1998). The fractal approach. A new tool for the spatial analysis of urban agglomerations. *Population*, 10(1), 205–240.
- [50] Gemert, J. C. van, Jain, M., Gati, E., & Snoek, C. G. M. (2015). APT: Action localization proposals from dense trajectories. *Proceedings of the British Machine Vision Conference 2015*, 177.1-177.12. <https://doi.org/10.5244/C.29.177>
- [51] Girdhar, R., Gkioxari, G., Torresani, L., Paluri, M., & Tran, D. (2017). *Detect-and-Track: Efficient Pose Estimation in Videos*. 10.
- [52] Girdhar, R., & Ramanan, D. (2017). Attentional Pooling for Action Recognition. *ArXiv:1711.01467 [Cs]*. Retrieved from <http://arxiv.org/abs/1711.01467>
- [53] Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). *Rich feature hierarchies for accurate object detection and semantic segmentation*. Retrieved from <https://arxiv.org/pdf/1311.2524.pdf>
- [54] Gkioxari, G., Girshick, R., Dollar, P., & He, K. (2017). *Detecting and Recognizing Human-Object Interactions*. 9.
- [55] Goodfellow, I. (2016). NIPS 2016 Tutorial: Generative Adversarial Networks. *ArXiv:1701.00160 [Cs]*. Retrieved from <http://arxiv.org/abs/1701.00160>

- [56] Goodfellow, I., Bengio, Y., & Courville, A. (2017). *Deep Learning*. Cambridge, Massachusetts: The MIT Press.
- [57] Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... Bengio, Y. (2014). Generative Adversarial Networks. *ArXiv:1406.2661 [Cs, Stat]*. Retrieved from <http://arxiv.org/abs/1406.2661>
- [58] Gopalakrishnan, K. (2018). Deep Learning in Data-Driven Pavement Image Analysis and Automated Distress Detection: A Review. *Data*, 3(3), 28. <https://doi.org/10.3390/data3030028>
- [59] Griffiths, D., & Boehm, J. (2018). RAPID OBJECT DETECTION SYSTEMS, UTILISING DEEP LEARNING AND UNMANNED AERIAL SYSTEMS (UAS) FOR CIVIL ENGINEERING APPLICATIONS. *ISPRS - International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, XLII-2*, 391–398. <https://doi.org/10.5194/isprs-archives-XLII-2-391-2018>
- [60] Guerra, J. C. V., Khanam, Z., Ehsan, S., Stolkin, R., & McDonald-Maier, K. (2018). Weather Classification: A new multi-class dataset, data augmentation approach and comprehensive evaluations of Convolutional Neural Networks. *ArXiv:1808.00588 [Cs]*. Retrieved from <http://arxiv.org/abs/1808.00588>
- [61] Guler, R. A., Neverova, N., & Kokkinos, I. (2018). *DensePose: Dense Human Pose Estimation In The Wild*. 10.
- [62] Guo, M., Chou, E., Huang, D.-A., Song, S., Yeung, S., & Fei-Fei, L. (2018). *Neural Graph Matching Networks for Fewshot 3D Action Recognition*. 17.
- [63] Guo, S., Huang, W., Zhang, H., Zhuang, C., Dong, D., Scott, M. R., & Huang, D. (2018). CurriculumNet: Weakly Supervised Learning from Large-Scale Web Images. *ArXiv:1808.01097 [Cs]*. Retrieved from <http://arxiv.org/abs/1808.01097>
- [64] Guo, Y., Liu, Y., Oerlemans, A., Lao, S., Wu, S., & Lew, M. S. (2016). Deep learning for visual understanding: A review. *Neurocomputing*, 187, 27–48. <https://doi.org/10.1016/j.neucom.2015.09.116>
- [65] He, H., Yang, D., Wang, S., Wang, S., & Li, Y. (2019). Road Extraction by Using Atrous Spatial Pyramid Pooling Integrated Encoder-Decoder Network and Structural Similarity Loss. *Remote Sensing*, 11(9), 1015. <https://doi.org/10.3390/rs11091015>
- [66] He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep Residual Learning for Image Recognition. *ArXiv:1512.03385v1*. Retrieved from <https://arxiv.org/pdf/1512.03385.pdf>
- [67] He, L., Wang, G., & Hu, Z. (2018). Learning Depth from Single Images with Deep Neural Network Embedding Focal Length. *IEEE Transactions on Image Processing*, 27(9), 4676–4689. <https://doi.org/10.1109/TIP.2018.2832296>
- [68] Helbich, M., Yao, Y., Liu, Y., Zhang, J., Liu, P., & Wang, R. (2019). Using deep learning to examine street view green and blue spaces and their associations with geriatric depression in Beijing, China. *Environment International*, 126, 107–117. <https://doi.org/10.1016/j.envint.2019.02.013>
- [69] Held, D., Thrun, S., & Savarese, S. (2016). Learning to Track at 100 FPS with Deep Regression Networks. *ArXiv:1604.01802 [Cs]*. Retrieved from <http://arxiv.org/abs/1604.01802>
- [70] Heppenstall, A. J., Crooks, A. T., See, L. M., & Batty, M. (Eds.). (2012). *Agent-Based Models of Geographical Systems*. <https://doi.org/10.1007/978-90-481-8927-4>
- [71] Hester, T., Vecerik, M., Pietquin, O., Lanctot, M., Schaul, T., Piot, B., ... Gruslys, A. (2017). Deep Q-learning from Demonstrations. *ArXiv:1704.03732 [Cs]*. Retrieved from <http://arxiv.org/abs/1704.03732>
- [72] Hong, S.-J., Han, Y., Kim, S.-Y., Lee, A.-Y., & Kim, G. (2019). Application of Deep-Learning Methods to Bird Detection Using Unmanned Aerial Vehicle Imagery. *Sensors*, 19(7), 1651. <https://doi.org/10.3390/s19071651>
- [73] Hou, R., Chen, C., & Shah, M. (2017). Tube Convolutional Neural Network (T-CNN) for Action Detection in Videos. *2017 IEEE International Conference on Computer Vision (ICCV)*, 5823–5832. <https://doi.org/10.1109/ICCV.2017.620>
- [74] Huang, G., Liu, Z., Weinberger, K. Q., & van der Maaten, L. (2017). Densely connected convolutional networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1, 3.
- [75] Ibrahim, M. R., Haworth, J., & Cheng, T. (2019). URBAN-i: From urban scenes to mapping slums, transport modes, and pedestrians in cities using deep learning and computer vision. *Environment and Planning B: Urban Analytics and City Science*, 239980831984651. <https://doi.org/10.1177/2399808319846517>
- [76] Insafutdinov, E., Andriluka, M., Pishchulin, L., Tang, S., Levinkov, E., Andres, B., & Schiele, B. (2016). ArtTrack: Articulated Multi-person Tracking in the Wild. *ArXiv:1612.01465 [Cs]*. Retrieved from <http://arxiv.org/abs/1612.01465>
- [77] Insafutdinov, E., Pishchulin, L., Andres, B., Andriluka, M., & Schiele, B. (2016). DeeperCut: A Deeper, Stronger, and Faster Multi-Person Pose Estimation Model. *ArXiv:1605.03170 [Cs]*. Retrieved from <http://arxiv.org/abs/1605.03170>
- [78] Isalgué, A., Coch, H., & Serra, R. (2007). Scaling laws and the modern city. *Physica A: Statistical Mechanics and Its Applications*, 382(2), 643–649. <https://doi.org/10.1016/j.physa.2007.04.019>
- [79] Isola, P., Zhu, J.-Y., Zhou, T., & Efros, A. A. (2016). Image-to-Image Translation with Conditional Adversarial Networks. *ArXiv:1611.07004 [Cs]*. Retrieved from <http://arxiv.org/abs/1611.07004>
- [80] Jégou, S., Drozdal, M., Vazquez, D., Romero, A., & Bengio, Y. (2016). The One Hundred Layers Tiramisu: Fully Convolutional DenseNets for Semantic Segmentation. *ArXiv:1611.09326 [Cs]*. Retrieved from <http://arxiv.org/abs/1611.09326>
- [81] Jiang, C., Xiao, J., Xie, Y., Tillo, T., & Huang, K. (2018). Siamese network ensemble for visual tracking. *Neurocomputing*, 275, 2892–2903. <https://doi.org/10.1016/j.neucom.2017.10.043>
- [82] Kale, G. V., & Patil, V. H. (2016). A Study of Vision based Human Motion Recognition and Analysis. *International Journal of Ambient Computing and Intelligence*, 7(2), 18.
- [83] Kang, J., Körner, M., Wang, Y., Taubenböck, H., & Zhu, X. X. (2018). Building instance classification using street view images. *ISPRS Journal of Photogrammetry and Remote Sensing*, 145, 44–59. <https://doi.org/10.1016/j.isprsjprs.2018.02.006>
- [84] Kang, K., Ouyang, W., Li, H., & Wang, X. (2016). Object Detection from Video Tubelets with Convolutional Neural Networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 817–825. <https://doi.org/10.1109/CVPR.2016.95>
- [85] Karras, T., Laine, S., & Aila, T. (2018). A Style-Based Generator Architecture for Generative Adversarial Networks. *ArXiv:1812.04948 [Cs, Stat]*. Retrieved from <http://arxiv.org/abs/1812.04948>
- [86] Kocabas, M., Karagoz, S., & Akbas, E. (2018). MultiPoseNet: Fast Multi-Person Pose Estimation using Pose Residual Network. *ArXiv:1807.04067 [Cs]*. Retrieved from <http://arxiv.org/abs/1807.04067>
- [87] Krause, J., Sugita, G., Baek, K., & Lim, L. (2018). *WTPlant (What's That Plant?): A Deep Learning System for Identifying Plants in Natural Images. Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval - ICMR '18*, 517–520. <https://doi.org/10.1145/3206025.3206089>
- [88] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Proceeding NIPS'12 Proceedings of the 25th International Conference on Neural Information Processing Systems, 1*, 1097–1105. Lake Tahoe, Nevada: Curran Associates Inc., USA ©2012.
- [89] Kuo, C.-C. J. (2016). Understanding convolutional neural networks with a mathematical model. *Journal of Visual Communication and Image Representation*, 41, 406–413.
- [90] Law, S., Paige, B., & Russell, C. (2018). Take a Look Around: Using Street View and Satellite Images to Estimate House Prices. *ArXiv:1807.07155 [Cs, Econ]*. Retrieved from <http://arxiv.org/abs/1807.07155>
- [91] Law, S., Seresinhe, C. I., Shen, Y., & Gutierrez-Roig, M. (2018). Street-Frontage-Net: Urban image classification using deep convolutional neural networks. *International Journal of Geographical Information Science*, 1–27. <https://doi.org/10.1080/13658816.2018.1555832>
- [92] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- [93] Li, P., Wang, D., Wang, L., & Lu, H. (2018). Deep visual tracking: Review and experimental comparison. *Pattern Recognition*, 76, 323–338. <https://doi.org/10.1016/j.patcog.2017.11.007>
- [94] Li, X., Jie, Z., Wang, W., Liu, C., Yang, J., Shen, X., ... Feng, J. (2017). FoveaNet: Perspective-aware Urban Scene Parsing. *ArXiv:1708.02421 [Cs]*. Retrieved from <http://arxiv.org/abs/1708.02421>
- [95] Li, X., Wang, Z. J. W., Yang, C. L. J., Chen, X. S. Z. L. Q., Yan, S., & Feng, J. (2017). FoveaNet: Perspective-aware Urban Scene Parsing. *ArXiv Preprint ArXiv:1708.02421*.
- [96] Li, Y., Lan, C., Xing, J., Zeng, W., Yuan, C., & Liu, J. (2016). Online Human Action Detection using Joint Classification-Regression Recurrent Neural Networks. *ArXiv:1604.05633 [Cs]*. Retrieved from <http://arxiv.org/abs/1604.05633>
- [97] Li, Z., Shen, H., Cheng, Q., Liu, Y., You, S., & He, Z. (2019). Deep learning based cloud detection for medium and high resolution remote sensing images of different sensors. *ISPRS Journal of Photogrammetry*

- and Remote Sensing, 150, 197–212. <https://doi.org/10.1016/j.isprsjprs.2019.02.017>
- [98] Lin, G., Milan, A., Shen, C., & Reid, I. (2017). Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [99] Lin, T.-Y., Goyal, P., Girshick, R., He, K., & Dollár, P. (2018). Focal Loss for Dense Object Detection. *ArXiv:1708.02002 [Cs]*. Retrieved from <http://arxiv.org/abs/1708.02002>
- [100] Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., ... Dollár, P. (2014). Microsoft COCO: Common Objects in Context. *ArXiv:1405.0312 [Cs]*. Retrieved from <http://arxiv.org/abs/1405.0312>
- [101] Liu, C., Tsow, F., Zou, Y., & Tao, N. (2016). Particle Pollution Estimation Based on Image Analysis. *PLOS ONE*, 11(2), e0145955. <https://doi.org/10.1371/journal.pone.0145955>
- [102] 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. *Computers, Environment and Urban Systems*, 65, 113–125. <https://doi.org/10.1016/j.compenvurbsys.2017.06.003>
- [103] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., & Berg, A. C. (2016). Ssd: Single shot multibox detector. *European Conference on Computer Vision*, 21–37. Springer.
- [104] Liu, W., Yang, Y., Wei, L., & School of Automation, China University of Geosciences. (2017). Weather Recognition of Street Scene Based on Sparse Deep Neural Networks. *Journal of Advanced Computational Intelligence and Intelligent Informatics*, 21(3), 403–408. <https://doi.org/10.20965/jaci.2017.p0403>
- [105] Liu, Y., Racah, E., Prabhat, Correa, J., Khosrowshahi, A., Lavers, D., ... Collins, W. (2016). Application of Deep Convolutional Neural Networks for Detecting Extreme Weather in Climate Datasets. *ArXiv:1605.01156 [Cs]*. Retrieved from <http://arxiv.org/abs/1605.01156>
- [106] Long, J., Shelhamer, E., & Darrell, T. (2015). *Fully Convolutional Networks for Semantic Segmentation*. arXiv:1411.4038v2 [cs.CV].
- [107] Maeda, H., Sekimoto, Y., Seto, T., Kashiyama, T., & Omata, H. (2018). Road Damage Detection Using Deep Neural Networks with Images Captured Through a Smartphone. *ArXiv:1801.09454 [Cs]*. Retrieved from <http://arxiv.org/abs/1801.09454>
- [108] Mahmud, S. M. S., Ferreira, L., Hoque, Md. S., & Tavassoli, A. (2017). Application of proximal surrogate indicators for safety evaluation: A review of recent developments and research needs. *IATSS Research*, 41(4), 153–163. <https://doi.org/10.1016/j.iatssr.2017.02.001>
- [109] Manen, S., Gygli, M., Dai, D., & Gool, L. V. (2017). PathTrack: Fast Trajectory Annotation with Path Supervision. *2017 IEEE International Conference on Computer Vision (ICCV)*, 290–299. <https://doi.org/10.1109/ICCV.2017.40>
- [110] Marcos, D., Volpi, M., Kellenberger, B., & Tuia, D. (2018). Land cover mapping at very high resolution with rotation equivariant CNNs: Towards small yet accurate models. *ISPRS Journal of Photogrammetry and Remote Sensing*, 145, 96–107. <https://doi.org/10.1016/j.isprsjprs.2018.01.021>
- [111] Mettes, P., van Gemert, J. C., & Snoek, C. G. M. (2016). Spot On: Action Localization from Pointly-Supervised Proposals. In B. Leibe, J. Matas, N. Sebe, & M. Welling (Eds.), *Computer Vision – ECCV 2016* (Vol. 9909, pp. 437–453). https://doi.org/10.1007/978-3-319-46454-1_27
- [112] Mirowski, P., Grimes, M. K., Malinowski, M., Hermann, K. M., Anderson, K., Teplyashin, D., ... Hadsell, R. (2018). Learning to Navigate in Cities Without a Map. *ArXiv:1804.00168 [Cs]*. Retrieved from <http://arxiv.org/abs/1804.00168>
- [113] Mnih, V., Badia, A. P., Mirza, M., Graves, A., Harley, T., Lillicrap, T. P., ... Kavukcuoglu, K. (2016). *Asynchronous Methods for Deep Reinforcement Learning*. 10.
- [114] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M. (2013). Playing Atari with Deep Reinforcement Learning. *ArXiv:1312.5602 [Cs]*. Retrieved from <http://arxiv.org/abs/1312.5602>
- [115] Mohamed, A. N., & Ali, M. M. (2013). HUMAN MOTION ANALYSIS, RECOGNITION AND UNDERSTANDING IN COMPUTER VISION: A REVIEW. *Journal of Engineering Sciences*, 41(5), 19.
- [116] Mohanty, S. P., Hughes, D. P., & Salathé, M. (2016). Using Deep Learning for Image-Based Plant Disease Detection. *Frontiers in Plant Science*, 7. <https://doi.org/10.3389/fpls.2016.01419>
- [117] Murcio, R., Masucci, A. P., Arcaute, E., & Batty, M. (2015). Multifractal to monofractal evolution of the London street network. *Physical Review E*, 92(6). <https://doi.org/10.1103/PhysRevE.92.062130>
- [118] Naik, N., Kominers, S. D., Raskar, R., Glaeser, E. L., & Hidalgo, C. A. (2017). Computer vision uncovers predictors of physical urban change. *Proceedings of the National Academy of Sciences*, 114(29), 7571–7576. <https://doi.org/10.1073/pnas.1619003114>
- [119] Naik, N., Philipoom, J., Raskar, R., & Hidalgo, C. (2014). Streetscore-predicting the perceived safety of one million streetscapes. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 779–785.
- [120] Naik, N., Raskar, R., & Hidalgo, C. A. (2016). Cities Are Physical Too: Using Computer Vision to Measure the Quality and Impact of Urban Appearance. *American Economic Review*, 106(5), 128–132. <https://doi.org/10.1257/aer.p20161030>
- [121] Narazaki, Y., Hoskere, V., Hoang, T. A., & Jr, B. F. S. (2017). *Vision-based Automated Bridge Component Recognition Integrated With High-level Scene Understanding*. 10.
- [122] Nguyen, Q. C., Sajjadi, M., McCullough, M., Pham, M., Nguyen, T. T., Yu, W., ... Tasdizen, T. (2018). Neighbourhood looking glass: 360° automated characterisation of the built environment for neighbourhood effects research. *Journal of Epidemiology and Community Health*, 72(3), 260–266. <https://doi.org/10.1136/jech-2017-209456>
- [123] Oliva, A., & Torralba, A. (2006). Chapter 2 Building the gist of a scene: The role of global image features in recognition. In *Progress in Brain Research* (Vol. 155, pp. 23–36). [https://doi.org/10.1016/S0079-6123\(06\)55002-2](https://doi.org/10.1016/S0079-6123(06)55002-2)
- [124] Paganini, M., de Oliveira, L., & Nachman, B. (2018). CaloGAN: Simulating 3D high energy particle showers in multilayer electromagnetic calorimeters with generative adversarial networks. *Physical Review D*, 97(1). <https://doi.org/10.1103/PhysRevD.97.014021>
- [125] Papandreou, G., Zhu, T., Kanazawa, N., Toshev, A., Tompson, J., Bregler, C., & Murphy, K. (2017). Towards Accurate Multi-person Pose Estimation in the Wild. *ArXiv:1701.01779 [Cs]*. Retrieved from <http://arxiv.org/abs/1701.01779>
- [126] Peng, C., Zhang, X., Yu, G., Luo, G., & Sun, J. (2017). Large Kernel Matters—Improve Semantic Segmentation by Global Convolutional Network. *ArXiv Preprint ArXiv:1703.02719*.
- [127] Pfister, T., Charles, J., & Zisserman, A. (2015). Flowing ConvNets for Human Pose Estimation in Videos. *ArXiv:1506.02897 [Cs]*. Retrieved from <http://arxiv.org/abs/1506.02897>
- [128] Priya, G., Paul, S. N., & Singh, Y. J. (2015). *Human walking motion detection and classification of actions from Video Sequences*. 3(1), 6.
- [129] Quercia, D., O’Hare, N. K., & Cramer, H. (2014). Aesthetic capital: What makes london look beautiful, quiet, and happy? *Proceedings of the 17th ACM Conference on Computer Supported Cooperative Work & Social Computing - CSCW ’14*, 945–955. <https://doi.org/10.1145/2531602.2531613>
- [130] Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. *ArXiv:1511.06434 [Cs]*. Retrieved from <http://arxiv.org/abs/1511.06434>
- [131] Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. *ArXiv:1506.02640 [Cs]*. Retrieved from <http://arxiv.org/abs/1506.02640>
- [132] Redmon, J., & Farhadi, A. (2017). YOLO9000: Better, Faster, Stronger. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 6517–6525. <https://doi.org/10.1109/CVPR.2017.690>
- [133] Redmon, J., & Farhadi, A. (2018). *YOLOv3: An Incremental Improvement*. 6.
- [134] Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., & Lee, H. (2016). Generative Adversarial Text to Image Synthesis. *ArXiv:1605.05396 [Cs]*. Retrieved from <http://arxiv.org/abs/1605.05396>
- [135] Reed, S. E., Akata, Z., Mohan, S., Tenka, S., Schiele, B., & Lee, H. (2016). *Learning What and Where to Draw*. 9.
- [136] Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., & Prabhat. (2019). Deep learning and process understanding for data-driven Earth system science. *Nature*, 566(7743), 195–204. <https://doi.org/10.1038/s41586-019-0912-1>
- [137] Ren, S., He, K., Girshick, R., & Sun, J. (2016). *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks*. arXiv:1506.01497v3.
- [138] Robie, A. A., Seagraves, K. M., Egnor, S. E. R., & Branson, K. (2017). Machine vision methods for analyzing social interactions. *The Journal of Experimental Biology*, 220(1), 25–34. <https://doi.org/10.1242/jeb.142281>
- [139] Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image Segmentation. *ArXiv:1505.04597 [Cs]*. Retrieved from <http://arxiv.org/abs/1505.04597>

- [140] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3), 211–252. <https://doi.org/10.1007/s11263-015-0816-y>
- [141] Saha, S., Singh, G., & Cuzzolin, F. (2017). AMTnet: Action-Micro-Tube Regression by End-to-end Trainable Deep Architecture. *ArXiv:1704.04952 [Cs]*. Retrieved from <http://arxiv.org/abs/1704.04952>
- [142] Saha, S., Singh, G., Sapienza, M., Torr, P. H. S., & Cuzzolin, F. (2016). Deep Learning for Detecting Multiple Space-Time Action Tubes in Videos. *ArXiv:1608.01529 [Cs]*. Retrieved from <http://arxiv.org/abs/1608.01529>
- [143] Salesses, P., Schechtner, K., & Hidalgo, C. A. (2013). The Collaborative Image of The City: Mapping the Inequality of Urban Perception. *PLoS ONE*, 8(7), e68400. <https://doi.org/10.1371/journal.pone.0068400>
- [144] Sayed, T., Zaki, M. H., & Autey, J. (2013). Automated safety diagnosis of vehicle–bicycle interactions using computer vision analysis. *Safety Science*, 59, 163–172. <https://doi.org/10.1016/j.ssci.2013.05.009>
- [145] Seresinhe, C. I., Preis, T., & Moat, H. S. (2017). Using deep learning to quantify the beauty of outdoor places. *Royal Society Open Science*, 4(7), 170170. <https://doi.org/10.1098/rsos.170170>
- [146] Sharma, A., Liu, X., Yang, X., & Shi, D. (2017). A patch-based convolutional neural network for remote sensing image classification. *Neural Networks*, 95, 19–28. <https://doi.org/10.1016/j.neunet.2017.07.017>
- [147] Shou, Z., Chan, J., Zareian, A., Miyazawa, K., & Chang, S.-F. (2017). CDC: Convolutional-De-Convolutional Networks for Precise Temporal Action Localization in Untrimmed Videos. *ArXiv:1703.01515 [Cs]*. Retrieved from <http://arxiv.org/abs/1703.01515>
- [148] Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *ArXiv Preprint ArXiv:1409.1556*.
- [149] Singh, G., Saha, S., Sapienza, M., Torr, P., & Cuzzolin, F. (2016). Online Real-time Multiple Spatiotemporal Action Localisation and Prediction. *ArXiv:1611.08563 [Cs]*. Retrieved from <http://arxiv.org/abs/1611.08563>
- [150] Sirirattapanol, C., Nagai, M., Witayangkurn, A., Pravinovongvuth, S., & Ekpanyapong, M. (2019). Bangkok CCTV Image through a Road Environment Extraction System Using Multi-Label Convolutional Neural Network Classification. *ISPRS International Journal of Geo-Information*, 8(3), 128. <https://doi.org/10.3390/ijgi8030128>
- [151] Soomro, K., & Shah, M. (2017). Unsupervised Action Discovery and Localization in Videos. *2017 IEEE International Conference on Computer Vision (ICCV)*, 696–705. <https://doi.org/10.1109/ICCV.2017.82>
- [152] Srivastava, S., Vargas-Muñoz, J. E., & Tuia, D. (2019). Understanding urban landuse from the above and ground perspectives: A deep learning, multimodal solution. *Remote Sensing of Environment*, 228, 129–143. <https://doi.org/10.1016/j.rse.2019.04.014>
- [153] Sun, Y., Liu, Y., Wang, G., & Zhang, H. (2017). Deep Learning for Plant Identification in Natural Environment. *Computational Intelligence and Neuroscience*, 2017, 1–6. <https://doi.org/10.1155/2017/7361042>
- [154] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., & Reed, S. (2015). *Going Deeper with Convolutions*. Retrieved from <https://www.cs.unc.edu/~wliu/papers/GoogLeNet.pdf>
- [155] Tian, K., Zhou, S., & Guan, J. (2017). DeepCluster: A General Clustering Framework Based on Deep Learning. In M. Ceci, J. Hollmén, L. Todorovski, C. Vens, & S. Džeroski (Eds.), *Machine Learning and Knowledge Discovery in Databases* (Vol. 10535, pp. 809–825). https://doi.org/10.1007/978-3-319-71246-8_49
- [156] van Hasselt, H., Guez, A., & Silver, D. (2015). *Deep Reinforcement Learning with Double Q-Learning*. 7.
- [157] van Veenstra, A. F., & Kotterink, B. (2017). Data-Driven Policy Making: The Policy Lab Approach. In P. Parycek, Y. Charalabidis, A. V. Chugunov, P. Panagiotopoulos, T. A. Pardo, Ø. Sæbø, & E. Tambouris (Eds.), *Electronic Participation* (pp. 100–111). Springer International Publishing.
- [158] Vanhoey, K., Dai, D., Van Gool, L., de Oliveira, C. E. P., Riemenschneider, H., Bódis-Szomoró, A., ... Kroeger, T. (2017). VarCity - the video: The struggles and triumphs of leveraging fundamental research results in a graphics video production. *ACM SIGGRAPH 2017 Talks on - SIGGRAPH '17*, 1–2. <https://doi.org/10.1145/3084363.3085085>
- [159] Viola, P., & Jones, M. (2001). Rapid object detection using a boosted cascade of simple features. *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference On, 1*, I–I. IEEE.
- [160] Voigtlaender, P., Krause, M., Sekar, B. B. G., Geiger, A., & Leibe, B. (2019). *MOTS: Multi-Object Tracking and Segmentation*. 10.
- [161] Wang, B., Zhao, W., Gao, P., Zhang, Y., & Wang, Z. (2018). Crack Damage Detection Method via Multiple Visual Features and Efficient Multi-Task Learning Model. *Sensors*, 18(6), 1796. <https://doi.org/10.3390/s18061796>
- [162] Wang, Lei, Xu, X., Dong, H., Gui, R., & Pu, F. (2018). Multi-Pixel Simultaneous Classification of PolSAR Image Using Convolutional Neural Networks. *Sensors*, 18(3), 769. <https://doi.org/10.3390/s18030769>
- [163] Wang, Limin, Qiao, Y., & Tang, X. (2015). Action recognition with trajectory-pooled deep-convolutional descriptors. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4305–4314. <https://doi.org/10.1109/CVPR.2015.7299059>
- [164] Wang, S., Quan, D., Liang, X., Ning, M., Guo, Y., & Jiao, L. (2018). A deep learning framework for remote sensing image registration. *ISPRS Journal of Photogrammetry and Remote Sensing*. <https://doi.org/10.1016/j.isprsjprs.2017.12.012>
- [165] Wang, W., Yang, S., He, Z., Wang, M., Zhang, J., & Zhang, W. (2018). Urban Perception of Commercial Activeness from Satellite Images and Streetscapes. *Companion of the The Web Conference 2018 on The Web Conference 2018 - WWW '18*, 647–654. <https://doi.org/10.1145/3184558.3186581>
- [166] Wang, Z., Schaul, T., Hessel, M., van Hasselt, H., Lanctot, M., & de Freitas, N. (2015). Dueling Network Architectures for Deep Reinforcement Learning. *ArXiv:1511.06581 [Cs]*. Retrieved from <http://arxiv.org/abs/1511.06581>
- [167] Weinzaepfel, P., Harchaoui, Z., & Schmid, C. (2015). Learning to Track for Spatio-Temporal Action Localization. *2015 IEEE International Conference on Computer Vision (ICCV)*, 3164–3172. <https://doi.org/10.1109/ICCV.2015.362>
- [168] Weinzaepfel, P., Martin, X., & Schmid, C. (2016). Human Action Localization with Sparse Spatial Supervision. *ArXiv:1605.05197 [Cs]*. Retrieved from <http://arxiv.org/abs/1605.05197>
- [169] Williams, D., Britten, A., McCallum, S., Jones, H., Aitkenhead, M., Karley, A., ... Graham, J. (2017). A method for automatic segmentation and splitting of hyperspectral images of raspberry plants collected in field conditions. *Plant Methods*, 13(1). <https://doi.org/10.1186/s13007-017-0226-y>
- [170] Wu, G., Lu, W., Gao, G., Zhao, C., & Liu, J. (2016). Regional deep learning model for visual tracking. *Neurocomputing*, 175, 310–323. <https://doi.org/10.1016/j.neucom.2015.10.064>
- [171] Wurm, M., Stark, T., Zhu, X. X., Weigand, M., & Taubenböck, H. (2019). Semantic segmentation of slums in satellite images using transfer learning on fully convolutional neural networks. *ISPRS Journal of Photogrammetry and Remote Sensing*, 150, 59–69. <https://doi.org/10.1016/j.isprsjprs.2019.02.006>
- [172] Xie, J., Girshick, R., & Farhadi, A. (2016). *Unsupervised Deep Embedding for Clustering Analysis*. 10.
- [173] Xu, H., Das, A., & Saenko, K. (2017). R-C3D: Region Convolutional 3D Network for Temporal Activity Detection. *ArXiv:1703.07814 [Cs]*. Retrieved from <http://arxiv.org/abs/1703.07814>
- [174] Yang, D., Liu, X., He, H., & Li, Y. (2019). Air-to-ground multimodal object detection algorithm based on feature association learning. *International Journal of Advanced Robotic Systems*, 16(3), 172988141984299. <https://doi.org/10.1177/1729881419842995>
- [175] Yang, M., Yu, K., Zhang, C., Li, Z., & Yang, K. (2018). *DenseASPP for Semantic Segmentation in Street Scenes*. 9.
- [176] Yang, Z., & Pun-Cheng, L. S. C. (2018). Vehicle detection in intelligent transportation systems and its applications under varying environments: A review. *Image and Vision Computing*, 69, 143–154. <https://doi.org/10.1016/j.imavis.2017.09.008>
- [177] Yu, F., & Koltun, V. (2015). Multi-scale context aggregation by dilated convolutions. *ArXiv Preprint ArXiv:1511.07122*.
- [178] Yu, H., Wu, Z., Wang, S., Wang, Y., & Ma, X. (2017). Spatiotemporal Recurrent Convolutional Networks for Traffic Prediction in Transportation Networks. *Sensors*, 17(12), 1501. <https://doi.org/10.3390/s17071501>
- [179] Zaki, M. H., & Sayed, T. (2013). A framework for automated road-users classification using movement trajectories. *Transportation Research Part C: Emerging Technologies*, 33, 50–73. <https://doi.org/10.1016/j.trc.2013.04.007>
- [180] Zaki, M. H., Sayed, T., Tageldin, A., & Hussein, M. (2013). Application of Computer Vision to Diagnosis of Pedestrian Safety Issues. *Transportation Research Record: Journal of the Transportation Research Board*, 2393(1), 75–84. <https://doi.org/10.3141/2393-09>

- [181]Zhang, B., Wang, L., Wang, Z., Qiao, Y., & Wang, H. (2016). Real-time Action Recognition with Enhanced Motion Vector CNNs. *ArXiv:1604.07669 [Cs]*. Retrieved from <http://arxiv.org/abs/1604.07669>
- [182]Zhang, F., Wu, L., Zhu, D., & Liu, Y. (2019). Social sensing from street-level imagery: A case study in learning spatio-temporal urban mobility patterns. *ISPRS Journal of Photogrammetry and Remote Sensing*, 153, 48–58. <https://doi.org/10.1016/j.isprsjprs.2019.04.017>
- [183]Zhang, F., Zhou, B., Liu, L., Liu, Y., Fung, H. H., Lin, H., & Ratti, C. (2018). Measuring human perceptions of a large-scale urban region using machine learning. *Landscape and Urban Planning*, 180, 148–160. <https://doi.org/10.1016/j.landurbplan.2018.08.020>
- [184]Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. (2016). StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks. *ArXiv:1612.03242 [Cs, Stat]*. Retrieved from <http://arxiv.org/abs/1612.03242>
- [185]Zhang, S., Wei, Z., Nie, J., Huang, L., Wang, S., & Li, Z. (2017). A Review on Human Activity Recognition Using Vision-Based Method. *Journal of Healthcare Engineering*, 2017, 1–31. <https://doi.org/10.1155/2017/3090343>
- [186]Zhang, X., Xia, G.-S., Lu, Q., Shen, W., & Zhang, L. (2018). Visual object tracking by correlation filters and online learning. *ISPRS Journal of Photogrammetry and Remote Sensing*, 140, 77–89. <https://doi.org/10.1016/j.isprsjprs.2017.07.009>
- [187]Zhao, H., Shi, J., Qi, X., Wang, X., & Jia, J. (2017). Pyramid scene parsing network. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2881–2890.
- [188]Zhao, J., Liu, X., Kuang, Y., Chen, Y. V., & Yang, B. (2018). Deep CNN-Based Methods to Evaluate Neighborhood-Scale Urban Valuation Through Street Scenes Perception. *2018 IEEE Third International Conference on Data Science in Cyberspace (DSC)*, 20–27. <https://doi.org/10.1109/DSC.2018.00012>
- [189]Zhao, Y., Xiong, Y., Wang, L., Wu, Z., Tang, X., & Lin, D. (2017). Temporal Action Detection with Structured Segment Networks. *ArXiv:1704.06228 [Cs]*. Retrieved from <http://arxiv.org/abs/1704.06228>
- [190]Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., & Torralba, A. (2017). Scene Parsing through ADE20K Dataset. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5122–5130. <https://doi.org/10.1109/CVPR.2017.544>
- [191]Zhu, H., Vial, R., & Lu, S. (2017). TORNADO: A Spatio-Temporal Convolutional Regression Network for Video Action Proposal. *2017 IEEE International Conference on Computer Vision (ICCV)*, 5814–5822. <https://doi.org/10.1109/ICCV.2017.619>
- [192]Zhu, Y., Lan, Z., Newsam, S., & Hauptmann, A. G. (2017). Hidden Two-Stream Convolutional Networks for Action Recognition. *ArXiv:1704.00389 [Cs]*. Retrieved from <http://arxiv.org/abs/1704.00389>
- [193]Zou, Z., Shi, Z., Guo, Y., & Ye, J. (2019). Object Detection in 20 Years: A Survey. *ArXiv:1905.05055v2*, 40.