# On Machine Learning and Structure for Mobile Robots

**Markus Wulfmeier**
University of Oxford
markus@robots.ox.ac.uk

## Abstract

Due to recent advances - compute, data, models - the role of learning in autonomous systems has expanded significantly, rendering new applications possible for the first time. While some of the most significant benefits are obtained in the perception modules of the software stack, other aspects continue to rely on known manual procedures based on prior knowledge on geometry, dynamics, kinematics etc. Nonetheless, learning gains relevance in these modules when data collection and curation become easier than manual rule design. Building on this coarse and broad survey of current research, the final sections aim to provide insights into future potentials and challenges as well as the necessity of structure in current practical applications.

## Contents

# 1 Introduction

This survey[1] provides an informal overview of current challenges and potentials of learning across various tasks of relevance in robotics and automation. In this context, similar to the long-term discussion on how much innate structure is optimal for artificial intelligence, there is the more short-term question of how to merge traditional programming and learning (e.g. described as *differentiable programming* or *software 2.0*) for more narrow applications in efficient, robust and safe automation. The question about structure as beneficial or limiting aspect becomes arguably easier to answer in the context of robotic near-term applications as we can simply *acknowledge our ignorance* (the missing knowledge about what will work best in the future) and focus on the present to benchmark and combine the most efficient and effective directions.

Existing solutions to many tasks in mobile robotics, such as localisation, mapping, or planning, focus on prior knowledge about the structure of our tasks and environments. This may include geometry or kinematic and dynamic models, which therefore have been built into traditional programs. However, recent successes and the flexibility of fairly unconstrained, learned models shift the focus of new academic and industrial projects. Successes in image recognition (ImageNet) as well as triumphs in reinforcement learning (Atari, Go, Chess) inspire like-minded research.

This survey is by no means complete, its purpose is to provide a high-level review with more details to be found in the respective references. While the references only represent a small subset of available work in each field, the overall review demonstrates the regularities, connections and principles underlying the tasks and common software solutions.

## 1.1 The Broad Question

Recently, discussions have come up about the potential relevance of reinforcement learning for deployable mobile robots. When hearing these questions, it seems easy to reject them as an executive's fear of missing out on another hyped technology. However, it is worth taking a moment to investigate these questions.

Deep learning predominantly has made its mark regarding applications in the perception pipeline of autonomous systems including pedestrian / car / cyclist / traffic sign detection, semantic segmentation, and other related tasks. While these perception systems heavily rely on learning; localisation, reasoning, and planning modules often continue to be the domain of carefully crafted rules and programs, exploiting geometric priors and intuitions. The design of these systems requires expert knowledge and repeated iteration between testing - in simulation as well as on the real platform - and refinement of hundreds if not thousands of heuristics.

While for example in the early DARPA challenges, robotic systems nearly completely relied on these structures, the paradigm is starting to shift [164]. Given the success in perception tasks, the natural question is: 'what else can we learn from data?'. Discussions about using (reinforcement) learning naturally arise in the context of reducing manual efforts and instead automatically learning decision patterns. The overall question now focusses on the general application of learning in further parts of our pipeline; with RL representing one of the potentially more *high risk, high reward* scenarios.

While ML has been able to improve our efficiency in addressing various tasks, it does not represent *free lunch*. Independent of all its advantages, machine learning delivers no magic tool. Its successful application commonly requires detailed domain knowledge, systems engineering and demands significant time for data collection and curation, experimental setup and safety arrangements.
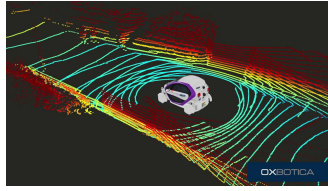
# 2 Learning for Autonomous Systems

Autonomous systems are generally modularised for the same reasons as any large software systems: reuseability, ease of testing, separation of responsibilities, interpretability, etc. Robots / autonomous systems are treated in this article as a collection of these modules, including: perception, localisation, mapping, tracking, prediction, planning, and control.

---

[1]A more colloquial version including the embedded videos can be found under https://markusrw.github.io/articles/tldr-on-ml-and-structure-for-robotics/ with a summary under http://ori.ox.ac.uk/on-learning-and-prior-structure/
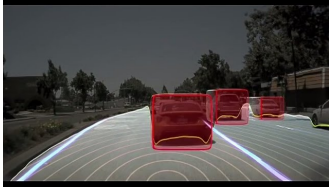
(a) Waymo (link)



(b) Oxford Robotics Institute / Oxbotica (link)



(c) Tesla (link)



(d) NVIDIA (link)



(e) Drive.ai (link)



(f) Toyota Research Institute (link)

The following paragraphs survey a subsection of work in each field, exemplifying the state of the art of learning based methods for these modules, followed by additional directions of relevance across the whole pipeline on uncertainty and introspection as well as representations. The following, final sections represent a more personal take on challenges and potentials. More resources on software systems and computer vision for autonomous platforms for interested readers can be found under [9, 189, 92] and generally in the mentioned references.
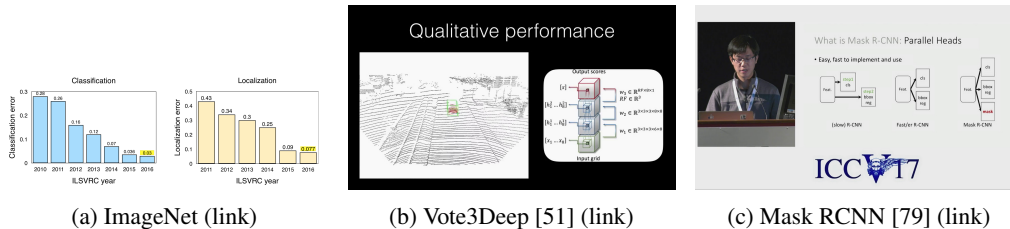
## 2.1  Perception

Current perception modules represent one of the principal success stories of deep learning in autonomous systems. Image classification, object detection, depth estimation, semantic segmentation, activity recognition are all principally dominated by deep learning [62, 34, 12] (a detailed survey of recent work can be found under [182]).

While classification benchmarks have been long-standing pillars of computer vision research, the ImageNet benchmark [161] in particular presents a cornerstone for the acceleration of progress in machine learning and in particular deep learning. A good share of models originally was developed specifically for this benchmark. ImageNet dataset as well as benchmark have massive forces for research on deep learning, which triumphed in all recent competitions.

The **detection** of traffic participants including pedestrians, cyclists and other vehicles [63, 35] relies predominantly on deep learning approaches for image [156, 25, 155] as well as LIDAR data [27, 51]. LIDAR can be understood as essentially a light based radar: the sensor's output being a long sequence of distance measurements. Notably, the natural structure of LIDAR significantly differs from image data and is unfit for the application of models designed for images. One essential challenge is that the same pointcloud can be represented by many different sequences and applied models have to be permutation invariant. Most early approaches are able to prevail by building on manually designed grids with predefined feature extractors for each cell [51, 74]. Replacing this kind of manual feature design, a more recent direction is the combination of low-level feature learning based on recurrent modules and high-level grid-based representations via convolutions for end-to-end training [149]. Further work relies on max-pooling as symmetric function over point-wise descriptors (treating the data as a set rather than as a sequence) and extensions to address local features at varying contextual scales [148, 150].

Similarly, **pixel-wise semantic and instance segmentation** [60, 137, 34, 62] as well as image-based **depth / disparity estimation** (mono and stereo) [166, 62, 65, 190, 107, 61] abide in the domain of deep learning based approaches. Interestingly however, ideas from geometric computer vision are making their way back into current research for the latter direction, e.g. in the form of reprojection losses [65]. Multi-task training with shared encoder segments has demonstrated additional improvements [131, 79] when parts of the architecture can be shared. Lastly, in general

(a) ImageNet (link)          (b) Vote3Deep [51] (link)          (c) Mask RCNN [79] (link)

**scene understanding**, deep learning has been beneficial for tasks such as the prediction of road attributes [168].

The applications in this section are predominantly mastered via learning. Innovation often focuses on architecture or loss function design, which exceed the scope of this review. Furthermore, a significant share of intriguing work takes place in product-oriented, more applied research, which is less published.
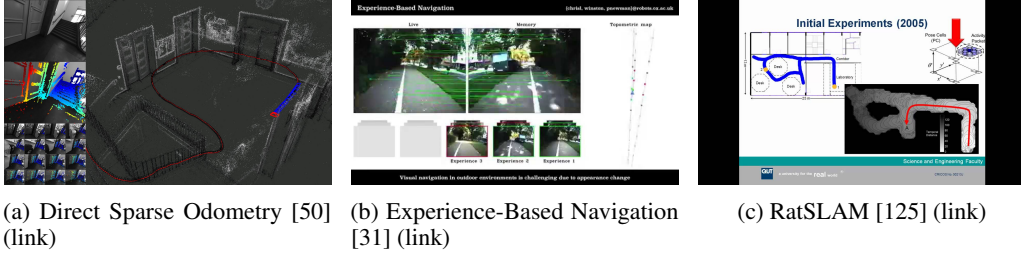
## 2.2 Localisation and Mapping

We will start this section by conveying condensed and highly simplified intuitions about the computations underlying current solutions to relative localisation, absolute localisation as well as the full SLAM (simultaneous localisation and mapping) problem.

In relative localisation, we commonly determine matching features in consecutive sensor measurements and, based on the changes in their coordinates, we compute our change in pose; the latter being accurately described via geometric rules (projection, triangulation, etc). Absolute localisation (localisation against a map) additionally involves a feature-matching process of current perception against locations in our map to determine the coarse location and potentially re-localise. For SLAM, we additionally build a map. While we determine our own location from the estimated positions of features, we estimate the position of new features with respect to the map to update and enhance. Additional refinement of the map for this joint optimisation problem is commonly formulated via iterative filtering and bundle adjustment techniques. The overall problem has many unmentioned challenges based on the efficient, robust realisation of these sub-tasks as well as in the complexity of real-world data including sensor noise, occlusions and dynamic environments.

Applications in localisation and mapping have provided challenging benchmarks for learning-based approaches. Geometric methods e.g. for **visual odometry** continue to outperform *end-to-end learning* [62] (end-to-end indicating here the learning of the complete odometry image-pose pipeline - in opposition to learning modules such as interest point descriptors [44]). Geometric methods have the benefit of incorporating our prior knowledge about exact geometric rules (e.g. regarding homographies and projections), which learning based methods will at best learn to approximate. While we are able to formulate exact equations e.g. for homography estimation, there are various tasks which commonly are solved more heuristically. The sub-optimal compression of available information, as in the context of feature descriptors, provides an opportunity for learning to minimise the loss of relevant information.

However, the actual gap between distinctly geometric or learned approaches for localisation is decreasing in practice due to consolidations of both directions. Recent work combines the flexibility of learned sub-systems and prior knowledge about task-dependent computations incorporating prior intuition from geometric CV [145, 200, 186, 26]. One example is given by the integration of auxiliary training losses to address the common drift problem of relative pose estimation [134]. Additionally to predicting accurate relative transforms, this objective can be applied to the integrated motion over multiple steps [145] to reduce accumulated drift. Furthermore, learning-based approaches provide the benefit of being independent of knowledge about (intrinsic camera) calibration as distorted images can be directly used [102].

**Absolute localisation** - relative to map instead of relative to our last position - commonly relies on such a map populated with features to localise against. Generally, in the context of current deep learning, most approaches utilise no explicit constraint on the type of computation [102, 20], there are however notable exceptions [217]). Recent work aims at harnessing geometric prior knowledge to obtain more informative training objectives [100] (with a survey under [99]).

4

(a) Direct Sparse Odometry [50] (link)   (b) Experience-Based Navigation [31] (link)   (c) RatSLAM [125] (link)

Early work on *neural approaches* to the full **SLAM** problem [33] is given by Milford and Wyeth [125] taking inspiration from computational models of the hippocampus of rodents (RatSLAM). More recently, Tateno and colleagues [181] apply learning to solve a sub-problem of SLAM: using a convolutional neural network as depth estimator (see Section 2.1) to overcome shortcomings of monocular SLAM regarding absolute scale. Another approach to *neural SLAM* is taken in [44] bu addressing another task within the SLAM systems and providing a fast deep learning based point tracking systems. An extension of their work (and combination with homography estimation [42]) extends past synthetic data and outperforms various learned and non-learned point detector/descriptor baselines on a range of tasks [43]. A final approach to neural SLAM is given by Zhang et al [217], who try to *embed procedures mimicking that of traditional Simultaneous Localization and Mapping (SLAM) into the soft attention-based addressing of external memory architectures, in which the external memory acts as an internal representation of the environment*. In essence, the authors aim at providing a model, which inherently encourages learning SLAM-like procedures. However, the evaluation focuses on limited toy examples in simulation.

To enable loop closures, **place recognition** addresses the recognition of previously visited locations based on their appearance and is a relevant part of the SLAM pipeline. This represents a general localisation sub-problem well-suited for learning-based approaches [6]. Commonly, it is targeted in the context of comparing the feature representations between potential candidates [176, 6] (the latter being a differentiable adaptation of the VLAD image descriptor [94]). Finally, appearance change in our environment continues to be one of the most challenging aspects for learning as well as geometric approaches [117]. Work on obtaining weather or lighting invariant representations to address this challenge is summarised in Section 2.7.
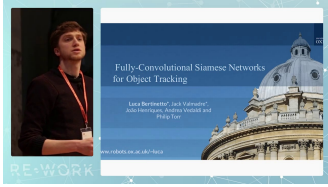
Predominantly, current combined methods for localisation focus on employing learning for sub-tasks which are only heuristically solved in traditional CV as well as utilising geometrically inspired structure and computations to *learn into*.
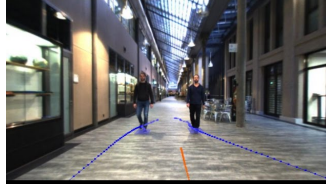
### 2.3 Tracking and Prediction

In essence, most object tracking pipelines can be divided into two steps: prediction of all tracks based on a prediction model and creation as well as update of the tracks based on current measurements; the latter depending on accurate assignments between current measurements and existing tracks.

Focusing on multi-object scenarios, one of the principal challenges of tracking lies in the **data association** problem: knowing which of the current detections corresponds to which established track. The most common methods for high-frequency tracking pipelines rely on simple distance-based associations between predicted position and current detections. However, in cluttered environments and the context of occlusions, additional information such as appearance is required to enable accurate tracking. Learning-based approaches have been successfully applied here e.g. to metric learning for appearance-based entity re-identification for pedestrians [114]. By applying objectives such as contrastive [75], triplet [84] and magnet loss [158] these methods learn a metric space where different instances of the same type reside closer together. Further methods train models for direct appearance-based tracking [106, 15] (more details in Section 3.1).

To associate existing tracks with new detections or provide position updates at potentially higher frequency than the received detections, we need to predict future positions. A classic, the (extended/unscented) Kalman filter [95, 198], is actually sufficient in most common situations. The simple **motion models** often underlying these methods (e.g. constant velocities) though will turn out unreliable in the context of long-term predictions and cluttered environments. For accurate

(a) Fully Convolutional Siamese Networks for Object Tracking [15] (link)



(b) Predict Actions to Act Predictably [144] (link)



(c) Deep Tracking (precursor to [41]) (link)

**prediction** of future trajectories, more information such as interactions with static scenery as well as other agents (cyclists, cars, pedestrians) needs to be considered.

Early work to represent these interactions with *social force* models [80] applies potential fields for modelling repulsive and attracting forces. Reciprocal Velocity Obstacles (RVO) were introduced as a computationally efficient extension with applications not only directly for prediction but also integrated into motion planning [194, 3] though the approach requires additional knowledge about the interacting entities. Recent work on learned predictive models employs deep neural networks to integrate information about static environment and dynamic environments [2, 52, 142]. In addition to flexibly utilising large quantities of raw sensor measurements, these approaches have been shown to be able to partially address the drift of trajectories by directly predicting multi-step sequences [142]. Notably, the KF itself has become a target for learning. BackprobKF [73] provides a fully differentiable architecture for state estimation which is evaluated on the KITTI visual odometry benchmark [62].

Given the perspective of robotics, we're often interested not just in the prediction of tracks for known and detected objects but the complete prediction of future states ,including aspects we do net explicitly handle in the detection module. Addressing this challenge, a different angle to tracking is given by approaches like 'Deep Tracking' [41] which predict complete future sensor observations (e.g. LIDAR and camera) [53, 122, 97, 197, 83]. These methods bypass the data association problem as well as the general detection challenge and can provide redundancy for the prediction of future observations, such as occupancy grids. However, learning generative models for the prediction of complete sensor measurements has so far proven particularly challenging.

In general, the prediction of future motion, in particular other agents' reactions, has great benefits for the following modules including motion planning [167, 144].

### 2.4 Planning and Control

Planning and control are the final components of our pipeline and the connecting modules to determine commands for actuation. A principal question for these modules is the type and source of supervision. While a significant share of currently deployed solutions builds on manually hand-crafted rules, learning provides a relevant alternative to prevent repeated hyperparameter and heuristic tuning for different environments and scenarios. Now, one solution for supervision signal can be through reinforcement learning, which - while representing a multi-faceted topic of its own - still needs to overcome many real-world challenges and simply is to intricate to cover as just a side aspect of this review. This section mostly focuses on **Learning from Demonstration** to provide supervision based on demonstrations of a task from human experts and other potential authorities.

**Behavioural Cloning** (BC) aims at directly mimicking expert behaviour to solve a task; essentially supervised optimisation of regression or classification models. BC can be integrated into existing pipelines to build on more abstract representations but most commonly has been investigated in the scenario of end-to-end learning based on raw inputs. These methods have been empirically demonstrated in constrained scenarios e.g. for lane keeping [18, 147, 128]. Given independence of the source of demonstration data, BC is not restricted to imitate human experts and can be applied with automatically generated trajectories [143].

However, this application of naively trained supervised models *in a non-iid scenario* comes with additional challenges as performed actions affect future input data. Small errors lead to the data distribution diverging from the training data, which commonly focuses on states along optimal

(a) Cost-Function Learning via IRL [213] (link)　　(b) Terrain Classification [11] (link)　　(c) IRL for Flying [1] (link)

trajectories, a phenomenon known as *covariate shift*. Model performance degrades, potentially shifting the input data even further from our training data distribution, causing a *compounding of errors* [160]. To prevent this result, we need to learn how to recover from suboptimal states, which are usually not part of given expert demonstrations. [18] addresses the problem in the context of lane keeping by synthetically generating off lane-centre states with additional cameras on the sides of the vehicle with corrected steering manoeuvres. One of the most common approaches is presented by DAgger [160] and extensions, which collect additional expert supervision during application of the model - however this leads to increased efforts for providing supervision [109, 96, 108]. Finally, this type of end-to-end modelling is limited in terms of interpretability, and can rely on larger amounts of training data than modular or abstracted approaches [171].
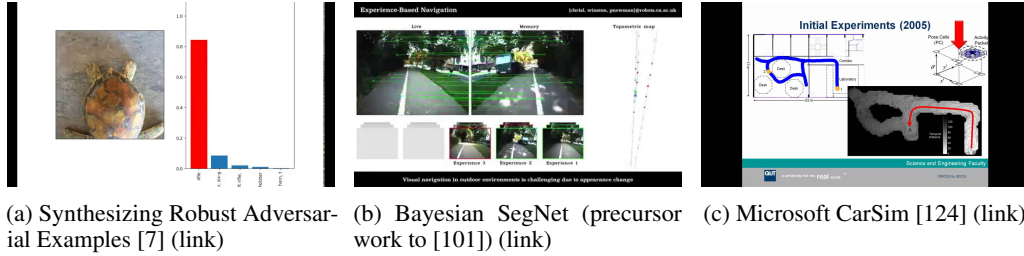
**Inverse Reinforcement Learning** (IRL) presents another popular approach to address the problem of *covariate shift* - by blending supervised learning with reinforcement learning (RL) or planning to learn robust models. IRL aims to infer expert preferences by optimising a reward function that generates agent behaviour similar to the expert demonstrations - instead of an imitating policy. In the context of probabilistic models, it can be understood as optimising a model that maximises the probability of the expert's trajectories [220].

While BC only learns accurate behaviour for expert-visited states, IRL extends to states visited by the RL or planning step and learns corrective behaviour when diverging from the original trajectories. Furthermore, recent work directly utilises human domain knowledge [213] to define behaviour for states not sufficiently encountered by either. However, IRL in comparison to BC relies on an accurate systems model to simulate behaviour or the possibility to sample on the real system. If both are not possible we can turn, in the context of mobile robotics, to another approach based on supervised learning: training supervised segmentation models for traversable terrain [11, 183].

Though the IRL problem is underconstrained, as many reward functions are able to describe the same optimal behaviour, various approaches have introduced simple assumptions to address the degeneracy and derive efficient, practical solutions [220, 152, 28]. Based on these, impressive successes of IRL include learning artistic flying manoeuvres for an RC helicopter [1] and predicting future motion for traffic participants such as cars and pedestrians [104, 221, 133].

A major benefit of IRL-based methods lies in the integration into existing systems. By applying methods to learn cost (negated reward) functions for existing motion planning systems [153, 213, 172], we can directly integrate learned models into deployable systems, which can straightforwardly be tested and benchmarked against existing, hand-crafted cost functions.

When treating the robot control problem as part of a multi-agent scenario, we aim to optimise our actions not just for internal goals but as well for interaction and the internal goals of other agents. This approach gains relevance in cases when the other agents are represented by humans, as is the case in most robot applications. Research on **interpretability** of models has a long-standing history and recently gained increased attention based on the massive complexity and - more importantly - real world relevance of deep learning [76, 115, 39]. The eminent aspect for control design is the interpretability of robot behaviour, urging us to *act predictably* when directly interacting with humans [173]. When planning motions in the direct vicinity of humans we benefit from providing non-verbal cues and human-like behaviour, enabling others to infer our driving style [163] and intentions [144, 88] to ensure convenient and stressfree interaction. Broader surveys on legible behaviour for human robot interaction can be found in [48, 49].

(a) Synthesizing Robust Adversarial Examples [7] (link)



(b) Bayesian SegNet (precursor work to [101]) (link)



(c) Microsoft CarSim [124] (link)

## 2.5 Safety, Uncertainty and Introspection

The following sections present aspects which are disjoint of the modular pipeline structure addressed before and cover more general, cross-module concerns and potentials of learning from data in robotics.

Notably, relevant performance metrics in active, sequential decision making such as robotics and autonomous platforms differ from the metrics commonly benchmarked for machine learning (such as classification accuracy, precision, recall, etc.). First, not all mistakes are equal: making a mistake confidently can be much more harmful than demonstrating **uncertainty** about the situation (e.g. the class of a pointcloud - pedestrian versus paper-bag). False negatives and false positives have massively different relevance for our modelling decisions. Second, we are able to act conservatively in the context of uncertainty and additionally probe the environment by collecting more data instead of forcing the model to make a confident prediction. Given the safety requirements for wide-scale application of autonomous vehicles [205], this approach is highly compelling.
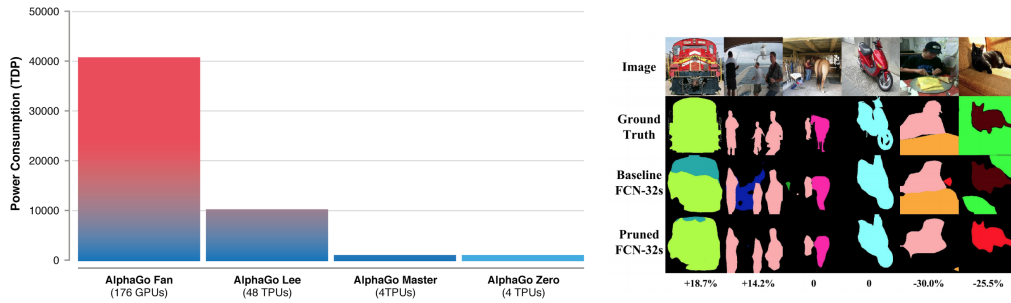
To determine the necessity of conservative behaviour, we depend on knowledge about the model's uncertainty for its predictions, which is commonly investigated as model **introspection**. [68] investigates SVMs, GPs and a number of other popular classification models (pre deep learning) and empirically support the intuition that *better introspection leads to improved decision making in the context of tasks such as autonomous driving or semantic map generation*. Furthermore, the authors indicate that commonly used detection metrics of precision and recall can be insufficient for describing model performance in safety-critical applications.

Furthermore, commonly used *pseudo probabilities* for introspecting model predictions, such as detection scores or softmax output, often do not suffice and supplemental uncertainty measures are required. *Bayesian uncertainty* modelling for deep learning has found interest long before the current AI summer [17, 119, 129, 58], but recently received more spotlight thanks to the application of deep learning in safety-critical environments. [101] investigates the use of *aleatoric and epistemic* uncertainty metrics in deep learning, respectively for describing noise inherent in the observations and model uncertainty. Notably, while we cannot easily affect observation noise on the software side, we can reduce model uncertainty by collecting more data. Both metrics are suited for different purposes. While aleatoric uncertainty becomes highly relevant in large scale applications where epistemic uncertainty can be neglected, epistemic uncertainty can be employed to determine covariate shift between training and application data. In addition, this detection of *novel* data can be addressed via generative models [202].

While the previously mentioned work aims at investigating and extending the capability of the model itself, a separate direction of research aims at **redundancy** and parallel streams of information to examine model predictions. Explicit introspection tools have been introduced for predicting the performance of perception modules [72], the whole pipeline from perception to planning [38] and control [55].

Finally, **saliency detection** presents another essential approach to investigate model predictions, identifying input sections of high relevance. These visualisations are essentially obtained by determining the minimal input changes required to change model predictions [54, 37, 169, 203].

In order to accurately understand strengths and weaknesses of our system and improve on the latter, we depend on thorough testing and refinement of our systems. Training and testing systems in **simulation** enables us to repeat and vary edge cases. In essence, it enables us to generate multiple orders of magnitude more driving data [204, 184, 188] at different granularities. Various datasets and simulators are openly available for research [47, 124, 157, 159, 207, 57, 170] with a more

(a) AlphaGo Power Consumption (source: businessinsider.com) (link)



(b) Learning to Prune in CNNs [87] (link)

comprehensive list given in [191]. While current simulators become increasingly accurate, the reality gap still persists and and emphasises the potential of related work on transfer learning and domain adaptation in Section 2.7.

Just like every other software program, learning-based approaches have their **vulnerabilities** and can be fooled. New types of *adversarial attacks* [214] - ways to mess with the input data to fool the model - as well as methods for defence have gained increased attention in the past years. Most impressively, recent work presents more general approaches and has shown adversarial examples with invariance with respect to 3D viewpoint [7] and attacked model [127]. Furthermore, [103] demonstrates the possibility of attacks on the training data set. Interesting work on defence against adversaries includes input transformations [70] and different encodings [5].

## 2.6 Knowledge Representation and Efficient Models

Applications on mobile platforms and embedded systems have increased the demand for **computationally efficient** systems with reduced memory footprint aiming at on-chip rather than off-chip placement. The situation has lead to improvements of over an order of magnitude reduction in parameters, FLOPS, and the corresponding increase in possible frame-rate compared to previous state-of-the-art models with only limited reduction in accuracy [86, 136, 85]. Notable, the basic building blocks such as convolutions have been adapted via the introduction of asymmetric [178] and separable convolutions [30] to increase parameter efficiency.

Work on model compression based on *pruning, trained quantization and Huffman coding* was able to reduce the memory footprint of existing architectures [77, 87]. Furthermore, instead of directly compressing the model, *knowledge distillation* enables the training of smaller models via mimicking the predictions of large state-of-the-art architectures without their computational footprint [22, 81]. The underlying intuition being that the logits extracted by the more powerful network include more information than the hard one-hot encodings, in particular about relations between classes. Recently, it has additionally been shown that distillation between networks of the same architecture can increase performance [56].

Finally, a role for learning can be found in *reducing computation and time requirements* [143]. In this context, one perspective on AlphaGoZero [174] covers the aspect of learning to imitate more expensive computations (here: MCTS), which enables the final, trained program to run faster and with lower computational requirements. While the version of AlphaGo that bested Lee Sedol had an estimated **power consumption** of approximately 1 MW (50,000 times as much power as the amount of power required for a human brain), AlphaGoZero, which beat the previous version 100-0, uses an order of magnitude less compute.

## 2.7 Transfer, Multimodal, and Representation Learning

Robotics platforms perceive their environment through a multitude of different sensors. Learning can aid to combine and analyse this flood of information. Integrated training with **different sensing modalities** enables us to capture joint distributions, deploy with restricted access to our sensor setup [132, 151, 192] and increases robustness to sensor failure [116].

9

(a) High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs [201] (link)

(b) Deep Multispectral Semantic Scene Understanding [192] (link)

(c) Deep Photo Style Transfer [118] (link)

The additional result of these high-throughput sensor setups is the generation of massive amounts of unsupervised data, exceeding our capability for manual dense annotation. Nonetheless, we can benefit from this overwhelming amount of data by splitting the problem in two: First, unsupervised learning of a representation from which we require less supervised data and second, provide supervision for determining the final mapping. By reducing the requirements for human annotation, **representation learning** (unsupervised learning) has significant potential, though has not yet found the same commercial success as supervised approaches. One of our challenges is that '*we [often] don't know what's a good representation*' [13]. Essentially, we're only able to **benchmark** in the context of surrogate tasks such as reconstruction accuracy or performance of subsequent classifier modules.

Various approaches aim at finding relevant representations based on such proxy objectives; the aspects in common for most approaches is the prediction or verification of **structure - spatial and temporal**. Predicting spatial structure includes the relative location of image patches [45], the order of shuffled image patches [135], image inpainting [138], or employing various *foundational [supervised] proxy 3D tasks* for learning a generic 3D representation [216]. The recent increase in compute capacity enables the extension of these ideas from the spatial to the temporal domain, thereby facilitating the use of temporal consistency and structure to learn representations from videos (e.g. [177]). Examples operate by verifying the temporal order of sequences [126], predicting future frame representations [196], or by predicting low-level motion-based clustering [137]. Further detailed views on representation learning can be found in survey form [14], in blog form on predictive versus representation learning [69], or in blog form on recent work and extensions to find training signals for RL [139].

One particularly relevant type of representation learning in this context is the field of unsupervised domain adaptation, which aims to induce *domain invariant embeddings* such as to increase the performance for a task in domains without annotated data. The metric for evaluation is clearly defined in this case by the supervised task. Similar to various approaches emphasised in this review, the benefit of domain-invariant representations lies in *addressing covariate shift*, the encounter of data outside our supervised training distribution during deployment of the model. Essentially, while acquiring widespread training and validation data is most beneficial [90], the effort is often impractical based on expenses and the challenge of considering all potential conditions. Recent approaches aim at empowering domain adaptation methods via deep models [59, 210, 19, 187, 36]; furthermore extending to adaptation in continually changing environments [211]. Additionally to domain adaptation in the feature space of a model, transforming images between different domains has become a promising direction [19, 219, 29, 201, 89]. Of course, fine tuning of pretrained models [215, 46] from different domains or different tasks remains nearly [212] always helpful if small amounts of supervised data are available for the application domain.

## 3  Challenges and Potentials

After the review of current work on learning in different modules as well as at their intersections, this section concludes with promising directions and potentials as well as some challenges ahead.

## 3.1 Bottom-Up

The principal strength of learning lies in applications in direct, reactive **perception** problems including classification, detection, segmentation, and related tasks. State-of-the-art models provide high accuracy with comparably little encoded structure (other than e.g. convolutions).

While there is no strong foundation for believing that **localisation** tasks are inherently much harder (or easier) to learn than e.g. object detection, geometric approaches simply provide a much stronger baseline in this field[62, 218]. Traditionally, CV addresses localisation very structured by utilising geometric knowledge about the mathematical rules underlying particular sub-problems (triangulation, homography estimation, etc). Instead of aiming to formalise rules for complete mappings from images to relative or absolute poses, odometry or SLAM systems extract and match features across sequences or between current perception and constructed map (while constructing said map), utilising prior knowledge about geometric properties, and improving pose graphs on the back-end, e.g. via filtering or bundle adjustment [33]. However, the first generation of learning-based methods addressed the problem with little to no structure, learning end-to-end mappings and based on pre-existing, annotated datasets (e.g. [102]). While given infinite data, covering the complete distribution of interest, and a flexible enough model, this approach can perform really well, the reality is often more constrained.

By combining both directions, the strength of geometric methods can have a substantial contribution towards more robust and reliable approaches. Various aspects of the task can be solved to perfect accuracy given geometric knowledge, while other aspects are affected by the real-world noise and can benefit via learning from data. Components that are likely to be focus on improvement via machine learning in future work include feature extraction e.g. for loop closure detection and re-localization or better point descriptors for sparse SLAM methods. Finally, deep learning can greatly improve the quality of map semantics - i.e. going beyond poses or pointclouds to a more complete understanding of the types and functionalities of objects or regions in the map. One particularly relevant direction lies in improved robustness (for example through better handling of dynamic objects and environmental changes [10]). Learning will additionally be beneficial when the assumptions underlying traditional approaches are invalid or we require faster methods [100]. Back at ICCV 2015, the question if deep learning would replace geometric CV for SLAM might have been received with significant scepticism [121] (and might still be), however most research has actually not tried to replace geometry but instead to enhance and augment, only learning parts of the system where cannot provide exact prior structure.

Commonly, the accuracy of perception and localisation systems can represent a **bottleneck** for overall performance and safety of a system as they provide the foundation for the rest of the pipeline, emphasising the relevance of even minute improvements. Part of the direct *consumers* of their output are the **tracking and prediction** modules. Traditional tracking approaches (e.g. EKF, UKF) often provide reasonable solutions for the space of robotics. These methods are straightforward to implement and fully suffice as long as tracking itself does not become the bottleneck. Applications in more complex, cluttered environment, where data association becomes more challenging, represent a prime example of the benefits of learning for appearance-based tracking [106, 105, 193, 15]. Furthermore, densely populated scenarios often lead to more intricate interactions. Learned elaborate interactive motion models (in Section 2.3) enable us to predict motion with higher accuracy in these environments and can be incorporated into existing trackers.

Traditionally, **motion planning**, at least for deployed platforms, has been one of the fields more resistant to learning-based approaches (in particular in safety-critical applications). Planning approaches represent structured procedures for reasoning [110, 146], utilising knowledge about e.g. kinematic and dynamic constraints, as well as geometric extension of platform and obstacles. As with localisation, parts of the planning computation are accurately modelled via known structure and equations (e.g. collision checking). In this context, learning focuses on more intuitive aspects, e.g. in improving prediction in interactive scenarios - such as highway lane merging - which requires to predict the reaction of other cars to potential actions. In essence, the more interactive and intuitive parts of driving, which are less governed by strict, easy-to-define rules, present opportunities for learning from data, where these forms of *common sense* and *intuition* are too complex for manual rules [40]. Recent work on imitation learning for driving for example outsources high-level planning and takes additional commands as input [32] to focus learning on what is more straightforward to learn. This approach does not learn to plan but essentially a reactive controller (based on raw images) and presents another example for merging learning with existing systems.

(a) Conditional Imitation Learning [32] (link)    (b) Cognitive Mapping and Planning [71] (link)

Route planning, as topological, high-level planning process, is commonly addressed via graph search (A* and related approaches). While there has been research on learning these kinds of programs as end-to-end approach in limited scenarios [67, 130], given current applications, the existing algorithms do not represent a bottleneck. However, it can be expected that the costs associated with edges and nodes for the route graph are well-suited for estimation from data.
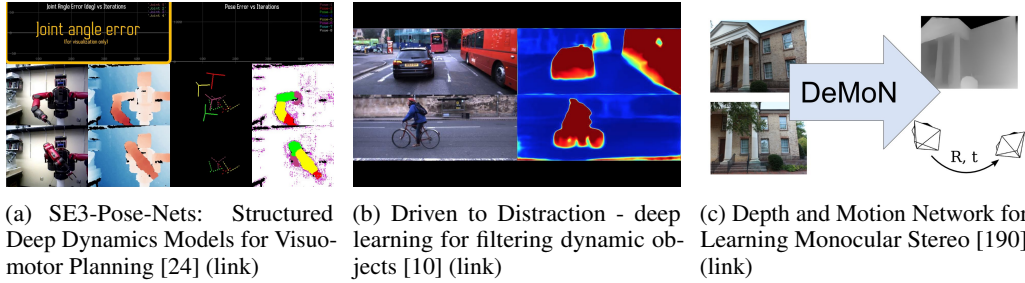
Focusing on the incorporation of learning for planning and control into existing modular software pipelines, another principal application lies in the characterisation of traversability and obstacles as well as the prediction of the reactions of other traffic participants. Hand-crafted cost functions for different kinds of terrains are commonly designed to help bridge the gap between perception and action and reduce the complexity of our environment representation to focus on the aspects we care about. Learning cost functions for *driving like human experts* (as e.g. in Section 2.4) is addressed via a sub-field of learning from demonstration [172, 213]. Furthermore, similar techniques can be applied to learn to predict reactive behaviours - interweaving planning and prediction models for dynamic environments [144]. Finally, in addition to manually defined cost functions, planning and reasoning systems commonly include many other heuristics, parameters determined during deployment to work in tested scenarios. Here, learning can play a major role in determining general rules on how to turn these knobs to adapt to new environments or different user preferences.

In the context of safety-critical applications, learning is suitable for the generation of parallel systems for **redundancy** and additional checks, modules that replicate functionality and enable the second-guessing of decisions. Optimally, the application of multiple redundant modules is not restricted to verification but can culminate in a framework for **learning from disagreement**, where disagreements between modules are not only detected but feed back into the optimisation of the overall system, such as through adaptation of a module's uncertainty for future predictions.

A recent example for this type of framework is given by Pei et al [140]. The authors devise an approach for sets of networks, retraining the modules that disagree with majority decisions. The underlying assumption, that the majority will always be right, is critical for success of the approach. Their training process aims at balancing two objectives: maximise the number of active neurons and trigger as many conflicts between the modules as possible. This objective is interestingly similar to basic ideas from software testing aiming to maximise code coverage. Variations of the idea aiming at adapting uncertainties and taking module uncertainty into account for a weighted majority represent promising further directions. Furthermore, it will be generally beneficial to address the transfer of various other concepts and paradigms which have been demonstrated successful and indispensable for software engineering.

## 3.2 Top-Down

Machine learning (in particular deep learning) has the capability to extract rules from massive amounts of data and the benefit of high flexibility: merging of arbitrary objectives, Lego-like capabilities for the reuse and combination of models. On the other hand, we have accurate mathematical formulations about the underlying math and programs to solve specific sub-problems e.g. for localisation and planning. Modern deep learning improves the ease for the integration of fixed and learned modules, enabling us to be standing on the shoulders of giants from both fields and build on known solutions.

(a) SE3-Pose-Nets: Structured Deep Dynamics Models for Visuo-motor Planning [24] (link)

(b) Driven to Distraction - deep learning for filtering dynamic objects [10] (link)

(c) Depth and Motion Network for Learning Monocular Stereo [190] (link)

While potentially trained as independent modules, the overall trajectory goes towards combinations which can be optimised as complete system via deterministic gradients as well as - if required - various stochastic gradient estimators (REINFORCE (or likelihood-ratio) trick [208], evolution strategies [8, 165], continuous relaxations such as Gumbel-Softmax or Concrete distribution and extensions [120, 93, 185] ).

Combinations of both approaches can provide significant advantages via **redundant** systems and often **complementary** properties. In essence, we aim to take *the best of both worlds* when merging two systems; similarly to how automation via ML aims to adapt and enhance job responsibilities [123] by addressing tasks complementary to human strengths. Ongoing directions include **optimising input data or correcting the output** of traditional programs. Examples for input improvement include learning image enhancement networks for traditional visual odometry methods [66]; output refinement includes pose correction updates for visual localisation [141] and refining dense reconstructions [180] as well as hand-crafted cost maps for motion planning [213].

Similarly, pure learning-based approaches benefit from **incorporating prior knowledge** about the underlying structure: implicit and explicit translation invariance [112, 209], *objectness* [23], temporal structure [82], planning procedures such as value iteration [179], further geometric properties [78, 91], structure that encourages SLAM-like computations [217] and access to SLAM - location and map - information for reinforcement learning [16]. Notably, the incorporation of structure can, under specific circumstances, even help when the incorporated models are inaccurate [206].

The combination of prior geometric knowledge and the flexibility of learning enables the reformulation of geometric properties to create **self-supervised objectives**. Examples are given by [195, 24] which utilise predictions for depth, segmentation masks and poses to differentiably warp frames in time to match image sections. In addition to training with externally supervised labels these approaches enable self-supervised learning e.g. via reprojection photometric errors. Lastly, the use of data augmentation represents the application of prior knowledge, structuring our wanted invariances via randomising the relevant aspects in our training data. A related survey on limits and potentials of deep learning in robotics can be found in [175].

When moving from manually defining features and computations to designing the most efficient structure to learn into, one question arises naturally: why not **learn everything** (the architecture [154, 21], the optimiser [199, 4], or complete programs [98]). However, required investments in data hygiene and annotation for many applications with potential for real world impact often render it more efficient, in terms of human effort, to port our prior knowledge into algorithmic structure. Leslie Kaelbling formulated this well during a panel session at CoRL2017: 'What structure can we build in that does not obstruct learning?'. The point being twofold, with respect to model and the optimisation procedure. If structure is a *necessary good* or *necessary evil* might be up to discussion [111, 113, 64, 162], but for now, practically, it is **necessary** as well as are the advantages of learning.

Building autonomous platforms, like addressing any other sufficiently complex and versatile software problem, results in a significant effort for **systems engineering and iterative testing and refinement**. The relative emphasis on learning or traditional programming blocks narrows down to the required effort and efficiency when creating reliable, safe and generalisable systems with either approach as well as the potential benefits of combination.

## Acknowledgements

## References

[1] P. Abbeel, A. Coates, and A. Y. Ng. Autonomous helicopter aerobatics through apprenticeship learning. *The International Journal of Robotics Research*, 29(13):1608–1639, 2010.

[2] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–971, 2016.

[3] J. Alonso-Mora, A. Breitenmoser, M. Rufli, P. Beardsley, and R. Siegwart. *Optimal Reciprocal Collision Avoidance for Multiple Non-Holonomic Robots*, pages 203–216. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.

[4] M. Andrychowicz, M. Denil, S. Gomez, M. W. Hoffman, D. Pfau, T. Schaul, B. Shillingford, and N. de Freitas. Learning to learn by gradient descent by gradient descent. *ArXiv e-prints*, 2016.

[5] Anonymous. Thermometer encoding: One hot way to resist adversarial examples. *International Conference on Learning Representations*, 2018.

[6] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic. Netvlad: Cnn architecture for weakly supervised place recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5297–5307, 2016.

[7] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok. Synthesizing Robust Adversarial Examples. *ArXiv e-prints*, 2017.

[8] T. Back. *Evolutionary algorithms in theory and practice: evolution strategies, evolutionary programming, genetic algorithms*. Oxford university press, 1996.

[9] Baidu. Baidu apollo. http://apollo.auto/, 2017.

[10] D. Barnes, W. Maddern, G. Pascoe, and I. Posner. Driven to Distraction: Self-Supervised Distractor Learning for Robust Monocular Visual Odometry in Urban Environments. *ArXiv e-prints*, 2017.

[11] D. Barnes, W. Maddern, and I. Posner. Find Your Own Way: Weakly-Supervised Segmentation of Path Proposals for Urban Autonomy. *ArXiv e-prints*, 2016.

[12] R. Benenson. Cv benchmarks - last update 2016. http://rodrigob.github.io/are_we_there_yet/build/ , 2016.

[13] Y. Bengio. Yoshua bengio interview. https://www.youtube.com/watch?v=pnTLZQhFpaE, 2017.

[14] Y. Bengio, A. Courville, and P. Vincent. Representation Learning: A Review and New Perspectives. *ArXiv e-prints*, 2012.

[15] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr. Fully-convolutional siamese networks for object tracking. In *European Conference on Computer Vision*, pages 850–865. Springer, 2016.

[16] S. Bhatti, A. Desmaison, O. Miksik, N. Nardelli, N. Siddharth, and P. H. S. Torr. Playing Doom with SLAM-Augmented Deep Reinforcement Learning. *ArXiv e-prints*, 2016.

[17] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.

[18] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.

[19] K. Bousmalis, A. Irpan, P. Wohlhart, Y. Bai, M. Kelcey, M. Kalakrishnan, L. Downs, J. Ibarz, P. Pastor, K. Konolige, S. Levine, and V. Vanhoucke. Using Simulation and Domain Adaptation to Improve Efficiency of Deep Robotic Grasping. *ArXiv e-prints*, 2017.

[20] E. Brachmann, F. Michel, A. Krull, M. Ying Yang, S. Gumhold, and c. Rother. Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[21] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. SMASH: One-Shot Model Architecture Search through HyperNetworks. *ArXiv e-prints*, 2017.

[22] C. Buciluǎ, R. Caruana, and A. Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541. ACM, 2006.

[23] A. Byravan and D. Fox. SE3-Nets: Learning Rigid Body Motion using Deep Neural Networks. *ArXiv e-prints*, 2016.

[24] A. Byravan, F. Leeb, F. Meier, and D. Fox. Se3-pose-nets: Structured deep dynamics models for visuomotor planning and control. *CoRR*, abs/1710.00489, 2017.

[25] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos. A Unified Multi-scale Deep Convolutional Neural Network for Fast Object Detection. *ArXiv e-prints*, 2016.

[26] L. Carlone and S. Karaman. Attention and Anticipation in Fast Visual-Inertial Navigation. *ArXiv e-prints*, 2016.

[27] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia. Multi-View 3D Object Detection Network for Autonomous Driving. *ArXiv e-prints*, 2016.

[28] J. Choi and K.-E. Kim. Map inference for bayesian inverse reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 1989–1997, 2011.

[29] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation. *ArXiv e-prints*, 2017.

[30] F. Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint arXiv:1610.02357*, 2016.

[31] W. Churchill and P. Newman. Experience-based navigation for long-term localisation. *The International Journal of Robotics Research*, 32(14):1645–1661, 2013.

[32] F. Codevilla, M. Müller, A. Dosovitskiy, A. López, and V. Koltun. End-to-end Driving via Conditional Imitation Learning. *ArXiv e-prints*, 2017.

[33] cometlabs. Slam systems overview. https://blog.cometlabs.io/teaching-robots-presence-what-you-need-to-know-about-slam-9bf0ca037553, 2017.

[34] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. 2016.

[35] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3223, 2016.

[36] G. Csurka. Domain adaptation for visual applications: A comprehensive survey. *arXiv preprint arXiv:1702.05374*, 2017.

[37] P. Dabkowski and Y. Gal. Real Time Image Saliency for Black Box Classifiers. *ArXiv e-prints*, 2017.

[38] S. Daftry, S. Zeng, J. A. D. Bagnell, and M. Hebert. Introspective perception: Learning to predict failures in vision systems. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2016)*. IEEE, July 2016.

[39] DARPA. Darpa explainable artificial intelligence. https://www.cc.gatech.edu/ alan-wags/DLAI2016/(Gunning)2017.

[40] E. Davis and G. Marcus. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM*, 58(9):92–103, 2015.

[41] J. Dequaire, P. Ondrúška, D. Rao, D. Wang, and I. Posner. Deep tracking in the wild: End-to-end tracking using recurrent neural networks. *The International Journal of Robotics Research*, page 0278364917710543, 2017.

[42] D. DeTone, T. Malisiewicz, and A. Rabinovich. Deep Image Homography Estimation. *ArXiv e-prints*, 2016.

[43] D. DeTone, T. Malisiewicz, and A. Rabinovich. SuperPoint: Self-Supervised Interest Point Detection and Description. *ArXiv e-prints*, 2017.

[44] D. DeTone, T. Malisiewicz, and A. Rabinovich. Toward Geometric Deep SLAM. *ArXiv e-prints*, 2017.

[45] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1422–1430, 2015.

[46] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. *ArXiv e-prints*, 2013.

[47] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. Carla: An open urban driving simulator. In *Conference on Robot Learning*, pages 1–16, 2017.

[48] A. D. Dragan. Legible robot motion planning. 2015.

[49] A. D. Dragan. Robot Planning with Mathematical Models of Human State and Action. *ArXiv e-prints*, 2017.

[50] J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

[51] M. Engelcke, D. Rao, D. Zeng Wang, C. Hay Tong, and I. Posner. Vote3Deep: Fast Object Detection in 3D Point Clouds Using Efficient Convolutional Neural Networks. *ArXiv e-prints*, 2016.

[52] T. Fernando, S. Denman, S. Sridharan, and C. Fookes. Soft + Hardwired Attention: An LSTM Framework for Human Trajectory Prediction and Abnormal Event Detection. *ArXiv e-prints*, 2017.

[53] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in Neural Information Processing Systems*, pages 64–72, 2016.

[54] R. Fong and A. Vedaldi. Interpretable Explanations of Black Boxes by Meaningful Perturbation. *ArXiv e-prints*, 2017.

[55] L. Fridman, B. Jenik, and B. Reimer. Arguing Machines: Perception-Control System Redundancy and Edge Case Discovery in Real-World Autonomous Driving. *ArXiv e-prints*, 2017.

[56] T. Furlanello, Z. Lipton, L. Itti, and A. Amandkumar. Born again neural networks. In *NIPS Workshop on Meta Learning*, 2017.

[57] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig. Virtual Worlds as Proxy for Multi-Object Tracking Analysis. *ArXiv e-prints*, 2016.

[58] Y. Gal. What my deep model doesn't know... http://mlg.eng.cam.ac.uk/yarin/website/blog_3d801aa532c1ce.html, 2016.

[59] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, V. Lempitsky, U. Dogan, M. Kloft, F. Orabona, and T. Tommasi. Domain-Adversarial Training of Neural Networks. *Journal of Machine Learning Research*, 17:1–35, 2016.

[60] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez. A Review on Deep Learning Techniques Applied to Semantic Segmentation. *ArXiv e-prints*, 2017.

[61] R. Garg, G. Carneiro, and I. Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *European Conference on Computer Vision*, pages 740–756. Springer, 2016.

[62] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013.

[63] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

[64] D. George, W. Lehrach, K. Kansky, M. Lázaro-Gredilla, C. Laan, B. Marthi, X. Lou, Z. Meng, Y. Liu, H. Wang, A. Lavin, and D. S. Phoenix. A generative vision model that trains with high data efficiency and breaks text-based captchas. *Science*, 2017.

[65] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017.

[66] R. Gomez-Ojeda, Z. Zhang, J. Gonzalez-Jimenez, and D. Scaramuzza. Learning-based Image Enhancement for Visual Odometry in Challenging HDR Environments. *ArXiv e-prints*, 2017.

[67] A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwińska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. Agapiou, A. Badia, K. M. Hermann, Y. Zwols, G. Ostrovski, A. Cain, H. King, C. Summerfield, P. Blunsom, K. Kavukcuoglu, and D. Hassabis. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538:471 EP –, 10 2016.

[68] H. Grimmett, R. Triebel, R. Paul, and I. Posner. Introspective classification for robot perception. *The International Journal of Robotics Research*, 35(7):743–762, 2016.

[69] R. Grosse. Predictive learning vs. representation learning. https://hips.seas.harvard.edu/blog/2013/02/04/predictive-learning-vs-representation-learning/, 2013.

[70] C. Guo, M. Rana, M. Cisse, and L. van der Maaten. Countering Adversarial Images using Input Transformations. *ArXiv e-prints*, 2017.

[71] S. Gupta, J. Davidson, S. Levine, R. Sukthankar, and J. Malik. Cognitive Mapping and Planning for Visual Navigation. *ArXiv e-prints*, 2017.

[72] C. Gurău, D. Rao, C. H. Tong, and I. Posner. Learn from experience: probabilistic prediction of perception performance to avoid failure. *The International Journal of Robotics Research*, 0(0):0278364917730603, 0.

[73] T. Haarnoja, A. Ajay, S. Levine, and P. Abbeel. Backprop KF: Learning Discriminative Deterministic State Estimators. *ArXiv e-prints*, May 2016.

[74] T. Hackel, N. Savinov, L. Ladicky, J. D. Wegner, K. Schindler, and M. Pollefeys. SEMAN-TIC3D.NET: a New Large-Scale Point Cloud Classification Benchmark. *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2017.

[75] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *Computer vision and pattern recognition, 2006 IEEE computer society conference on*, volume 2, pages 1735–1742. IEEE, 2006.

[76] P. Hall, W. P. Sri, and S. Ambati. Ideas on interpreting machine learning. https://www.oreilly.com/ideas/ideas-on-interpreting-machine-learning, 2017.

[77] S. Han, H. Mao, and W. J. Dally. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. *ArXiv e-prints*, 2015.

[78] A. Handa, M. Bloesch, V. Patraucean, S. Stent, J. McCormac, and A. Davison. gvnn: Neural Network Library for Geometric Computer Vision. *ArXiv e-prints*, 2016.

[79] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. *ArXiv e-prints*, 2017.

[80] D. Helbing and P. Molnar. Social force model for pedestrian dynamics. *Physical review E*, 51(5):4282, 1995.

[81] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[82] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[83] S. Hoermann, M. Bach, and K. Dietmayer. Dynamic Occupancy Grid Prediction for Urban Autonomous Driving: A Deep Learning Approach with Fully Automatic Labeling. *ArXiv e-prints*, 2017.

[84] E. Hoffer and N. Ailon. Deep metric learning using triplet network. In *International Workshop on Similarity-Based Pattern Recognition*, pages 84–92. Springer, 2015.

[85] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *ArXiv e-prints*, 2017.

[86] G. Huang, S. Liu, L. van der Maaten, and K. Q. Weinberger. CondenseNet: An Efficient DenseNet using Learned Group Convolutions. *ArXiv e-prints*, 2017.

[87] Q. Huang, K. Zhou, S. You, and U. Neumann. Learning to Prune Filters in Convolutional Neural Networks. *ArXiv e-prints*, 2018.

[88] S. H. Huang, D. Held, P. Abbeel, and A. D. Dragan. Enabling Robots to Communicate their Objectives. *ArXiv e-prints*, 2017.

[89] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. *arXiv preprint arXiv:1703.06868*, 2017.

[90] Intel. Intel mobileye plans for autonomous fleet. https://newsroom.intel.com/news/intel-mobileye-integration-plans-build-fleet-autonomous-test-cars/, 2017.

[91] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pages 2017–2025, 2015.

[92] J. Janai, F. Güney, A. Behl, and A. Geiger. Computer Vision for Autonomous Vehicles: Problems, Datasets and State-of-the-Art. *ArXiv e-prints*, 2017.

[93] E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.

[94] H. Jégou, M. Douze, C. Schmid, and P. Pérez. Aggregating local descriptors into a compact image representation. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 3304–3311. IEEE, 2010.

[95] S. J. Julier and J. K. Uhlmann. A new extension of the kalman filter to nonlinear systems. In *Int. symp. aerospace/defense sensing, simul. and controls*, volume 3, pages 182–193. Orlando, FL, 1997.

[96] G. Kahn, T. Zhang, S. Levine, and P. Abbeel. Plato: Policy learning using adaptive trajectory optimization. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 3342–3349. IEEE, 2017.

[97] N. Kalchbrenner, A. v. d. Oord, K. Simonyan, I. Danihelka, O. Vinyals, A. Graves, and K. Kavukcuoglu. Video pixel networks. *arXiv preprint arXiv:1610.00527*, 2016.

[98] N. Kant. Recent Advances in Neural Program Synthesis. *ArXiv e-prints*, 2018.

[99] A. Kendall. Reprojection losses: Deep learning surpassing classical geometry in computer vision? https://alexgkendall.com/computer_vision/Reprojection_losses_geometry_computer_vision/, 2017.

[100] A. Kendall and R. Cipolla. Geometric loss functions for camera pose regression with deep learning. *arXiv preprint arXiv:1704.00390*, 2017.

[101] A. Kendall and Y. Gal. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? *ArXiv e-prints*, 2017.

[102] A. Kendall, M. Grimes, and R. Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015.

[103] P. W. Koh and P. Liang. Understanding Black-box Predictions via Influence Functions. *ArXiv e-prints*, 2017.

[104] H. Kretzschmar, M. Spies, C. Sprunk, and W. Burgard. Socially compliant mobile robot navigation via inverse reinforcement learning. *The International Journal of Robotics Research*, page 0278364915619772, 2016.

[105] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, L. Cehovin, G. Fernández, T. Vojir, G. Hager, G. Nebehay, and R. Pflugfelder. The visual object tracking vot2015 challenge results. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 1–23, 2015.

[106] M. e. a. Kristan. The visual object tracking vot2016 challenge results. 2016.

[107] Y. Kuznietsov, J. Stückler, and B. Leibe. Semi-Supervised Deep Learning for Monocular Depth Map Prediction. *ArXiv e-prints*, 2017.

[108] M. Laskey, J. Lee, R. Fox, A. Dragan, and K. Goldberg. Dart: Noise injection for robust imitation learning. In *Conference on Robot Learning*, pages 143–156, 2017.

[109] M. Laskey, S. Staszak, W. Y.-S. Hsieh, J. Mahler, F. T. Pokorny, A. D. Dragan, and K. Goldberg. Shiv: Reducing supervisor burden in dagger using support vectors for efficient learning from demonstrations in high dimensional state spaces. In *Robotics and Automation (ICRA), 2016 IEEE International Conference on*, pages 462–469. IEEE, 2016.

[110] S. M. LaValle. Rapidly-exploring random trees: A new tool for path planning. 1998.

[111] LeCun and Manning. Deep learning, structure and innate priors - a discussion between yann lecun and christopher manning. http://www.abigailsee.com/2018/02/21/deep-learning-structure-and-innate-priors.html, 2018.

[112] Y. LeCun et al. Generalization and network design strategies. *Connectionism in perspective*, pages 143–155, 1989.

[113] Y. LeCun and G. Marcus. Debate: "does ai need more innate machinery?" (yann lecun, gary marcus). https://www.youtube.com/watch?v=vdWPQ6iAkT4, 2017.

[114] S. Liao, Y. Hu, X. Zhu, and S. Z. Li. Person re-identification by local maximal occurrence representation and metric learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.

[115] Z. C. Lipton. The Mythos of Model Interpretability. *ArXiv e-prints*, 2016.

[116] G.-H. Liu, A. Siravuru, S. Prabhakar, M. Veloso, and G. Kantor. Learning End-to-end Multimodal Sensor Policies for Autonomous Navigation. *ArXiv e-prints*, 2017.

[117] S. Lowry, N. Sünderhauf, P. Newman, J. J. Leonard, D. Cox, P. Corke, and M. J. Milford. Visual place recognition: A survey. *IEEE Transactions on Robotics*, 32(1):1–19, 2016.

[118] F. Luan, S. Paris, E. Shechtman, and K. Bala. Deep Photo Style Transfer. *ArXiv e-prints*, 2017.

[119] D. J. MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.

[120] C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.

[121] T. Malisiewicz. The future of real-time slam and deep learning vs slam. http://www.computervisionblog.com/2016/01/why-slam-matters-future-of-real-time.html, 2016.

[122] M. Mathieu, C. Couprie, and Y. LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.

[123] McKinsey. Human + machine: A new era of automation in manufacturing. https://www.mckinsey.com/business-functions/operations/our-insights/human-plus-machine-a-new-era-of-automation-in-manufacturing, 2017.

[124] Microsoft. Microsoft carsim. https://www.microsoft.com/en-us/research/blog/autonomous-car-research/, 2017.

[125] M. J. Milford, G. F. Wyeth, and D. Prasser. Ratslam: a hippocampal model for simultaneous localization and mapping. In *Robotics and Automation, 2004. Proceedings. ICRA'04. 2004 IEEE International Conference on*, volume 1, pages 403–408. IEEE, 2004.

[126] I. Misra, C. L. Zitnick, and M. Hebert. Shuffle and Learn: Unsupervised Learning using Temporal Order Verification. *ArXiv e-prints*, 2016.

[127] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard. Universal adversarial perturbations. *arXiv preprint arXiv:1610.08401*, 2016.

[128] U. Muller, J. Ben, E. Cosatto, B. Flepp, and Y. L. Cun. Off-road obstacle avoidance through end-to-end learning. In *Advances in neural information processing systems*, pages 739–746, 2006.

[129] R. M. Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.

[130] A. Neelakantan, Q. V. Le, and I. Sutskever. Neural Programmer: Inducing Latent Programs with Gradient Descent. *ArXiv e-prints*, 2015.

[131] D. Neven, B. De Brabandere, S. Georgoulis, M. Proesmans, and L. Van Gool. Fast Scene Understanding for Autonomous Driving. *ArXiv e-prints*, 2017.

[132] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng. Multimodal deep learning. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 689–696, 2011.

[133] Q. P. Nguyen, B. K. H. Low, and P. Jaillet. Inverse reinforcement learning with locally consistent reward functions. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 1738–1746. Curran Associates, Inc., 2015.

[134] D. Nistér, O. Naroditsky, and J. Bergen. Visual odometry. In *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, volume 1, pages I–I. Ieee, 2004.

[135] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016.

[136] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:1606.02147*, 2016.

[137] D. Pathak, R. Girshick, P. Dollár, T. Darrell, and B. Hariharan. Learning Features by Watching Objects Move. *ArXiv e-prints*, 2016.

[138] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016.

[139] G. Patrini. In search of the missing signals. http://giorgiopatrini.org/posts/2017/09/06/in-search-of-the-missing-signals/, 2017.

[140] K. Pei, Y. Cao, J. Yang, and S. Jana. Deepxplore: Automated whitebox testing of deep learning systems. *arXiv preprint arXiv:1705.06640*, 2017.

[141] V. Peretroukhin and J. Kelly. DPC-Net: Deep Pose Correction for Visual Localization. *ArXiv e-prints*, 2017.

[142] M. Pfeiffer, G. Paolo, H. Sommer, J. Nieto, R. Siegwart, and C. Cadena. A Data-driven Model for Interaction-aware Pedestrian Motion Prediction in Object Cluttered Environments. *ArXiv e-prints*, Sept. 2017.

[143] M. Pfeiffer, M. Schaeuble, J. Nieto, R. Siegwart, and C. Cadena. From Perception to Decision: A Data-driven Approach to End-to-end Motion Planning for Autonomous Ground Robots. *ArXiv e-prints*, 2016.

[144] M. Pfeiffer, U. Schwesinger, H. Sommer, E. Galceran, and R. Siegwart. Predicting Actions to Act Predictably: Cooperative Partial Motion Planning with Maximum Entropy Models . In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016. arxiv preprint: http://arxiv.org/abs/1607.02329.

[145] S. Pillai and J. J. Leonard. Towards Visual Ego-motion Learning in Robots. *ArXiv e-prints*, 2017.

[146] M. Pivtoraiko and A. Kelly. Efficient constrained path planning via search in state lattices. In *International Symposium on Artificial Intelligence, Robotics, and Automation in Space*, pages 1–7, 2005.

[147] D. A. Pomerleau. Neural network based autonomous navigation. In *Vision and Navigation*, pages 83–93. Springer, 1990.

[148] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation.

[149] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. *ArXiv e-prints*, 2016.

[150] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5105–5114, 2017.

[151] D. Rao, M. D. Deuge, N. Nourani–Vatani, S. B. Williams, and O. Pizarro. Multimodal learning and inference from visual and remotely sensed data. *The International Journal of Robotics Research*, 36(1):24–43, 2017.

[152] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich. Maximum margin planning. In *Proceedings of the 23rd international conference on Machine learning*, pages 729–736. ACM, 2006.

[153] N. D. Ratliff, D. Silver, and J. A. Bagnell. Learning to search: Functional gradient techniques for imitation learning. *Autonomous Robots*, 27(1):25–53, 2009.

[154] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le. Regularized Evolution for Image Classifier Architecture Search. *ArXiv e-prints*, 2018.

[155] J. Redmon and A. Farhadi. YOLO9000: Better, Faster, Stronger. *ArXiv e-prints*, 2016.

[156] J. Ren, X. Chen, J. Liu, W. Sun, J. Pang, Q. Yan, Y.-W. Tai, and L. Xu. Accurate single stage detector using recurrent rolling convolution. *arXiv preprint arXiv:1704.05776*, 2017.

[157] S. R. Richter, Z. Hayder, and V. Koltun. Playing for benchmarks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2213–2222, 2017.

[158] O. Rippel, M. Paluri, P. Dollar, and L. Bourdev. Metric learning with adaptive density discrimination. *arXiv preprint arXiv:1511.05939*, 2015.

[159] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. Lopez. The SYNTHIA Dataset: A large collection of synthetic images for semantic segmentation of urban scenes. 2016.

[160] S. Ross, G. J. Gordon, and J. A. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the International Conference on Artifical Intelligence and Statistics*, 2010.

[161] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.

[162] S. Sabour, N. Frosst, and G. E Hinton. Dynamic Routing Between Capsules. *ArXiv e-prints*, 2017.

[163] D. Sadigh, S. Sastry, S. A. Seshia, and A. D. Dragan. Planning for autonomous cars that leverage effects on human actions. In *Robotics: Science and Systems*, 2016.

[164] B. Salesky. A decade after darpa: Our view on the state of the art in self-driving cars (bryan salesky,argoai). https://medium.com/self-driven/a-decade-after-darpa-our-view-on-the-state-of-the-art-in-self-driving-cars-3e8698e6afe8, 2017.

[165] T. Salimans, J. Ho, X. Chen, and I. Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.

[166] D. Scharstein, R. Szeliski, and H. Hirschmüller. Middlebury stereo vision. http://vision.middlebury.edu/stereo/, 2001.

[167] E. Schmerling, K. Leung, W. Vollprecht, and M. Pavone. Multimodal Probabilistic Model-Based Planning for Human-Robot Interaction. *ArXiv e-prints*, 2017.

[168] A. Seff and J. Xiao. Learning from Maps: Visual Common Sense for Autonomous Driving. *ArXiv e-prints*, 2016.

[169] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization. *ArXiv e-prints*, 2016.

[170] A. Shafaei, J. J. Little, and M. Schmidt. Play and Learn: Using Video Games to Train Computer Vision Models. *ArXiv e-prints*, 2016.

[171] S. Shalev-Shwartz and A. Shashua. On the Sample Complexity of End-to-end Training vs. Semantic Abstraction Training. *ArXiv e-prints*, 2016.

[172] K. Shiarlis, J. Messias, and S. Whiteson. Rapidly exploring learning trees. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 1541–1548. IEEE, 2017.

[173] S. S. Shwartz. The need for human-like dirving (section of video). https://youtu.be/FovLsAFiIJU?t=1m31s, 2016.

[174] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.

[175] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford, and P. Corke. The Limits and Potentials of Deep Learning for Robotics. *ArXiv e-prints*, 2018.

[176] N. Sunderhauf, S. Shirazi, A. Jacobson, F. Dayoub, E. Pepperell, B. Upcroft, and M. Milford. Place recognition with convnet landmarks: Viewpoint-robust, condition-robust, training-free. *Proceedings of Robotics: Science and Systems XII*, 2015.

[177] syntropy.ai. Dimensional reduction via sequential data. https://medium.com/syntropy-ai/dimensional-reduction-via-sequential-data-798d4c3510d9, 2017.

[178] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the Inception Architecture for Computer Vision. *ArXiv e-prints*, 2015.

[179] A. Tamar, Y. Wu, G. Thomas, S. Levine, and P. Abbeel. Value Iteration Networks. *ArXiv e-prints*, 2016.

[180] M. Tanner, S. Saftescu, A. Bewley, and P. Newman. Meshed Up: Learnt Error Correction in 3D Reconstructions. *ArXiv e-prints*, 2018.

[181] K. Tateno, F. Tombari, I. Laina, and N. Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. *arXiv preprint arXiv:1704.03489*, 2017.

[182] the M tank. A year in computer vision. http://www.themtank.org/a-year-in-computer-vision, 2018.

[183] S. Thrun, M. Montemerlo, and A. Aron. Probabilistic terrain analysis for high-speed desert driving. In *Robotics: Science and Systems*, pages 16–19, 2006.

[184] N. Y. Times. Nyt - virtual reality driverless cars. https://www.nytimes.com/2017/10/29/business/virtual-reality-driverless-cars.html, 2017.

[185] G. Tucker, A. Mnih, C. J. Maddison, J. Lawson, and J. Sohl-Dickstein. Rebar: Low-variance, unbiased gradient estimates for discrete latent variable models. In *Advances in Neural Information Processing Systems*, pages 2624–2633, 2017.

[186] M. Turan, Y. Almalioglu, H. Araujo, E. Konukoglu, and M. Sitti. Deep EndoVO: A Recurrent Convolutional Neural Network (RCNN) based Visual Odometry Approach for Endoscopic Capsule Robots. *ArXiv e-prints*, 2017.

[187] E. Tzeng, C. Devin, J. Hoffman, C. Finn, X. Peng, S. Levine, K. Saenko, and T. Darrell. Towards adapting deep visuomotor representations from simulated to real environments. *arXiv preprint arXiv:1511.07111*, 2015.

[188] Uber. Uber atg datavisualisation. https://eng.uber.com/atg-dataviz/, 2017.

[189] Udacity. Udacity selfdriving car project. https://github.com/udacity/self-driving-car, 2017.

[190] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox. DeMoN: Depth and Motion Network for Learning Monocular Stereo. *ArXiv e-prints*, 2016.

[191] UnrealCV. Synthetic computer vision. https://github.com/unrealcv/synthetic-computer-vision, 2017.

[192] A. Valada, G. L. Oliveira, T. Brox, and W. Burgard. Deep multispectral semantic scene understanding of forested environments using multimodal fusion. In *International Symposium on Experimental Robotics*, pages 465–477. Springer, 2016.

[193] J. Valmadre, L. Bertinetto, J. F. Henriques, A. Vedaldi, and P. H. S. Torr. End-to-end representation learning for Correlation Filter based tracking. *ArXiv e-prints*, 2017.

[194] J. Van den Berg, M. Lin, and D. Manocha. Reciprocal velocity obstacles for real-time multi-agent navigation. In *Robotics and Automation, 2008. ICRA 2008. IEEE International Conference on*, pages 1928–1935. IEEE, 2008.

[195] S. Vijayanarasimhan, S. Ricco, C. Schmid, R. Sukthankar, and K. Fragkiadaki. SfM-Net: Learning of Structure and Motion from Video. *ArXiv e-prints*, 2017.

[196] C. Vondrick, H. Pirsiavash, and A. Torralba. Anticipating visual representations from unlabeled video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 98–106, 2016.

[197] C. Vondrick and A. Torralba. Generating the future with adversarial transformers. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[198] E. A. Wan and R. Van Der Merwe. The unscented kalman filter for nonlinear estimation. In *Adaptive Systems for Signal Processing, Communications, and Control Symposium 2000. AS-SPCC. The IEEE 2000*, pages 153–158. Ieee, 2000.

[199] J. X. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Z. Leibo, R. Munos, C. Blundell, D. Kumaran, and M. Botvinick. Learning to reinforcement learn. *ArXiv e-prints*, 2016.

[200] S. Wang, R. Clark, H. Wen, and N. Trigoni. DeepVO: Towards End-to-End Visual Odometry with Deep Recurrent Convolutional Neural Networks. *ArXiv e-prints*, 2017.

[201] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs. *ArXiv e-prints*, 2017.

[202] W. Wang, A. Wang, A. Tamar, X. Chen, and P. Abbeel. Safer Classification by Synthesis. *ArXiv e-prints*, 2017.

[203] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas. Dueling Network Architectures for Deep Reinforcement Learning. *ArXiv e-prints*, Nov. 2015.

[204] Waymo. Inside waymo's secrect testing and simulation facilities. https://www.theatlantic.com/technology/archive/2017/08/inside-waymos-secret-testing-and-simulation-facilities/537648/, 2017.

[205] Waymo. Waymo safety report. https://storage.googleapis.com/sdc-prod/v1/safety-report/waymo-safety-report-2017-10.pdf?utm_source=The+Comet+Newsletter, 2017.

[206] T. Weber, S. Racanière, D. P. Reichert, L. Buesing, A. Guez, D. Jimenez Rezende, A. Puig-domènech Badia, O. Vinyals, N. Heess, Y. Li, R. Pascanu, P. Battaglia, D. Silver, and D. Wierstra. Imagination-Augmented Agents for Deep Reinforcement Learning. *ArXiv e-prints*, 2017.

[207] Y. Z. S. Q. Z. X. T. S. K. Y. W. A. Y. Weichao Qiu, Fangwei Zhong. Unrealcv: Virtual worlds for computer vision. *ACM Multimedia Open Source Software Competition*, 2017.

[208] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

[209] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, 2016.

[210] M. Wulfmeier, A. Bewley, and I. Posner. Addressing appearance change in outdoor robotics with adversarial domain adaptation. In *Proceedings of the IEEE International Conference on Intelligent Robots and Systems*, 2017.

[211] M. Wulfmeier, A. Bewley, and I. Posner. Incremental Adversarial Domain Adaptation for Continually Changing Environments. *ArXiv e-prints*, 2017.

[212] M. Wulfmeier, I. Posner, and P. Abbeel. Mutual Alignment Transfer Learning. *ArXiv e-prints*, 2017.

[213] M. Wulfmeier, D. Rao, D. Z. Wang, P. Ondruska, and I. Posner. Large-scale cost function learning for path planning using deep inverse reinforcement learning. *The International Journal of Robotics Research*, 36(10):1073–1087, 2017.

[214] xix.ai. Adversarial attacks. https://blog.xix.ai/how-adversarial-attacks-work-87495b81da2d, 2017.

[215] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014.

[216] A. R. Zamir, T. Wekel, P. Agrawal, C. Wei, J. Malik, and S. Savarese. Generic 3d representation via pose estimation and matching. In *European Conference on Computer Vision*, pages 535–553. Springer, 2016.

[217] J. Zhang, L. Tai, J. Boedecker, W. Burgard, and M. Liu. Neural SLAM. *ArXiv e-prints*, 2017.

[218] J. Zhu. Image gradient-based joint direct visual odometry for stereo camera. In *Int. Jt. Conf. Artif. Intell*, pages 4558–4564, 2017.

[219] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. *ArXiv e-prints*, 2017.

[220] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey. Maximum entropy inverse reinforcement learning. In *AAAI*, pages 1433–1438, 2008.

[221] B. D. Ziebart, N. Ratliff, G. Gallagher, C. Mertz, K. Peterson, J. A. Bagnell, M. Hebert, A. K. Dey, and S. Srinivasa. Planning-based Prediction for Pedestrians. In *Proceedings of the 2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, IROS'09, pages 3931–3936, Piscataway, NJ, USA, 2009. IEEE Press.