

DROW: Real-Time Deep Learning based Wheelchair Detection in 2D Range Data

Lucas Beyer[†], Alexander Hermans[†] and Bastian Leibe

Abstract— We introduce the DROW detector, a deep learning based object detector operating on 2D range data. Laser scanners are lighting invariant, provide accurate 2D range data, and typically cover a large field of view, making them interesting sensors for robotics applications. So far, research on detection in laser 2D range data has been dominated by hand-crafted features and boosted classifiers, potentially losing performance due to suboptimal design choices. We propose a Convolutional Neural Network (CNN) based detector for this task. We show how to effectively apply CNNs for detection in 2D range data, and propose a depth preprocessing step and a voting scheme that significantly improve CNN performance. We demonstrate our approach on wheelchairs and walkers, obtaining state of the art detection results. Apart from the training data, none of our design choices limits the detector to these two classes, though. We provide a ROS node for our detector and release our dataset containing 464k laser scans, out of which 24k were annotated.

I. INTRODUCTION

Many autonomous robots are equipped with a 2D laser scanner, typically used for navigation-related tasks including the detection of people [1], [34] and objects [18]. Laser scanners are widely used due to their typically large field of view and their invariance to lighting and environmental conditions. While early detection methods used simple heuristics such as fitting lines and circles [35], in the past few years hand crafted features, coupled with learned classifiers, have dominated laser based detection. Within this paradigm, a range of successful people [1], [30], [20], mobility aid [34] and road obstacle [18] detectors have been developed to support mobile robot navigation [8] and autonomous driving [28]. Even though the aforementioned models obtain respectable results, the general consensus seems to be that the information provided by 2D range data is not sufficient to reliably perform detection in a single scan, leading to approaches relying on sensor fusion [30], multi-layered sensor setups [20], [29], or temporal integration of information by tracking. We show that detection based on single 2D range scan can actually perform well.

In this paper, we specifically focus on the detection of wheelchairs and walkers. This is motivated by a service robot application in an elderly care facility, where many people rely on their mobility aids. Since the presence of a mobility aid can significantly change a person’s appearance, both in laser and in camera data, they are less reliably detected by existing people detectors. However, especially people relying

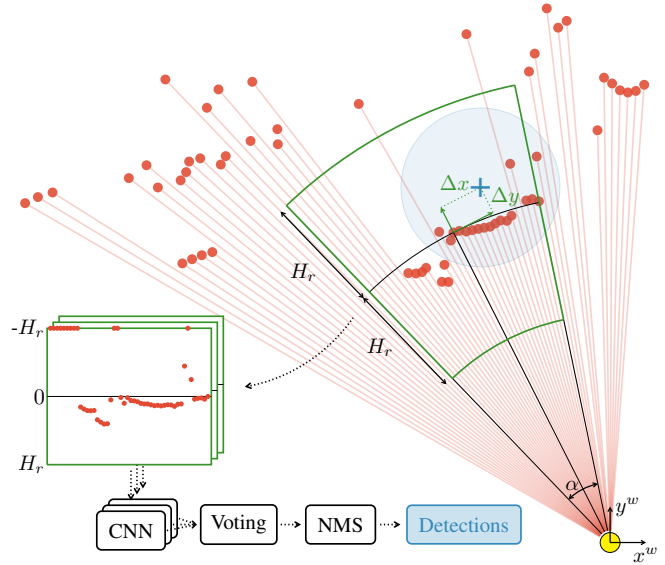


Fig. 1: An overview of our approach. We extract windows with a fixed real world size regardless of the laser coverage. Values are centered and clamped to the range of the window. Each window is classified by a CNN and, if an object is found, a weighted vote is cast for its relative offset $(\Delta x, \Delta y)$. Finally, the votes are consolidated using a non-maximum suppression (NMS) scheme to yield detection centroids.

on those aids will have a harder time avoiding an approaching robot, making a reliable detection all the more important. Weinrich *et al.* [34] face a similar task and propose the Gandalf detector for detecting people, wheelchairs, and walkers in individual laser scans. They introduce a new feature set for laser segments and classify these with an AdaBoost classifier. The resulting detector performs reasonably well and we include it as a baseline.

In this paper, we introduce the DROW (Distance ROBust Wheelchair/Walker) detector. Driven by our application scenario, we focus on the detection of wheelchairs and walkers. However, there are no design decisions which restrict our approach to these objects and we are confident it will generalize to detection of persons or other objects, given sufficient annotated training data. Code and data needed to train and run our detector, as well as code to annotate data for new tasks, will be made available upon publication.

In computer vision, deep learning has recently become the new best practice, replacing hand-crafted features by learned ones and overhauling the state of the art in many tasks [13], [32], [5]. Specifically, Convolutional Neural Networks (CNNs) have been very successful at challenging

[†]Equal contribution. Ordering determined by a last minute coin flip.

The work in this paper was funded by the EU project STRANDS (ICT-2011-600623). All authors are at the Visual Computing Institute, RWTH Aachen University. e-mail: last@vision.rwth-aachen.de

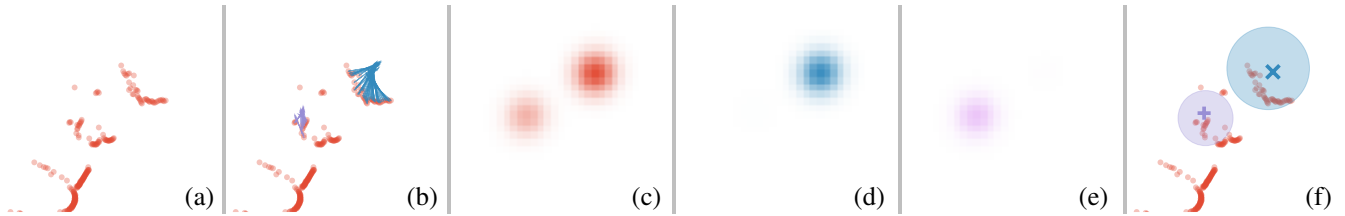


Fig. 2: Non-maximum suppression: (a) a part of an input scan, (b) the votes shown as arrows, (c) the joint voting grid, (d) the wheelchair voting grid, (e) the walker voting grid, (f) resulting detections.

tasks. In this paper, we show how CNNs can be applied for object detection in laser data, alleviating the need for feature engineering and enabling drastic improvements.

While CNN-based image-level detectors like Multi-Box [16] and YOLO [24] could in principle be applied to a 2D laser scan, we found that doing so naively is not effective. Since the spatial density of laser points varies greatly with distance, the fixed perceptive field of a CNNs covers widely different scales, making learning difficult. To make use of the spatial information a laser sensor provides, we propose a preprocessing stage that cuts out and normalizes a fixed real-world extent window around each laser point. All of those windows are fed through a CNN which can cast votes for object locations. These votes are then turned into individual detections using a non-maximum suppression scheme. We show that both depth preprocessing and voting are essential components of our approach, an overview of which can be seen in Fig. 1. Our approach does not require background subtraction and runs with a frame rate of ~ 75 fps on a modern desktop machine using a single core and a GPU. On our robot, it easily keeps up with the laser frame rate (~ 13 fps) on a laptop GPU.

To summarize, we make the following contributions:

- We introduce the DROW detector, a CNN based wheelchair and walker detector for 2D range data which, by effectively making use of the provided depth-information, achieves state of the art results.
- We publish our dataset, containing 464k raw scans of which 24k have been annotated with wheelchair and walker centroids.
- We provide ROS components of our detector and related service modules, including trained models.

II. APPROACH

Our approach consists of three steps: preprocessing, which cuts out a resampled window around every laser point and computes detection locations in a local coordinate system, a CNN classifying said windows and predicting relative detection locations, and finally a voting and non-maximum suppression scheme turning predictions into detections.

A. 2D Range Data Preprocessing

In order to apply a CNN for detection, the network’s receptive field must cover a large part of the object. The problem with laser scans is that nearby objects cover a large amount of laser beams, whereas distant objects are only hit by a handful of beams. This means the CNN’s receptive field

must cover most of the laser, which makes it very prone to overfitting to the training scenes’ backgrounds.

To circumvent this problem, and at the same time make use of the real-world scale information that laser data provides, we propose to evaluate the CNN in a depth-guided sliding-window fashion. This means that we preprocess the data such that objects have approximately the same representation at every distance, thus alleviating the need to implicitly learn completely different representations at varying distances: Around each beam, we cut out a window of real-world extent ℓ , thus spanning an angle $\alpha = 2 \sin^{-1}(\frac{\ell}{2r})$ and containing a variable amount of measurements depending on the distance r at which the current beam hits an obstacle. We then resample the measurements inside this window linearly at n fixed samples. When applied to such a window, the network’s receptive field always covers the same real-world extent, regardless of the distance r . In addition, we center the window around the current point, clamp any values outside a $\pm H_r$ hull in order to remove distant clutter and finally project the values into $[-1, 1]$. In total, this means $\max(-H_r, \min(x - r, H_r))/H_r$ is applied to each point x of the window around the point at depth r , as illustrated in Fig. 1. A detailed analysis regarding which of these preprocessing operations contribute the most to DROW’s performance can be found in Section III-E.

B. Prediction

For each window, and thus for each laser point with its context, a CNN both classifies whether this window belongs to an object class of interest through a SoftMax output and, if so, votes for the center location of that object through a regression output. As the 2D range data is inherently rotation-invariant, we do not want to perform voting in absolute (x, y) coordinates. Instead, we learn offsets $(\Delta x, \Delta y)$ in a coordinate system centered and aligned at the current laser point, as shown in Fig. 1.

C. Voting and Non-Maximum Suppression

The predictions for every window need to be consolidated into detection centers. This is achieved by making the CNN’s predictions vote into regular grids spanning the laser’s field of view. Let $p(O|w) = \sum_{c \in \mathcal{C}} p(c|w)$ be the total probability of the window w seeing an object of interest, where \mathcal{C} are the classes to be detected.

If $p(O|w)$ exceeds a predefined voting threshold T , the window will cast a vote into a class-agnostic grid with weight $p(O|w)$ as well as into each class-specific grid with weight $p(c|w)$. After all windows have potentially cast votes, each

grid is blurred with a Gaussian filter and non-maximum suppression is performed on the class-agnostic grid. For each maximum found in that grid, a detection is predicted at the cell’s center using the class which has the highest sum of votes in said cell. The reason for this voting scheme, as opposed to treating each class separately, is to avoid detections of both classes at the same position. Fig. 2 shows an example of the different steps involved. Votes cast from the raw laser points in (a) are shown in (b); (c) - (e) show the three voting grids and (f) shows the two resulting detections.

III. EXPERIMENTAL EVALUATION

For the evaluation of our approach we first introduce our new dataset. We outline the details of our training procedure and of the evaluation methods. After the general evaluation of DROW and a comparison to baselines, we perform several ablation studies to show how much individual parts of our detector contribute to the overall performance.

A. Dataset

Although Weinrich *et al.* [34] provide their datasets, their recordings are limited to scenes with a single wheelchair or walker. While this might suffice for learning the few parameters of fairly constrained features and classifiers, we want to learn features from scratch and thus require more general and varied data. This is why we recorded a little over 10 hours of data at an elderly care facility.

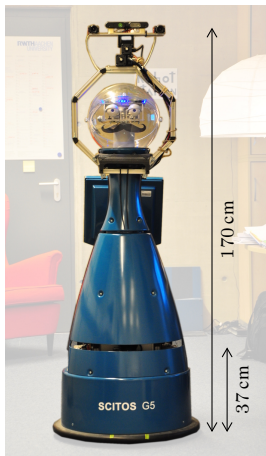


Fig. 3: Our robot Karl.

The software infrastructure is based on ROS and all data was stored in rosbags. To ensure seeing enough wheelchairs and walkers, one person permanently roamed around using either while the robot was recording. Our recordings consist of both (1) natural scenes where we recorded the everyday life in the facility including all kinds of clutter such as flower pots, chairs, furniture, rolling beds, . . . , and (2) artificial recordings where we drove around in certain patterns as to bring variation to the dataset. Apart from the resident’s wheelchairs, we included as many wheelchair and walker models as possible, including one motorized wheelchair.

1) *Recording Setup:* We recorded the data using a SCITOS G5 robot equipped with a SICK S300 laser scanner mounted at approximately 37 cm above the ground plane. The laser was configured to record at nearly 13 Hz and span a field of view of 225° at a resolution of 0.5° , totalling 450 measurements. For annotation purposes, we also recorded RGB-D data from an ASUS Xtion mounted on the robot’s head (see Fig. 3). Unfortunately, we are not allowed to publish the video streams for privacy

TABLE I: Dataset overview

	Train	Validation	Test	Total
Sequences	78	30	5	113
Scans	341 138	74 744	48 131	464 013
Annotated Scans	17 665	3 919	2 428	24 012
Wheelchairs	14 455	5 595	1 970	22 020
Walkers	2 047	219	581	2 847
Wheelchairs by distance				
Walkers by distance				

2) *Dataset Statistics:* We split the resulting dataset into a train, validation and test set. To create these sets, we split the care facility into four non-overlapping areas, three of which were assigned to the train, test and validation sets, respectively. The fourth, the entrance hall, was split into temporally disjoint sequences which were distributed over the train and validation sets. Based on this split, we can show how well the approach generalizes to never before seen areas. Table I shows an overview of our dataset, as well as statistics of the subset we annotated. The wheelchair and walker counts refer to individual detections as opposed to instances, and the bar plots show their distribution over the distance. Each bar represents a 1 m slice in the distance (15 bars for up to 15 m), clearly showing that the majority of observed mobility aids are encountered within 1 m to 6 m of the robot.

3) *Annotation:* We annotate our data with wheelchair and walker centers. Since the dataset contains 464k raw laser scans, we devise an annotation scheme that keeps effort manageable while still covering the full extent of the sequences as well as allowing temporal approaches to be developed in the future. Instead of all the laser scans, we annotate small batches throughout every sequence as follows: A batch consists of 100 frames, out of which we annotate every 5th frame, resulting in 20 annotated frames per batch.¹ Within a sequence, we only annotate every 4th batch, leading to a total of 5% of the laser scans being annotated. We wrote a custom annotation tool based on matplotlib [11], which

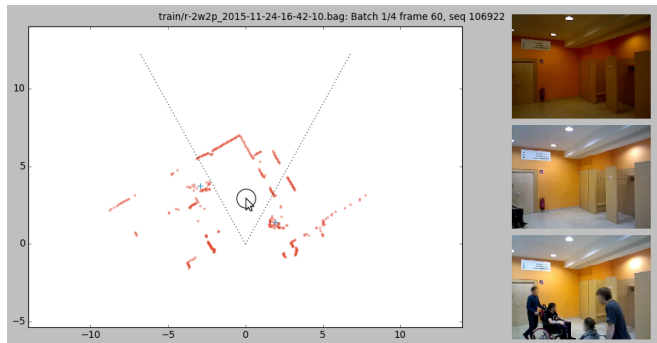


Fig. 4: An example view of our annotation tool. The dotted cone shows the Xtion’s field of view and on the right the first, current, and last image in the batch is shown. The blue cross shows annotated wheelchairs.

¹This allows experimenting with interpolation between the annotations, even though we didn’t do so in this work.

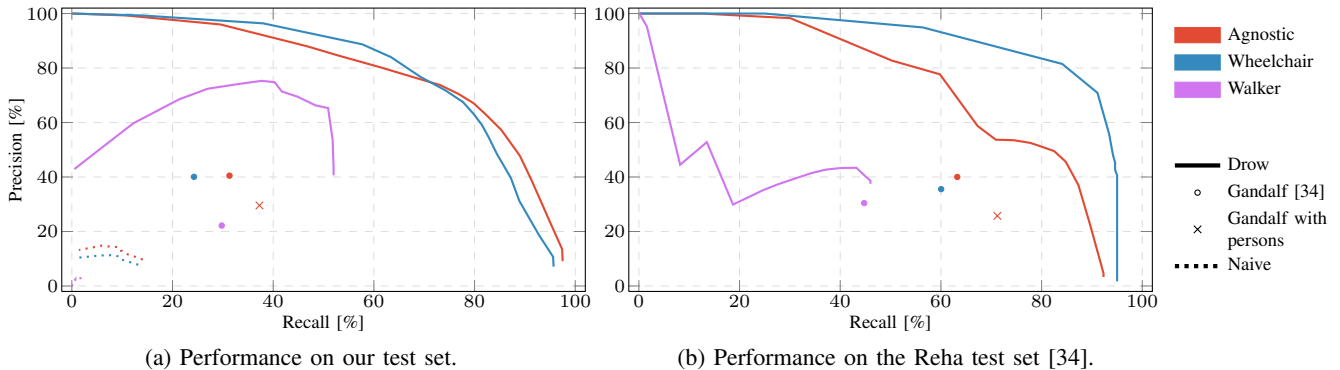


Fig. 5: Performance comparison of DROW, the Gandalf detector [34], and a naive deep learning baseline on two test sets (a) and (b). As explained in Section III-D, the naive baseline is not applicable on the Reha test set.

loads sequences of scans with corresponding RGB images from the head camera and automatically finds the batches that need to be annotated. To aid the user in annotation, we show the first, current, and last image of the current batch. An example view of this annotation tool can be seen in Fig. 4. A 1.2m circle around the mouse pointer, indicating the average wheelchair size, helps the user click on locations of wheelchair or walker centers. By using all of this supportive information to get an understanding of the scene, we were able to annotate most mobility aids, but based on the limited camera view, we likely still missed a few.

B. Training Procedure

We train a CNN which, for each preprocessed window, predicts whether it is indicative of a nearby wheelchair or walker and, if it is, predicts the offset of its center. We do so in a single pass by using a network with two output layers: a three-way SoftMax differentiating between background, wheelchair and walker, and a two-dimensional linear regression output. We optimize the network by minimizing the sum of a negative log-likelihood criterion on the SoftMax output and a root-mean-square error on the regression output. The regression targets are computed as $(\Delta x, \Delta y)$ in each window’s local coordinate system as shown in Fig. 1. The class-labels are determined by the type of the closest detection to the window’s center-point, with a maximal Euclidean distance of 0.6m for wheelchairs and 0.4m for walkers. When there is no nearby annotation, no error is backpropagated for the regression output and the network is thus free to predict any offset.

The architecture of our network, inspired by the popular VGGnet [27]², is as follows: Conv 5@64, Conv 5@64, Max 2, Conv 5@128, Conv 3@128, Max 2, Conv 5@256, Conv 3@5, where Conv $n@c$ represents a convolution with c filters of size n and Max p the maximum operation on a window of p values. Batch-normalization [12], dropout [31] of 0.25, and ReLU nonlinearities [9] are applied between all convolutional layers. For an input window resampled to 48 values, the network outputs a vector of length five, three of which are sent to a SoftMax and the other two are the

regression outputs Δx and Δy . All weights are initialized similarly to the scheme proposed in [26], except for those of the output layer, which are initialized to 0. For training the network, we use the AdaDelta optimizer with $\rho = 0.95$ and $\epsilon = 10^{-7}$ and train until approximate convergence.³ During training, we add small multiplicative random noise to the regression targets and we flip each window and its target with probability 0.5. We implemented our CNN in Theano [2].

For the voting scheme, we add multiplicative class-weights to the votes and re-normalize them, thus voting for class c with $\frac{w_c p(c|w)}{\sum_{i \in C} w_i p(i|w)}$ instead of $p(c|w)$. We then optimize all voting hyper-parameters (class-weights w , grid resolution b and blur-size σ) using hyperopt [4] to maximize $\max_T f_{\text{wheelchair}}(T) + f_{\text{walker}}(T)$ on our validation set, $f_c(T)$ being the F1-score of class c using detection threshold T . Interestingly, the best values, $w_{\text{bg}} = 0.38$, $w_{\text{wheelchair}} = 0.60$, $w_{\text{walker}} = 0.49$, $b = 0.1$ m and $\sigma = 2.93$ are not far off from our initial guesses.

C. Approach Evaluation

We trained our model on our training set and computed hyperparameters on our validation set as described in the previous section. In order to evaluate real-world performance of DROW, we now look at the precision and recall curves on our own test set and on the publicly available Reha test set of [34]⁴ recorded with a similar robot. Recall that our test set was recorded in a never before seen part of the care facility, thus a method which learns location priors or background models will fail. The result shown in Fig. 5 demonstrates that DROW generalizes very well and has significantly higher precision and recall than Gandalf [34], both on our test set and on the Reha test set.

In our application scenario, the actual detection of a mobility aid is more important than their correct classification, hence we also evaluate a class-agnostic performance. For this, we ignore the class of all detections and annotations during evaluation.

³This means that we started training in the evening and stopped it the next morning. The loss-curve was almost flat in log-space.

⁴Detection annotations for the Reha dataset were provided by the authors. We skip the Home test set as it does not contain mobility aids.

²VGG is a well-performing architecture in vision. Other architectures may perform better, but that’s beyond the scope of this paper.

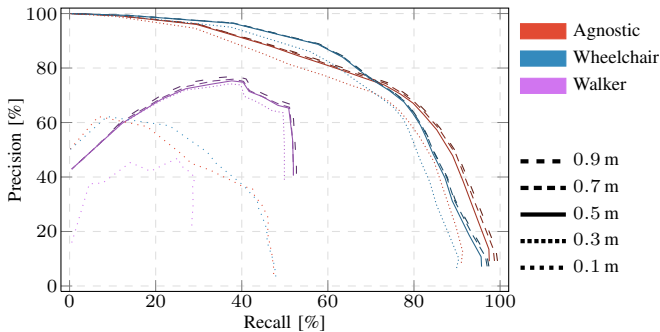


Fig. 6: Performance of DROW for different evaluation radii. The high overlap for radii above 0.3 shows that our predictions are well localized.

The precision and recall curves for each class are computed by varying the voting threshold T . We use an evaluation radius of 0.5m, meaning a detection is matched to a ground-truth if it falls within 0.5m of that ground-truth’s center and has the correct class. An annotation can be matched with at most one detection, all remaining detections of that class are false positives and all unmatched annotations of that class are false negatives. To see how well localized DROW’s detections are, we show precision-recall curves for varying evaluation radii in Fig.6. The cluttering of all curves for acceptance radii of 0.3m and above mean that our detections are well localized. During annotation, it became clear to us that the precision of the annotations is limited, meaning that evaluation at radii as small as 0.1m is strongly affected by labeling noise.

We also analyze how well DROW behaves with respect to the distance from the laser scanner. For this, we start by ignoring all detections beyond 0.1m from the laser. We then slowly grow that radius, taking into account more and more detections, and plot how the precision and recall at a fixed threshold evolve in Fig.7. DROW performs very well across the board, but especially so in mid-range which is crucial for navigation and planning. Unsurprisingly, the curves do not change much beyond a distance of 10m, as only few mobility aids were observed that far.

D. Baselines

We compare our proposed method to two baselines in Fig.5: the publicly available Gandalf detector [34] and a naive deep learning baseline.

None of the precision-recall curves in [34] quantify Gandalf’s actual detection performance, they only focus on parts of the system independently. To obtain comparable detection precision-recall values, we evaluate Gandalf⁵ using the evaluation protocol described above. As provided, the Gandalf detector has no single threshold to be tuned and thus results in a single point in our plots. For the class-agnostic case, we plot the results for the cases when Gandalf person detections are kept (\bullet) and discarded (\times), respectively, showing a trade-off between precision and recall.

⁵https://github.com/neurob/gandalf_detector

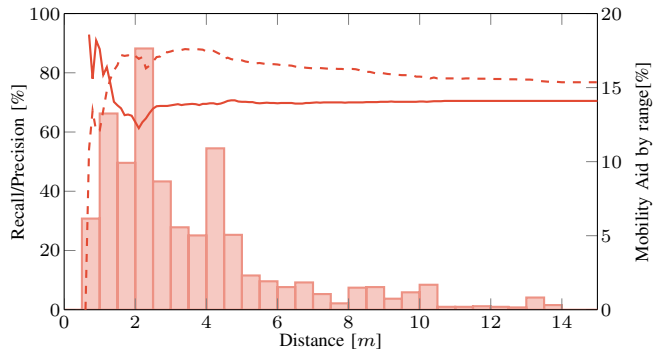


Fig. 7: Precision (—) and Recall (---) at $T = 0.5$ in the class-agnostic case up to a certain distance in meters. The histogram (■, y axis on the right) shows the percentage of all annotations at a certain range.

Note that the detector performance we obtained with the code from [34] is significantly lower than the one reported in the original paper. The performance could be improved by an additional covariance-based merging step briefly mentioned in [34] which is, however, neither described in detail in the paper, nor included in the provided detector code. We were thus not able to reproduce the results presented in [34]. However, extrapolating from the false positives and missing detections we observed using the provided code, even an optimistic bound on this method’s performance would place it below our detector’s curve.

As a naive deep learning baseline, we train another CNN, similarly to YOLO [24], that directly predicts up to two detections (sufficient for 95.0% of the frames) in normalized (x, y) coordinates, based on a full scan. The results of this baseline experiment are shown as dotted lines in Fig.5 (a), but are missing from the Reha test set as the scans have a different amount of laser beams. It should be noted that we spent a considerable amount of time getting the naive CNN baseline to work as well as possible and follow best-practice: carefully chosen receptive-fields, batch-normalization, careful initialization, AdaDelta optimizer, etc. However, as can easily be seen, just applying deep learning to 2D range data in such a naive way leads to abysmal results.

E. Ablation Studies

In this section, we analyze how each of our design decisions affects detection by systematically removing or substituting each part that defines our approach. Fig.8 summarizes these experiments. Overall, the complete DROW detector (—) outperforms all other configurations, *i.e.* each single decision contributes to the high performance. For each of the following experiments, we retrain the full network from scratch.

1) *Preprocessing*: Our depth preprocessing can be dissected into three operations: centering, clamping and resampling. Not centering the input window (—) has the smallest, albeit non-negligible, effect. Removing the clamping (—) hurts overall performance as much as removing the centering. Without the resampling step (—) which ensures that the network’s receptive field keeps a constant real-world

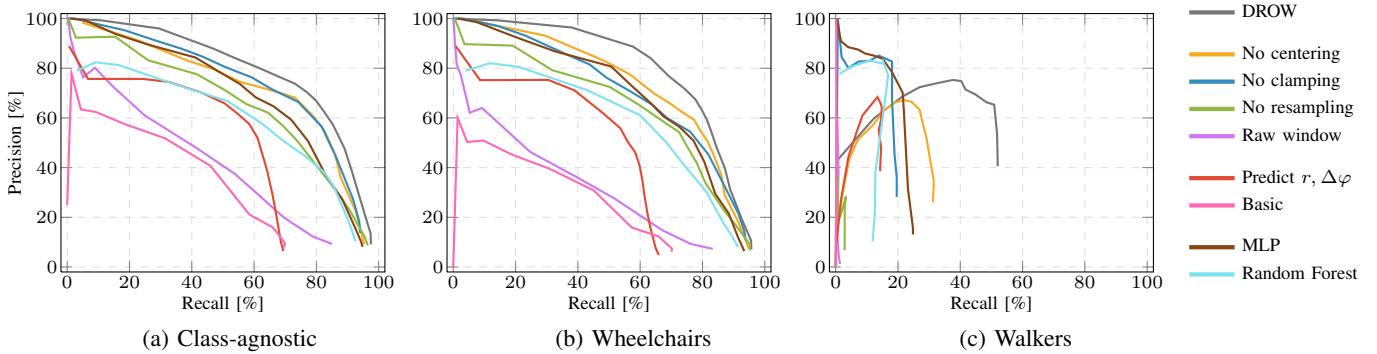


Fig. 8: Ablation studies. We retrain our approach with slight modifications to show how the performance changes.

size, performance drops significantly, especially for walkers. Considering that the CNN has to learn completely different representations across different distances, the overall performance is still surprisingly good. The fact that walkers are almost undetected and yet the agnostic performance is decent suggests a high confusion. Finally, dropping all of the above steps (—) performs the worst in all measures, showing that our preprocessing is indeed crucial.

In total, all these experiments clearly suggest that each of the three preprocessing steps are important.

2) *Voting*: The comparison to the naive deep learning baseline in Fig. 5a already showed that voting is essential for proper detection results.

In order to see the effect of regressing local offsets ($\Delta x, \Delta y$), we train a version of the network which regresses ($\Delta \varphi, r$): an angular offset and an absolute distance to the laser scanner (—). For this experiment, we also had to remove centering (*c.f.* —), as otherwise r is impossible to predict. Without centering, the network could in principle learn to “pass through” the centerpoint’s distance and add it as a bias in the last layer. The drastically degraded results suggest that this is a much more difficult problem that the learning procedure fails to solve.

To ensure that this difficulty is not caused by a bad interaction with our preprocessing, we also train a network to make these predictions on raw, unprocessed data (—). This model performing worst of all shows that even a plain voting-based network alone does not perform well if no care is taken on both the input and output.

3) *Model*: Since every window is sent through the model individually, and all windows have the same size due to the resampling step, there is no inherent reason for the model to

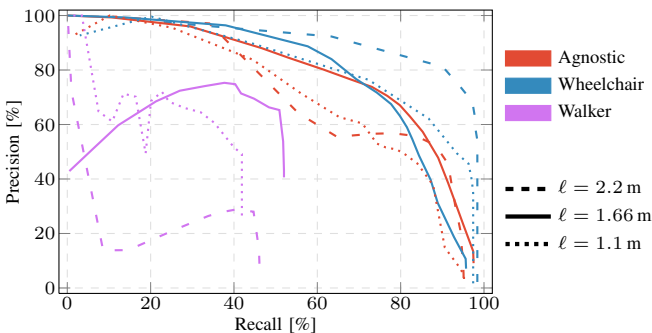


Fig. 9: Effect of the window size ℓ on the performance.

be a CNN. To show that the spatial prior encoded in a CNN is useful, we train a three hidden-layer perceptron (—) with 2048 hidden units each, ReLU non-linearities, dropout, and batch-normalization. As in the other experiments, we leave all other design decisions unchanged. The result shows that, although the MLP is a viable alternative, the CNN clearly outperforms it.

As an additional baseline we trained a regression forest (—) that regresses the three class probabilities and the local offsets ($\Delta x, \Delta y$). It was trained with the same training data as the CNN and the MLP (including the flip augmentation). We used the scikit-learn [23] implementation of the regression forest with all default settings and 50 trees. Training took 7.5 hours (using 5 threads) and resulted in a model of 6.9GB compared to our CNN model of 1.1MB. Nevertheless, the forest model performs even slightly worse than the MLP in all cases.

4) *Window Size*: So far, all evaluations have been performed with a real-world window spanning $\ell = 1.66$ m, motivated by the typical extent of wheelchairs being 1.2 m. To verify that this is a reasonable choice, we trained a network on both larger (2.20 m) and smaller (1.10 m) windows. The results of this experiment can be seen in Fig. 9. The class-agnostic prediction performance is best for our choice of $\ell = 1.66$ m, but wheelchairs alone could be detected considerably better by using a window of 2.2 m. These results suggest that further improvements might be obtained by providing the network with multiple scales of input windows.

IV. RELATED WORK

Most related to our approach is that of Weinrich *et al.* [34] who propose a distance invariant feature for laser based detection. They detect wheelchairs, walkers, and persons in 2D range data by first detecting larger jumps in the distance. Once such a jump is found, they create a window with a fixed real-world extent, covering subsequent scan points. The window is discretized into equally sized segments, for each of which a clamped distance relative to the window depth is computed. Each segment is characterized by the min, max and average depth which are concatenated for all segments to form the feature vector. These are then classified by an AdaBoost classifier. While their feature design has a high similarity to our depth preprocessing step, their windows are only created when a jump in distance is found, where we

create a window for every scan point. The biggest difference, however, is that they directly predict a detection center and object class for every window, while we vote for centers, a step in our detector which we have shown to be essential.

A few more publications [7], [14], [17] identify the need for a distance-invariant representation which we highlight in this work. The solution of creating a re-sampled depth image from a 3D point-cloud by ray-tracing [7], however, involves multiple significantly more complex operations than our proposed simple pre-processing and no timings were provided. The other commonly-used solution is the creation of a 3D occupancy-grid [14], [17]. Such an approach effectively uses $N + 1$ -dimensional inputs for data on an N -dimensional manifold, while our approach keeps the input N -dimensional. Arguably, the latter makes learning easier and predictions more robust [3]. In addition, multiple different heuristics exist for “filling holes” in those occupancy-grids [17], whereas our pre-processing doesn’t produce holes in the first place. In the end, both data representations may be made to work, but we believe our lower-dimensional representation is simpler and more effective.

Others have used voting for detection in laser scans, Mozos *et al.* [20] in a multi-height laser setup and Spinello *et al.* [29] in a layered fashion based on 3D range data. Both, related to [15], learn shape models from data and cast votes from the different laser layers to detection centroids. However, similar to Weinrich *et al.* [34], they rely on jump distance segmentation and only the segments can cast votes for detections. In [33], Wang *et al.* prove that detection using a voting scheme corresponds to sliding-window detection on a sparse grid for linear models. While they achieve good detections on point clouds using hand-crafted features and a linear SVM, it has been shown time and again in the computer vision literature (and in this paper) that learned non-linear features and classifiers outperform hand-crafted features and linear classifiers by a large margin. In line with our findings, voting has recently been shown to work well in combination with various other deep learning approaches [19], [25], [36]. It would thus be interesting to establish a relationship between deep, non-linear voting detectors such as DROW and sliding-window detectors on sparse inputs as shown in [33].

Detection in range data is not limited to persons or mobility aids. Around the time at which we uploaded the preprint⁶ of this paper, Ondruska *et al.* [22] shows how to do “tracking” of pedestrians, cyclists, buses, cars, and road obstacles in 2D range data. Although seemingly similar to our work, their input is an occupancy-grid (i.e. $N + 1$ -dimensional) for which they predict a labelled occupancy-grid, as opposed to discrete detections or tracks with IDs. Based on the biases for static grid cells in their RNN, it remains to be seen how well their approach works on a mobile robot.

Merdrignac *et al.* [18] also detect cars, bicyclists and static road obstacles in 2D range data. They hypothesize that more

information can be extracted from range data than done before and aim to achieve this by designing a large set of hand-crafted features. We agree, but we instead learn feature representations from the data directly.

V. DISCUSSION

A few interesting aspects of our detector should be highlighted. Firstly, given the experimental evaluation it becomes clear that we perform well with respect to wheelchairs, but our walker performance leaves some room for improvement. One reason for this could be the fact that our training set is rather biased towards wheelchairs. However, the class-agnostic precision-recall curves are very similar to the ones of wheelchairs in most of our experiments. This suggests that one significant problem for walker detection is a high confusion with wheelchairs. The fact that whenever DROW detects a mobility aid of either class, it is well localized, leads us to the conclusion that the network “understood” walkers but needs to see more of them.

Secondly, laser-based detection is often dismissed as too sparse or too difficult in a single frame. We hope our results refute this fear, even though we do not dismiss the performance gain that could likely be achieved by tracking our detections over time.

A commonly suggested alternative to laser-based detections is the use of a vision-based detector and RGB(-D) cameras, which produce much richer data. Several approaches have tried this before [21], [10], [6], but none of them is general enough to be applied in all scenarios. They typically rely on geometric primitives [21] or very specific camera setups [6]. Furthermore, commonly used RGB-D sensors such as the Asus Xtion, mounted on a human sized robot, have a too narrow field of view to perceive wheelchairs both far away, as well as when they come close to the robot. Vision-based detectors would thus need many cameras to cover the same field of view that is covered by a single laser scanner. Additionally, since mobility aids themselves come in many different appearances, training such a detector robustly would require a vastly larger dataset. Thus, a laser-based detector, as we have shown to be feasible here, is the more effective and efficient approach.

VI. CONCLUSIONS

In this paper we introduced the DROW detector, a fast deep learning based detector for wheelchairs and walkers from 2D range data. We have proposed a depth preprocessing and a voting scheme, both of which enable CNNs to vastly outperform naive CNN detection baselines and obtain state of the art results compared to a previous method. We performed a thorough experimental evaluation, justifying all our major design choices. To achieve all this, we recorded a large dataset in which we annotated wheelchair and walker centroids, and which we believe will enable further research. We are convinced that our detector generalizes to other classes given training data and will thus be useful to the community. Upon acceptance of the paper, we will publish our code, including a ROS node and the annotated dataset.

⁶<https://arxiv.org/abs/1603.02636v1>

REFERENCES

- [1] Kai O Arras, Óscar Martínez Mozos, and Wolfram Burgard. Using Boosted Features for the detection of People in 2D Range Data. In *ICRA*, 2007.
- [2] Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, James Bergstra, Ian J. Goodfellow, Arnaud Bergeron, Nicolas Bouchard, and Yoshua Bengio. Theano: new features and speed improvements. Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop, 2012.
- [3] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *PAMI*, 2013.
- [4] James Bergstra, Daniel Yamins, and David Daniel Cox. Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. In *ICML*, 2013.
- [5] Lucas Beyer, Alexander Hermans, and Bastian Leibe. Biternion Nets: Continuous Head Pose Regression from Discrete Training Labels. In *GCPD*, 2015.
- [6] Cyril Cauchois, Fabrice De Chaumont, Bruno Marhic, Laurent Delahoche, and Mélanie Delafosse. Robotic Assistance: an Automatic Wheelchair Tracking and Following Functionality by Omnidirectional Vision. In *IROS*, 2005.
- [7] Mark De Deuge, A Quadros, C Hung, and B Douillard. Unsupervised Feature Learning for Classification of Outdoor 3D Scans. In *ACRA*, 2013.
- [8] Christian Dondrup, Nicola Bellotto, Ferdian Jovan, Marc Hanheide, et al. Real-time Multisensor People Tracking for Human-Robot Spatial Interaction. In *ICRA*, 2015.
- [9] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep Sparse Rectifier Neural Networks. In *AISTATS*, 2011.
- [10] Chun-Rong Huang, Pau-Choo Chung, Kuo-Wei Lin, and Sheng-Chieh Tseng. Wheelchair detection using cascaded decision tree. *T-ITB*, 14(2):292–300, 2010.
- [11] J. D. Hunter. Matplotlib: A 2D graphics environment. *CS&E*, 9(3):90–95, 2007.
- [12] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *NIPS*, 2012.
- [14] Kevin Lai, Liefeng Bo, and Dieter Fox. Unsupervised Feature Learning for 3D Scene Labeling. In *ICRA*, 2014.
- [15] Bastian Leibe, Aleš Leonardis, and Bernt Schiele. Robust Object Detection with Interleaved Categorization and Segmentation. *IJCV*, 77(1-3):259–289, 2008.
- [16] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, and Scott Reed. SSD: Single Shot MultiBox Detector. *arXiv preprint arXiv:1512.02325*, 2015.
- [17] Daniel Maturana and Sebastian Scherer. VoxNet: A 3D Convolutional Neural Network for Real-Time Object Recognition. In *IROS*, 2015.
- [18] Pierre Merdrignac, Evangeline Pollard, and Fawzi Nashashibi. 2D Laser Based Road Obstacle Classification for Road Safety Improvement. In *ARSO*, 2015.
- [19] Fausto Milletari, Seyed-Ahmad Ahmadi, Christine Kroll, Annika Plate, Verena Rozanski, Juliana Maiostre, Johannes Levin, Olaf Dietrich, Birgit Ertl-Wagner, Kai Bötzel, et al. Hough-CNN: Deep Learning for Segmentation of Deep Brain Regions in MRI and Ultrasound. *arXiv preprint arXiv:1601.07014*, 2016.
- [20] Oscar Martinez Mozos, Ryo Kurazume, and Tsutomu Hasegawa. Multi-Part People Detection Using 2D Range Data. *IJSR*, 2(1):31–40, 2010.
- [21] Ashish Myles, Niels Da Vitoria Lobo, and Mubarak Shah. Wheelchair Detection in a Calibrated Environment. In *ACCV*, 2002.
- [22] Peter Ondruska, Julie Dequaire, Dominic Zeng Wang, and Ingmar Posner. End-to-End Tracking and Semantic Segmentation Using Recurrent Neural Networks. In *RSS, Workshop on Limits and Potentials of Deep Learning in Robotics*, 2016.
- [23] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine Learning in Python. *JMLR*, 12:2825–2830, 2011.
- [24] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified, Real-Time Object Detection. *arXiv preprint arXiv:1506.02640*, 2015.
- [25] Gernot Riegler, David Ferstl, Matthias Rütger, and Horst Bischof. Hough Networks for Head Pose Estimation and Facial Feature Localization. *BMVC*, 2014.
- [26] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *ICLR*, 2014.
- [27] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *ICLR*, 2015.
- [28] Sayanan Sivaraman and Mohan Manubhai Trivedi. Looking at Vehicles on the Road: A Survey of Vision-Based Vehicle Detection, Tracking, and Behavior Analysis. *TITS*, 14(4):1773–1795, 2013.
- [29] Luciano Spinello, Kai Oliver Arras, Rudolph Triebel, and Roland Siegwart. A Layered Approach to People Detection in 3D Range Data. In *AAAI*, 2010.
- [30] Luciano Spinello and Roland Siegwart. Human Detection using Multimodal and Multidimensional Features. In *ICRA*, 2008.
- [31] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *JMLR*, 15(1):1929–1958, 2014.
- [32] Jonathan J Tompson, Arjun Jain, Yann LeCun, and Christoph Bregler. Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation. In *NIPS*, 2014.
- [33] Dominic Zeng Wang and Ingmar Posner. Voting for voting in online point cloud object detection. In *Proceedings of Robotics: Science and Systems*, 2015.
- [34] Christoph Weinrich, Tim Wengefeld, Christof Schroeter, and Horst-Michael Gross. People Detection and Distinction of their Walking Aids in 2D Laser Range Data based on Generic Distance-Invariant Features. In *RO-MAN*, 2014.
- [35] Jodo Xavier, Marquidia Pacheco, Daniel Castro, António Ruano, and Urbano Nunes. Fast Line, Arc/Circle and Leg Detection from Laser Scan Data in a Player Driver. In *ICRA*, 2005.
- [36] Yuanpu Xie, Xiangfei Kong, Fuyong Xing, Fujun Liu, Hai Su, and Lin Yang. Deep Voting: A Robust Approach Toward Nucleus Localization in Microscopy Images. In *MICCAI*, 2015.

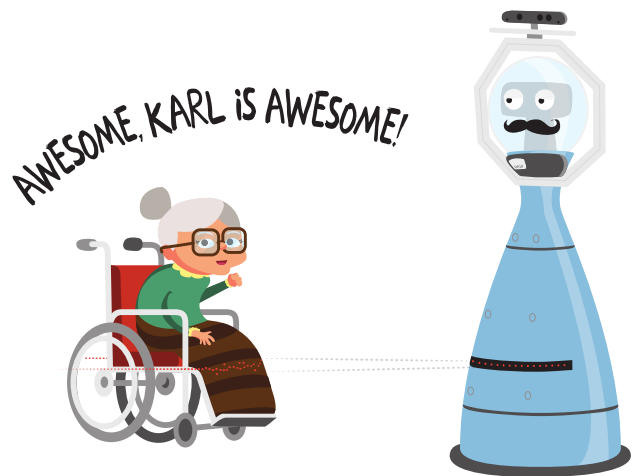


Image credit: Thibault Jan Beyer