

Visual Domain Adaptation for Monocular Depth Estimation on Resource-Constrained Hardware

Julia Hornauer¹, Lazaros Nalpantidis², and Vasileios Belagiannis¹

¹Ulm University, Ulm, Germany, {first.last}@uni-ulm.de

²DTU – Technical University of Denmark, Kgs. Lyngby, Denmark, lanalpa@elektro.dtu.dk

Abstract

Real-world perception systems in many cases build on hardware with limited resources to adhere to cost and power limitations of their carrying system. Deploying deep neural networks on resource-constrained hardware became possible with model compression techniques, as well as efficient and hardware-aware architecture design. However, model adaptation is additionally required due to the diverse operation environments. In this work, we address the problem of training deep neural networks on resource-constrained hardware in the context of visual domain adaptation. We select the task of monocular depth estimation where our goal is to transform a pre-trained model to the target's domain data. While the source domain includes labels, we assume an unlabelled target domain, as it happens in real-world applications. Then, we present an adversarial learning approach that is adapted for training on the device with limited resources. Since visual domain adaptation, i.e. neural network training, has not been previously explored for resource-constrained hardware, we present the first feasibility study for image-based depth estimation. Our experiments show that visual domain adaptation is relevant only for efficient network architectures and training sets at the order of a few hundred samples. Models and code are publicly available¹.

1. Introduction

Machine learning on resource-constrained hardware has emerged to an important research direction with applications on robotics [14], autonomous driving [9], and surveillance systems [27]. Executing learning-based algorithms directly on-site reduces the system latency, preserves the data privacy, and makes the system more reliable because

¹https://github.com/jhornauer/embedded_domain_adaptation

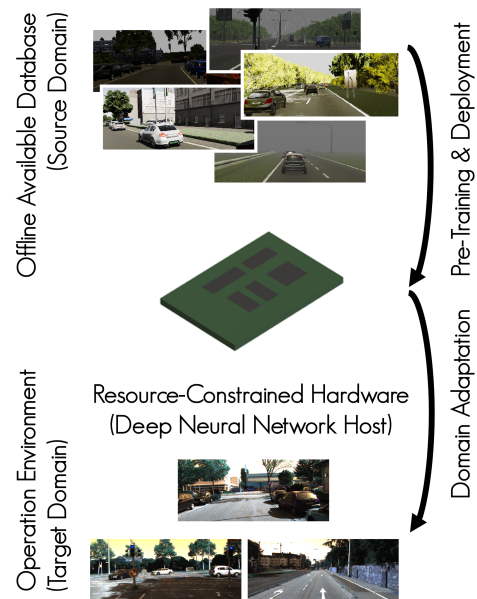


Figure 1: Resource-constrained hardware can be used for the deployment of deep neural networks. However, it is often required to further adapt the model because of the diverse operating environments. We study the problem of training a deep neural network on the embedded hardware in the context of domain adaption for monocular depth estimation. The visualized images are from the databases KITTI [10] and Virtual KITTI [4].

of the independence from external factors such as remote servers and communication networks. Moreover, resource-constrained hardware, such as embedded devices, is significantly less expensive than workstations or cloud services. Nevertheless, scaling up machine learning to real-world applications using resource-constrained hardware remains a challenge.

In practice, deep neural networks, the horsepower in

the field, are very successful models for deployment on resource-constrained systems. To reach real-time performance, the model complexity and memory footprint are reduced with network compression [3], as well as efficient [12] and hardware-aware [30] network architecture design. These approaches assume that further network training is not necessary after deployment. However, the integration of deep neural networks on low-cost embedded devices makes them more ubiquitous and at the same time exposes them to many and more diverse operation environments. Thus, on-device model adaptation is required for those devices to perform as expected. Training deep neural networks on the resource-constrained hardware, though, has not been addressed yet (illustrated in Fig. 1).

In this work, we address the problem of *training deep neural networks on resource-constrained hardware* in the context of visual domain adaptation. Our testbed is monocular image-based depth estimation where the model adaptation from the source to the target domain happens in an unsupervised manner. We assume a pre-trained model that resulted from the source domain data and it has been trained with supervision. The model pre-training takes place on a standard workstation. For the target domain, we suppose that data collection is possible, e.g. a mobile agent, but ground-truth depth maps are not available. Then, we present an adversarial learning approach [16] that is adapted for training on the resource-constrained hardware. Given the hardware limitations, we employ an efficient network architecture [29] for depth estimation. Besides, we also consider a complex architecture [17] for comparison reasons. Since domain adaptation has not been previously explored for resource-constrained devices, we present the first feasibility study for the perception task of monocular depth estimation. We analyze the training process of the deep neural network regarding the data and training set size, model complexity, and energy consumption. Our experiments show that visual domain adaptation on resource-constrained hardware—and thus deep neural network training—is meaningful only for efficient network architectures and training sets at the order of a few hundred samples.

2. Related Works

Embedded depth estimation deals with the depth prediction from a single image. In [17], [7], and [6], promising results are shown based on deep neural networks for regression. However, none of these methods is designed for usage on resource-constrained devices. Instead, complex models have been proposed, which make hardware deployment challenging. In [29], [26], [23] and [25], the effectiveness of different lightweight depth estimation network architectures is demonstrated on embedded devices, such as the Raspberry Pi and NVIDIA Jetson TX2. Poggi *et al.* [26] propose the lightweight architecture PyD-Net with a pyramidal fea-

ture extractor to train in an unsupervised manner for CPU processor usage. Wofk *et al.* [29] design their neural network, FastDepth, with depth-wise separable convolutions. Oh *et al.* [23] propose a Repetition-Reduction block within the encoder, and a condensed decoding connection block for feature propagation to the decoder, in an encoder-decoder architecture. Peluso *et al.* [25] demonstrate their accuracy-driven quantization-aware training method adapted for the ARMv7 core on PyD-Net [26]. Nevertheless, their approaches address only the problem of deployment.

Visual domain adaptation refers to the generalization of a network trained on a source domain to some related target domain [5]. In domain adaptation for depth-estimation, a major issue is the annotation of dense depth maps. In [2], [32] and [31] image-to-image translation is used to exploit data generation for addressing the problem. Instead of image-to-image translation Kundu *et al.* [16] explore an adversarial domain adaptation setting with two discriminators and different regularization techniques to obtain a monocular depth estimation model originally trained on synthetic data. Lasinger *et al.* [18] target the generalization ability towards different domains by training their MiDaS network with multiple datasets of different scenes and environments. As most datasets differ in their depth ground truth representation, they create an objective that is unaffected by the various label types. Aleotti *et al.* [1] create a large-scale dataset, called WILD, by making predictions on images of different environments with the pre-trained, large-scale MiDaS [18] network. The resulting dataset is used to train lightweight models with a high generalization ability for deployment on handheld devices. We consider the adversarial domain adaptation [16] as suitable for training on our hardware with limited resources. The aforementioned works accomplished great success in depth estimation and domain adaptation, but they rely on costly training of complex neural networks or only efficient model deployment on custom hardware.

Resource-constrained hardware training has not been considered at all. For example, lightweight architectures have been proposed in [29], [26] and [23] with the aim to deploy models in real-time on embedded devices, where there is only a minor performance loss. More general, Zhang *et al.* [30] propose to use hardware-aware neural network search to adapt the model for deployment to the dedicated hardware. Li *et al.* [19] propose linear learning rate scheduling with regard to limited training duration in terms of iteration. Although the design and deployment of hardware-aware and hardware-efficient networks has been studied in the past, the problem of training directly on the device with limited resources has not been addressed in those works.

3. Approach

We present the problem, our system and the domain adaptation algorithm for image-based depth estimation on resource-constrained hardware.

3.1. Problem formulation

Let $f_\theta : \mathbb{R}^{w \times h \times 3} \rightarrow \mathbb{R}^{w \times h \times 1}$ to be the function that maps the image \mathbf{x} with dimensions $w \times h \times 3$ to the depth map \mathbf{y} with dimensions $w \times h \times 1$, where the function is represented by a deep neural network with parameters θ . The model f_θ is trained with supervision on the database $\mathcal{S} = \{(\mathbf{x}, \mathbf{y})_i\}_{i=1}^{|\mathcal{S}|}$, which we refer to as the source domain. Consider now a different domain that is expressed by the data collection $\mathcal{T} = \{(\mathbf{x})_i\}_{i=1}^{|\mathcal{T}|}$, referred to as the target domain. The target domain represents the operation environment. In the target domain, we assume not to have access to the ground-truth depth map \mathbf{y} . Furthermore, only a resource-constrained hardware system, e.g. embedded device, is available for the model deployment. Then, our task is to adapt the model parameters θ to the target domain without supervision, by relying only on the collected set \mathcal{T} and the limited resources. Given the constrained hardware system and the visual domain adaptation task, we examine the training feasibility of the deep neural network w.r.t the data and training set size, as well as model complexity and the energy consumption.

3.2. Resource-constrained hardware

We consider the NVIDIA Jetson Nano for our evaluations. The processing unit consists of the ARM-A57 CPU with 4GB RAM and the 128-core CUDA Maxwell GPU. The system is running a Linux-based operating system provided by NVIDIA and stored on a 128 GB SD card with up to 100 MB/s transfer speeds. The SD card memory is sufficient for the operating system, the executed libraries, developed algorithms and stored data. Finally, the domain adaptation is implemented in the PyTorch [24] framework. Note that we reckoned with the Raspberry Pi 4 for our experiments as well. However, the available processing power is not sufficient for training image-based deep neural networks in a reasonable time, since it is not equipped with a GPU.

3.3. Visual domain adaption with limited resources

We build the domain adaptation framework based on AdaDepth [16], an adversarial domain adaptation approach for depth estimation. This method is chosen because it relies on a less expensive setup with two discriminators instead of image-to-image translation as in the related approaches [2], [32] and [31]. We assume that the model f_θ is composed of the encoder ϕ and the decoder ψ network,

such that:

$$f_\theta(\mathbf{x}) = \psi(\phi(\mathbf{x})). \quad (1)$$

The encoder ϕ maps the input image \mathbf{x} to a latent space, whereas the decoder ψ maps the latent space to the pixel-wise depth prediction.

At first, the encoder-decoder model f_θ is pre-trained on the source domain database \mathcal{S} . Similarly to AdaDepth [16], ψ is shared between the two domains, and thus it is not adapted. The domain adaptation transforms only the source domain encoder ϕ_s to the target domain encoder ϕ_t . In practice, only a subset of the encoder parameters will be adapted, as we discuss later in Sec. 3.4.

During training, we rely on the latent space discriminator and the depth map discriminator to distinguish between the latent space and the depth maps of the source and target domain respectively. The latent space discriminator $LD(\cdot)$ is trained to predict the domain of the latent space representations $\phi_s(\mathbf{x}_s)$ and $\phi_t(\mathbf{x}_t)$. Similar to LSGAN [20], we define the objective as:

$$\begin{aligned} \mathcal{L}_{LD} = & \mathbb{E}_{\mathbf{x}_s \sim \mathcal{S}} [\gamma (LD(\phi_s(\mathbf{x}_s)) - 1)^2] + \\ & \mathbb{E}_{\mathbf{x}_t \sim \mathcal{T}} [\gamma (LD(\phi_t(\mathbf{x}_t)))^2] + \\ & \mathbb{E}_{\mathbf{x}_t \sim \mathcal{T}} [(1 - \gamma) (LD(\phi_t(\mathbf{x}_t)) - 1)^2], \end{aligned} \quad (2)$$

where $LD(\phi_t(\mathbf{x}_t)) - 1$ stands for the adversary, i.e. setting the target domain as source domain (indicated by 1) and $\gamma \in \{0, 1\}$ is used for updating the discriminators. The discriminator is updated when $\gamma = 1$, while the encoder ϕ_t is updated for $\gamma = 0$ in a second step.

In addition, the depth map discriminator $DD(\cdot)$ takes the image and the depth map as input for identifying the domain. It is trained to distinguish the source ground-truth depth map \mathbf{y}_s from the target's domain predicted depth map $\psi(\phi_t(\mathbf{x}_t))$. The objective function for the depth map discriminator is given by:

$$\begin{aligned} \mathcal{L}_{DD} = & \mathbb{E}_{(\mathbf{x}_s, \mathbf{y}_s) \sim \mathcal{S}} [\gamma (DD(\mathbf{x}_s, \mathbf{y}_s) - 1)^2] + \\ & \mathbb{E}_{\mathbf{x}_t \sim \mathcal{T}} [\gamma (DD(\mathbf{x}_t, \psi(\phi_t(\mathbf{x}_t))))^2] + \\ & \mathbb{E}_{\mathbf{x}_t \sim \mathcal{T}} [(1 - \gamma) (DD(\mathbf{x}_t, \psi(\phi_t(\mathbf{x}_t))) - 1)^2]. \end{aligned} \quad (3)$$

Finally, we rely on the *domain consistency regularization* loss from AdaDepth [16] as a measure for the prevention of mode collapse. It minimizes the distance of the source and target latent space representation based on the target images. It is given by:

$$\mathcal{L}_{reg} = \mathbb{E}_{\mathbf{x}_t \sim \mathcal{T}} [(1 - \gamma) \|\phi_s(\mathbf{x}_t) - \phi_t(\mathbf{x}_t)\|_1]. \quad (4)$$

The adversarial training starts from the pre-trained model f_θ on the source domain data. As it progresses, ϕ_t is trained to produce samples that seem to be originating from the source domain. The minimization of all objectives is expressed as:

$$\arg \min_{\gamma=1, LD, DD} \min_{\gamma=0, \phi_t} \mathcal{L}_{LD} + \mathcal{L}_{DD} + \lambda \mathcal{L}_{reg}, \quad (5)$$

where the hyper-parameter λ controls the influence of the regularization term. For every iteration, the first minimization updates the parameters of the depth map discriminator $DD(\cdot)$ and latent space discriminator $LD(\cdot)$, while the second minimization updates the parameters of the target encoder ϕ_t . The regularizer, finally, is applied only once during the second minimization where $\gamma = 0$. The training process takes place on the resource-constrained hardware.

3.4. Neural network architectures for depth estimation

Two main limitations of the resource-constrained hardware are the computing power and the available memory. The standard depth estimation network architectures are computationally complex and memory demanding [17]. To address this problem, there have been recently proposed lightweight architectures for depth estimation [29]. For our experiments, we consider both a lightweight and a complex network architecture: FastDepth and ResNet-UpProj respectively.

Lightweight architecture FastDepth [29] is built with MobileNet [12] as the encoder and a lightweight decoder with depth-wise separable convolutions, followed by nearest-neighbor interpolation for up-sampling. The network counts 3.93M parameters in total [29]. Additional skip connections between the encoder and the decoder are added for feature propagation to compensate for the small number of parameters. In line with the concept of [16], only the last four layers of the encoder are trained during the domain adaptation.

Complex architecture Laina *et al.* design a encoder-decoder depth estimation network with ResNet-50 as the encoder and up-projection blocks within the decoder [17]. The architecture is implemented with five up-projections blocks, similar to FastDepth [29], and has 63.6M parameters in total. This is a significantly larger number of parameters compared to FastDepth. Similar to [16], the 5-th ResNet block is only adjusted during the adaptation.

At last, the depth map discriminator $DD(\cdot)$ follows the PatchGAN [13] network structure, while the latent space discriminator $LD(\cdot)$ is a convolutional discriminator which we later present based on the evaluation.

4. Evaluation

In this section, we present the findings of our approach on visual domain adaptation on the resource-constrained hardware for the demanding task of monocular depth estimation. In our evaluation, we consider the scenarios of indoor, as well as outdoor environments where we rely on

four standard benchmarks for depth estimation. In each scenario, we study the factors of image resolution and training set size, model complexity and energy consumption during training on the device with the limited resources. We report the mean performance after five runs for each experiment.

4.1. Indoor and outdoor benchmarks

In the indoor evaluation, the NYU Depth v2 dataset [21] serves as the source domain database, while the target domain is represented by the DIML/CVL RGB-D data set [15]. In the outdoor scenario, the synthetic Virtual KITTI (vKITTI) [4] is the source domain and the target domain is the KITTI database [10].

From NYU Depth v2 to DIML/CVL RGB-D (indoors)

The NYU Depth v2 (source domain) is an indoor dataset taken at 464 different scenes by a depth-sensing camera. The scenes are split into 249 scenes for training and 215 for testing. The images have a resolution of 480x640. We rely on the training set for creating the pre-trained model. The DIML/CVL RGB-D database (target domain) provides a large-scale indoor dataset with 220k training images taken by a Microsoft Kinect v2 at 283 different scenes of 18 different categories. In addition, there is a smaller set of 1500 training images and 500 samples for testing. The images and aligned depth maps have a resolution of 756x1344. The large-scale set utilization is not realistic because of the limited computational power and memory space, and thus we select the smaller set for the target domain.

From Virtual KITTI to KITTI (outdoors)

The synthetic virtual KITTI (vKITTI) dataset is used as the source domain for the outdoor scenario. It consists of 21260 synthetic image-depth pairs with a resolution of 375x1242, which we make use for training our networks. We rely only on the left images of the image-depth pairs. Since the maximum depth of KITTI is 80m, the ground truth depth is clipped to this value as maximum. The two examined architectures are pre-trained on vKITTI. On the other hand, the KITTI dataset is a real-world computer vision benchmark (target domain) with 42382 rectified stereo pairs. RGB image pairs with corresponding velodyne points are provided in the raw data, where we rely on the left images. The Eigen-split [7] is used for the evaluation since it is a common evaluation protocol in the literature [16], [26], [2], [31]. Eigen *et al.* divide the data into 22600 samples for training, 888 samples for validation and 697 samples for testing.

4.2. Training protocol & implementation

In both settings, the FastDepth [29] and ResNet-UpProj [17] architectures are first trained on the source domains \mathcal{S} . The stochastic gradient descent (SGD) solver is

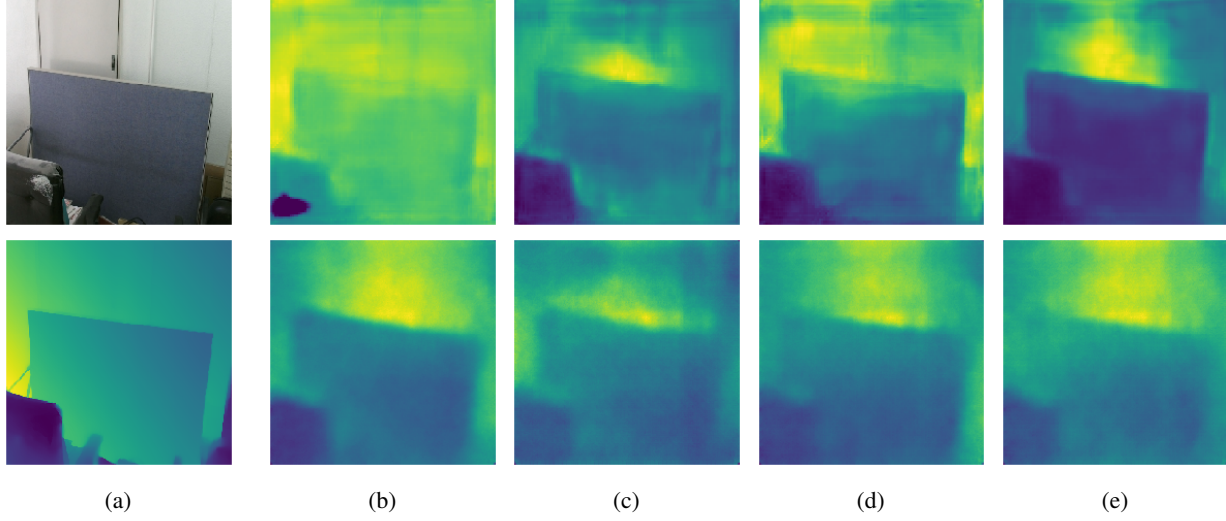


Figure 2: Visual results on DIML/CVL RGB-D. (a): RGB (top) & Ground truth (bottom). (b) - (e): Training FastDepth (top) & ResNet (bottom) with a resolution of 224x224. (b) Source-only; (c) Domain adaptation for $|\mathcal{A}| = |\mathcal{B}| = 100$; (d) Domain adaptation for $|\mathcal{A}| = |\mathcal{B}| = 500$, (e) Domain adaptation for $|\mathcal{A}| = |\mathcal{B}| = 1000$.

Table 1: Evaluation of the indoor experiment training FastDepth and ResNet-UpProj with a resolution of 224x224. The results of the models trained on the source domain (source) and domain adaption with subsets $|\mathcal{A}| = |\mathcal{B}| \in \{100, 500, 1000\}$ are listed. Depth prediction results are evaluated on the DIML/CVL RGB-D test set. For accuracy δ higher is better, for RMSE lower is better. We also report the peak power consumption, average energy consumption per epoch, multiply-accumulate operations (MACs), as well as the average training duration per epoch (in minutes) and the inference time (in milliseconds) both for the Jetson device and for a PC workstation, for comparison.

Architecture	Training Data	δ_1	δ_2	δ_3	RMSE	Power [W]	Energy [Wh]	MACs [G]	Training [min]		Inference [ms]	
Fast-Depth	Source	0.493	0.847	0.958	0.824							
	1000	0.560	0.856	0.953	0.801	12.4	1.6		13.5	2.3		
	500	0.563	0.861	0.954	0.796	11.8	0.8	0.76	7	1.2	33	10
	100	0.562	0.862	0.955	0.803	11.3	0.2		2.5	0.2		
	Source	0.444	0.816	0.947	0.872							
ResNet-UpProj	1000	0.576	0.857	0.941	0.777	11.9	9.4		61.5	4.8		
	500	0.578	0.858	0.942	0.755	11.9	4.6	32.25	31	2.4	610	44
	100	0.573	0.861	0.951	0.749	11.9	1.0		6.5	0.5		
	Source											

used to optimize the \mathcal{L}_1 distance between the input image \mathbf{x} and the depth map \mathbf{y} . Data augmentation, similar to [29], is applied too. Random color jitter, random rotation, random scaling and center cropping is applied before the images are downsampled to a specific resolution. The final resolution for indoors is 224x224, while for outdoors we consider a lower 256x512 and a higher 288x704 resolution. The source domain training takes place on a PC workstation². Then, the domain adaptation on both settings is performed using subsets of the training data $\mathcal{A} \subseteq \mathcal{S}$ and $\mathcal{B} \subseteq \mathcal{T}$ such that $|\mathcal{A}| = |\mathcal{B}|$. The subsets \mathcal{A} and \mathcal{B} of the training

²The workstation has a 6-core processor with 16GB RAM and 6GB GPU memory.

data \mathcal{S} and \mathcal{T} are selected randomly for every run. Different domain adaptation versions are trained using subsets of increasing size on the NVIDIA Jetson Nano. For all settings λ is set to 0.7. To imitate data collection on a real-world environment, we select small target domain sets such that $|\mathcal{A}| = |\mathcal{B}| \in \{100, 500, 1000\}$. Finally, the performance is evaluated on a target test set, which has not been observed during the domain adaptation.

Indoors The discriminator $LD(\cdot)$ consists of three convolutions with kernel size 3, each followed by a leaky rectified linear activation, and one last linear layer. After the last two convolutions dropout with probability 0.6 is ap-

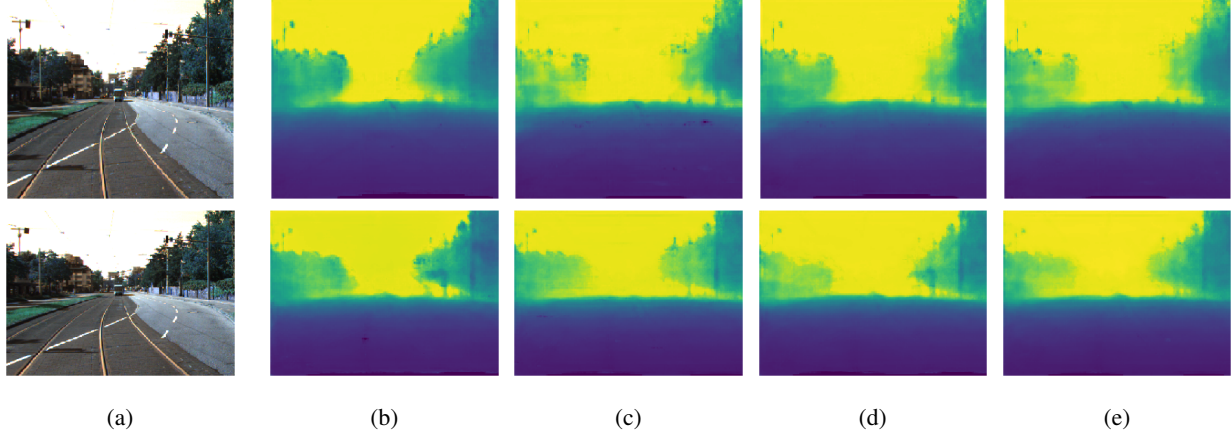


Figure 3: Visual results on KITTI test image training FastDepth with a resolution of 256x512 (top) and 288x704 (bottom). (a) RGB; (b) Source-only; (c) Domain adaptation for $|\mathcal{A}| = |\mathcal{B}| = 100$; (d) Domain adaptation for $|\mathcal{A}| = |\mathcal{B}| = 500$; (e) Domain adaptation for $|\mathcal{A}| = |\mathcal{B}| = 1000$.

plied. The two discriminators and the target encoder are optimized using the ADAM solver with $\beta_1 = 0.5$ and $\beta_2 = 0.999$. For FastDepth the learning rates are set to 0.0002 and for ResNet-UpProj to 0.00002. The augmentation is maintained during the adaptation in the indoor setting. Up to 25 epochs are necessary for the domain adaptation with FastDepth, i.e. 25 epochs for 100/500 samples and 15 epochs for the rest. With the ResNet-UpProj architecture it is 15 and 10 epochs respectively. One limiting factor of the embedded device is the available GPU memory. This leads to a maximum batch size depending on the model size. Training the adaptation with FastDepth a batch size of 16 per domain is selected. This is not possible with the larger ResNet-UpProj architecture, where the adaptation is trained with sub-batches of 2 samples per batch to maintain a resolution of 224x224.

Outdoors The discriminator $LD(\cdot)$ of the indoor adaptation is adopted, but the kernels in the convolutions are replaced by (4,7), (3,5) and (3,5) because of the differing resolution. The adaptation is performed for the resolutions 256x512 and 288x704. The outdoors resolution is too large for ResNet-UpProj, where we cannot run inference or training due to the limited resources. We rely only on FastDepth that is trained with the SGD optimizer with momentum set to 0.9 during the adaptation. The learning rates are set to $1e-4$ or $1e-5$ depending on the specific setup. The adaptation training requires up to 10 epochs, i.e. 10 epochs for 100/500 samples, and 5 epochs for the rest. The maximum possible batch size depends on the network parameters number. For both resolutions, FastDepth batch size was 4.

4.3. Evaluation metrics and baselines

The depth prediction is evaluated on the standard metrics of accuracy δ_1 , δ_2 and δ_3 , as well as the root-mean-square error (RMSE), similar to the literature [32], [31]. We report results for both network architectures by training only on the source domain (*source*) and afterwards with domain adaptation. As an addition, for domain adaptation in the outdoor setting we follow another standard protocol in which the predicted depth maps are scaled by $median(y_{gt})/median(y_{pred})$ as in [16], [33]. The evaluation on the resource-constrained hardware is based on the peak power consumption, average energy consumption per epoch, multiply-accumulate operations (MACs), as well as the average training duration per epoch and the inference time both for the NVIDIA Jetson Nano (Jetson) device and for a PC workstation, as reference.

4.4. Results and discussion

We present the results for the indoor domain adaptation, followed by the outdoor configuration.

Indoors adaptation (NYU to DIML/CVL) In the indoor evaluation, four different subsets with $|\mathcal{A}| = |\mathcal{B}| \in \{100, 500, 1000\}$ of the training data are randomly selected to conduct the domain adaptation. Table 1 shows the performance results of FastDepth and ResNet-UpProj architectures on the DIML/CVL RGB-D test set and presents hardware-specific metrics and the model complexity. At first, it is clear that the domain adaptation is always helpful compared to only applying the model trained in the source domain (*source*). Next, the lightweight FastDepth architecture functions with up to 1000 training data with faster training time per epoch than the ResNet-UpProj architec-

Table 2: Evaluation of the outdoor experiment training FastDepth with a resolution of 256x512 and 288x704. The results of the models trained on the source domain (source) and domain adaption with subsets $|\mathcal{A}| = |\mathcal{B}| \in \{100, 500, 1000\}$ are listed. The depth prediction results are evaluated on the KITTI test set. For accuracy δ higher is better, for RMSE lower is better. We also report the peak power consumption, average energy consumption per epoch, multiply-accumulate operations (MACs), as well as the average training duration per epoch (in minutes) and the inference time (in milliseconds) both for the Jetson device and for a PC workstation, for comparison.

Resolution	Training Data	δ_1	δ_2	δ_3	RMSE	Power [W]	Energy [Wh]	MACs [G]	Training [min]		Inference [ms]	
									Jetson	PC	Jetson	PC
256x512	Source	0.549	0.790	0.893	8.505							
	1000	0.649	0.825	0.906	9.337	11.8	3.2		23.4	2.3		
	500	0.637	0.822	0.905	9.385	11.8	1.6	1.98	11	1.3	37	10
	100	0.630	0.817	0.901	9.449	11.4	0.3		2.4	0.3		
288x704	Source	0.540	0.752	0.857	9.070							
	1000	0.652	0.825	0.906	9.567	11.8	5.2		34.5	3.1		
	500	0.643	0.813	0.896	9.883	11.9	2.6	3.07	17	1.6	38	10
	100	0.647	0.818	0.900	9.906	11.9	0.4		3.5	0.3		

Table 3: Evaluation of the outdoor experiment training FastDepth with a resolution of 256x512 and 288x704. For domain adaptation the results are reported with *sample-wise median scaling* as in [16] and [33]. The results of the models trained on the source domain (source) and domain adaptation with subsets $|\mathcal{A}| = |\mathcal{B}| \in \{100, 500, 1000\}$ are listed. The depth prediction results are evaluated on the KITTI test set. For accuracy δ higher is better, for RMSE lower is better.

Resolution	Training Data	δ_1	δ_2	δ_3	RMSE
256x512	Source	0.549	0.790	0.893	8.505
	1000	0.649	0.863	0.938	7.813
	500	0.650	0.863	0.938	7.788
	100	0.637	0.854	0.930	7.942
288x704	Source	0.540	0.752	0.857	9.070
	1000	0.635	0.862	0.939	7.924
	500	0.626	0.856	0.934	7.952
	100	0.629	0.855	0.934	8.168

ture. Moreover, the inference time of FastDepth is significantly faster. Also, the difference between the two architectures is large in the MAC operations. It is clear that the complex ResNet-UpProj architecture is more suitable for powerful computers. This is also inferred when comparing with the workstation’s training and inference time. Overall, relying on between 500 to 1000 training samples for adaptation results in a good balance between performance, training time and energy consumption. Finally, the complexity of the architecture does not play an important role in the energy consumption, which remains comparable for both models, considering that the more complex architec-

ture converges faster. Fig. 2 illustrates the visual depth map results for one test sample. In this figure, for both architectures, the improvement of the model adaptation is visible (c) compared to training only on the source domain (b). Especially the lightweight architecture shows a significant improvement: the objects become distinguishable from the background and the object in the lower-left corner is predicted as closest.

Adaptation vKITTI to KITTI We follow the same configuration for this experiment as well. In Table 2, we report the results of the domain adaptation in the resource-constrained hardware. In addition, we report the domain adaptation results with median-scaling in Table 3. We rely now on two image resolutions, where the depth prediction results are in a similar range. The higher resolution adds considerable training time, but the inference time is similar. Furthermore, the energy consumption is not significantly affected by the input resolution. On the other hand, the complex ResNet-UpProj architecture is not capable of running on the NVIDIA Jetson Nano due to the memory limitation with either image resolutions. Finally, the results after adaptation overall improve for both protocols. Moreover, our visual results show improvement from source only to domain adaptation, as shown in Fig. 3.

Discussion Both indoor and outdoor evaluations demonstrate the feasibility of domain adaptation on the resource-constrained hardware in a meaningful period of time. The image resolution and number of samples of course affects the training time, which can go up to 61.5 minutes for ResNet (indoors) and 34.5 minutes for FastDepth (outdoors). The energy consumption is not a concern for any of our experiments, while the inference time is architecture

dependent. For instance, we reach 10 milliseconds with FastDepth for both experiments. Thus, we conclude that training directly on the embedded device is possible with adversarial training. Given that monocular depth estimation is a demanding task, we expect perception tasks such as human trajectory estimation [11] and gesture recognition [28] or other robotics applications [22, 8] to be easier transferable to the resource-constrained hardware. Developing the complete perception of an autonomous agent on resource-constrained hardware is part of our future work. The main benefit will be to reduce the overall energy demands, while maintaining reliable performance.

5. Conclusion

We presented the first feasibility study on training deep neural networks on resource-constrained hardware in the context of visual domain adaptation. Our testbed is monocular depth estimation, where domain adaptation is accomplished without supervision. We extended an adversarial learning approach to function on the device with limited resources. In two evaluations using four standard databases, we have shown that domain adaptation on the resource-constrained hardware is manageable for lightweight architectures based on a few hundred samples from the target domain. Our study indicates that the deployment-hardware needs always to be considered along with the training algorithm, neural network architecture and the type of supervision to scale up machine learning to real-world applications.

References

- [1] Filippo Aleotti, Giulio Zaccaroni, Luca Bartolomei, M. Poggi, Fabio Tosi, and S. Mattoccia. Real-time single image depth perception in the wild with handheld devices. *ArXiv*, abs/2006.05724, 2020. 2
- [2] A. Atapour-Abarghouei and T. P. Breckon. Real-time monocular depth estimation using synthetic data with domain adaptation via image style transfer. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2800–2810, 2018. 2, 3, 4
- [3] Vasileios Belagiannis, Azade Farshad, and Fabio Galasso. Adversarial network compression. In *Computer Vision - ECCV 2018 Workshops - Munich, Germany, September 8-14, 2018, Proceedings, Part IV*, volume 11132 of *Lecture Notes in Computer Science*, pages 431–449. Springer, 2018. 2
- [4] Yohann Cabon, Naila Murray, and Martin Humenberger. Virtual KITTI 2, 2020. 1, 4
- [5] G. Csurka. Domain adaptation for visual applications: A comprehensive survey. *ArXiv*, abs/1702.05374, 2017. 2
- [6] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2650–2658, 2015. 2
- [7] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *ArXiv*, abs/1406.2283, 2014. 2, 4
- [8] Nico Engel, Stefan Hoermann, Markus Horn, Vasileios Belagiannis, and Klaus Dietmayer. Deeplocalization: Landmark-based self-localization with deep neural networks. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pages 926–933. IEEE, 2019. 8
- [9] Alexander Frickenstein, Manoj-Rohit Vemparala, Jakob Mayr, Naveen-Shankar Nagaraja, Christian Unger, Federico Tombari, and Walter Stechele. Binary dad-net: Binarized driveable area detection network for autonomous driving. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2295–2301. IEEE, 2020. 1
- [10] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. *International Journal of Robotics Research (IJRR)*, 2013. 1, 4
- [11] Irtiza Hasan, Francesco Setti, Theodore Tsesmelis, Vasileios Belagiannis, Sikandar Amin, Alessio Del Bue, Marco Cristani, and Fabio Galasso. Forecasting people trajectories and head poses by jointly reasoning on tracklets and vislets. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1267–1278, 2019. 8
- [12] A. Howard, Menglong Zhu, Bo Chen, D. Kalenichenko, W. Wang, Tobias Weyand, M. Andreetto, and H. Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. *ArXiv*, abs/1704.04861, 2017. 2, 4
- [13] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and A. Efros. Image-to-image translation with conditional adversarial networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976, 2017. 4
- [14] Jae In Kim, DongWook Kim, Matthew Krebs, Young Soo Park, and Yong-Lae Park. Force sensitive robotic end-effector using embedded fiber optics and deep learning characterization for dexterous remote manipulation. *IEEE Robotics and Automation Letters*, 4(4):3481–3488, 2019. 1
- [15] Youngjung Kim, Hyungjoo Jung, D. Min, and K. Sohn. Deep monocular depth estimation via integration of global and local predictions. *IEEE Transactions on Image Processing*, 27:4131–4144, 2018. 4
- [16] J. N. Kundu, P. K. Uppala, A. Pahuja, and R. V. Babu. AdaDepth: Unsupervised content congruent adaptation for depth estimation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2656–2665, 2018. 2, 3, 4, 6, 7
- [17] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 239–248, 2016. 2, 4
- [18] Katrin Lasinger, Ren’e Ranftl, K. Schindler, and V. Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, PP, 2020. 2
- [19] Mengtian Li, Ersin Yumer, and Deva Ramanan. Budgeted training: Rethinking deep neural network training under resource constraints. *ICLR*, 2020. 2

- [20] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017. 3
- [21] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *ECCV*, 2012. 4
- [22] Christian Nissler, Nikoleta Mouriki, Claudio Castellini, Vasileios Belagiannis, and Nassir Navab. Omg: introducing optical myography as a new human machine interface for hand amputees. In *2015 IEEE International Conference on Rehabilitation Robotics (ICORR)*, pages 937–942. IEEE, 2015. 8
- [23] S. Oh, H. S. Kim, J. Lee, and J. Kim. RRNet: Repetition-reduction network for energy efficient depth estimation. *IEEE Access*, 8:106097–106108, 2020. 2
- [24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*, pages 8026–8037, 2019. 3
- [25] V. Peluso, A. Cipolletta, A. Calimera, M. Poggi, F. Tosi, and S. Mattoccia. Enabling energy-efficient unsupervised monocular depth estimation on armv7-based platforms. In *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 1703–1708, 2019. 2
- [26] M. Poggi, F. Aleotti, F. Tosi, and S. Mattoccia. Towards real-time unsupervised monocular depth estimation on cpu. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5848–5854, 2018. 2, 4
- [27] Weihua Sheng, Yongsheng Ou, Duy Tran, Eyosiyas Tadesse, Meiqin Liu, and Gangfeng Yan. An integrated manual and autonomous driving framework based on driver drowsiness detection. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 4376–4381. IEEE, 2013. 1
- [28] Julian Wiederer, Arij Bouazizi, Ulrich Kressel, and Vasileios Belagiannis. Traffic control gesture recognition for autonomous vehicles. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10676–10683. IEEE, 2020. 8
- [29] D. Wofk, F. Ma, T. Yang, S. Karaman, and V. Sze. FastDepth: Fast monocular depth estimation on embedded systems. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 6101–6108, 2019. 2, 4, 5
- [30] Xiaofan Zhang, Haoming Lu, Cong Hao, Jiachen Li, Bowen Cheng, Y. Li, Kyle Rupnow, Jinjun Xiong, Thomas Huang, Humphrey Shi, W. Hwu, and D. Chen. SkyNet: a hardware-efficient method for object detection and tracking on embedded systems. *ArXiv*, abs/1909.09709, 2020. 2
- [31] S. Zhao, H. Fu, M. Gong, and D. Tao. Geometry-aware symmetric domain adaptation for monocular depth estimation. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9780–9790, 2019. 2, 3, 4, 6
- [32] C. Zheng, T. Cham, and J. Cai. T2Net: Synthetic-to-realistic translation for solving single-image depth estimation tasks. In *ECCV*, 2018. 2, 3, 6
- [33] Tinghui Zhou, M. Brown, Noah Snavely, and D. Lowe. Unsupervised learning of depth and ego-motion from video. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6612–6619, 2017. 6, 7