

CA-SpaceNet: Counterfactual Analysis for 6D Pose Estimation in Space

Shunli Wang^{1,2†}, Shuaibing Wang^{1,2†}, Bo Jiao^{1,2}, Dingkang Yang^{1,2},
Liuzhen Su^{1,2}, Peng Zhai^{1,2}, Chixiao Chen^{1,2*} and Lihua Zhang^{1,3,2,4*}

Abstract—Reliable and stable 6D pose estimation of uncooperative space objects plays an essential role in on-orbit servicing and debris removal missions. Considering that the pose estimator is sensitive to background interference, this paper proposes a counterfactual analysis framework named CA-SpaceNet to complete robust 6D pose estimation of the space-borne targets under complicated background. Specifically, conventional methods are adopted to extract the features of the whole image in the factual case. In the counterfactual case, a non-existent image without the target but only the background is imagined. Side effect caused by background interference is reduced by counterfactual analysis, which leads to unbiased prediction in final results. In addition, we also carry out low-bit-width quantization for CA-SpaceNet and deploy part of the framework to a Processing-In-Memory (PIM) accelerator on FPGA. Qualitative and quantitative results demonstrate the effectiveness and efficiency of our proposed method. To our best knowledge, this paper applies causal inference and network quantization to the 6D pose estimation of space-borne targets for the first time. The code is available at <https://github.com/Shunli-Wang/CA-SpaceNet>.

I. INTRODUCTION

As the eye of the spacecraft, vision-based navigation system is a crucial technology in many unmanned space missions. 6D pose estimation of space-borne objects is the premise of the navigation system. Fig. 1(a)&(b) shows two practical applications of 6D pose estimation in space: automatic docking and debris removal missions. Compared with terrestrial applications, many factors should be considered, such as harsh imaging conditions caused by the lack of atmospheric scattering and limited computing resources and power consumption. The robust and efficient 6D pose estimator is the key to ensuring the regular operation of on-orbit service.

In recent years, there have been some studies [1], [2], [3], [4], [5] in space engineering and computer vision community to explore 6D pose estimation of space-borne targets. Although considerable performance has been achieved, many methods directly migrate models from terrestrial to space scene without considering the particularity of space mission. In addition, these works mainly focus on the performance improvement of the model and ignore the power consumption and latency of the actual deployment on real spacecraft.

This work is supported by Shanghai Municipal Science and Technology Major Project (2021SHZDZX0103) and NSFC Grant (61974033).

¹Academy for Engineering & Technology, Fudan University. ²Engineering Research Center of AI and Robotics, Ministry of Education, China. ³Jilin Provincial Key Laboratory of Intelligence Science & Engineering, China. ⁴Artificial Intelligence and Unmanned Systems Engineering Research Center of Jilin Province, China

[†] The first two authors contributed equally to this work.

* Corresponding author, Email: {lihuaazhang, cxchen}@fudan.edu.cn

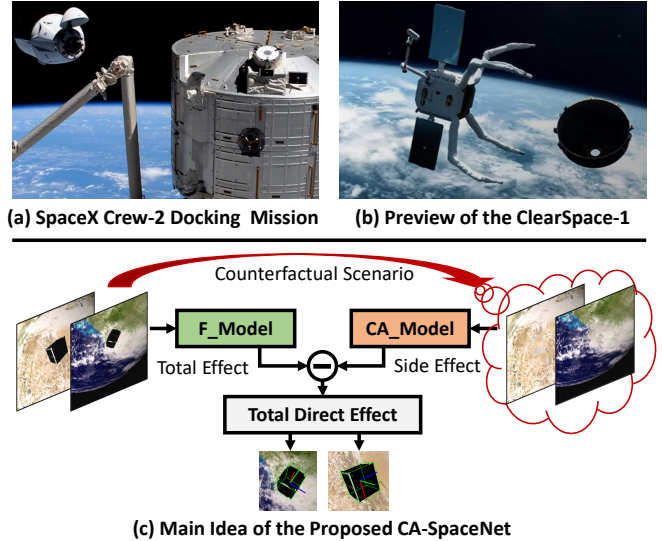


Fig. 1. Practical applications of the 6D pose estimation in many space missions, such as automatic docking and debris removal. (a) shows a screenshot of the docking in NASA’s SpaceX Crew-2 mission performed in 2021. (b) demonstrates a preview of the ClearSpace-1 satellite proposed by ESA and ClearSpace company which will be launched in 2025. (c) The complicated background of aerial images will interfere with the stability of the 6D pose estimator. Therefore, this paper introduces counterfactual analysis to the 6D pose estimation task in space and proposes the CA-SpaceNet framework. By imagining an image without the target (*i.e.*, Side Effect), the CA-SpaceNet can decouple the pure target features (*i.e.*, Total Direct Effect) from the raw features (*i.e.*, Total Effect) through counterfactual analysis to obtain more accurate estimation results.

To address these challenges, this paper proposes a Counterfactual Analysis SpaceNet (CA-SpaceNet) framework to handle complicated background information in aerial images. As demonstrated in Fig. 1(c), the CA-SpaceNet introduces counterfactual analysis to the 6D pose estimation task and constructs factual and counterfactual paths. In the factual path, the whole image will be sent to the F_Model to complete feature extraction and results in the factual features (*i.e.*, Total Effect). In the counterfactual path, an image without target but the background is imagined. This non-existent image will be sent to the CF_Model to complete feature extraction and results in the counterfactual features (*i.e.*, Side Effect). With the power of causal inference [6], [7], the CA-SpaceNet can remove the harmful background interference from factual features and generate accurate pose results, which cannot be easily identified by traditional methods. Secondly, to fill the gap in actual deployment on the low-power consumption hardware of the 6D pose estimator, this paper quantizes the CA-SpaceNet into a low-bit-width model

and explores the impact of quantizing different modules on the final performance. A part of the quantized network is implemented on FPGA. Extensive experimental results demonstrate the high performance of the CA-SpaceNet. Latency testing on FPGA confirms the efficiency of low-bit-width quantization and the accelerator architecture.

The main contributions of this paper are as follows:

- We propose a framework named CA-SpaceNet, which is robust to the interference of complicated background information by introducing counterfactual analysis to the 6D pose estimation task in space.
- Our approach outperforms state-of-the-arts on the challenging SwissCube dataset and achieves competitive results on the SPEED dataset.
- We quantize the CA-SpaceNet into a low-bit-width model and deploy a part of the quantized network into a Processing-In-Memory (PIM) chip on FPGA. Low latency proves the feasibility of our method.

As far as we know, it is the first time that the causal inference method and network quantization are explored to address the 6D pose estimation task in space. Robust performance and high efficiency confirm the effectiveness of our method and deployment.

II. RELATED WORK

6D Pose Estimation in Space: Monocular-based 6D pose estimation is a fundamental task in computer vision. According to stages, these methods can be roughly divided into two categories: two-stage and one-stage methods. Two-stage methods [8], [9], [10], [11] complete the keypoints detection firstly (usually adopt corners of the 3D object bounding box) and then solve the 6D pose by the 3D-to-2D correspondences through a PnP solver [12]. There is no keypoint detection process in one-stage methods [13], [14], [15]. The pose information of the object will be transformed into unit quaternion and 3-D translation vector, and the model directly regresses these parameters.

6D pose estimation of space-borne targets plays an important role in satellite on-orbit services and on-board vision-based navigation systems [16], [17]. Compared with terrestrial applications, strict navigation and restricted computation resources put forward higher requirements for the pose estimation model in space. The space engineering community has explored this problem. Traditional methods [1], [18], [19] first find an initial state, *i.e.*, *a priori*, and then use the iterative algorithms to solve the best pose solution that minimizes a specific error criterion pose via hand-crafted feature points. D’Amico *et al.* [2] and Sharma *et al.* [20] proposed some special methods of hand-crafted features to avoid the provision of the initial state based on authentic images captured during the PRISMA mission [2], [21]. Although a series of improvements increase the performance, there is still a huge gap between these optimization-based methods and ideal models.

With the proposal of some large-scale datasets in 6D pose estimation in space [3], [14], [4], some deep neural networks (DNNs) based methods [22], [5], [14], [4] are proposed.

Most of these methods directly modify DNNs to the space scene and do not consider the intrinsic characteristics of space tasks. However, Hu *et al.* [4] considered the extensive depth range and proposed the WDR model, which achieved superior performance on the proposed SwissCube dataset. Inspired by [4], the proposed CA-SpaceNet aims to reduce the interference of complex backgrounds and obtain unbiased pose estimation results through a counterfactual analysis strategy.

Counterfactual Analysis: Counterfactual analysis originates from psychology, which explores that human beings have the ability to evaluate outcomes that did not occur but could have occurred under different conditions [7]. As a powerful way for testing cause-and-effect relationships, counterfactual analysis has been widely used in politics, economics, and epidemiology [23], [24], [25], [26]. Recently, the computer vision community has paid more attention to the application of counterfactual analysis in many visual tasks such as long-tailed visual recognition [27], action anticipation [28], scene graph generation (SGG) [29], and visual question answering (VQA) [30]. Zhang *et al.* [28] presented a counterfactual analysis framework for the egocentric action anticipation (EAA) task. Through the construction of factual and counterfactual cases, side effect caused by semantic labels is reduced, which leads to accurate action anticipation results. The utilization of counterfactual analysis in these methods can improve the performance and interpretability of the model simultaneously. This paper alleviates the problem that the 6D pose estimation model is easily affected by the background through counterfactual analysis and improves the stability of the model.

Low-bit-width Quantization for DNNs: Although DNNs have achieved excellent results in various tasks, the vast computational cost hinders the deployment of these models. Researchers have to trade off the performance and the cost of deployment, especially in special scenes where the computing resources are strictly limited. Much work has explored the lightweight of DNNs, such as network pruning [31], [32], knowledge distillation [33], [34], and quantization [35], [36]. DNNs run with low precision operations during inference provide power and memory advantages over full precision, and it also benefits low-bit-width artificial intelligence chip design [37], [38]. The main idea of quantization is to map full precision floating-point numbers to lower precision (8-bit or lower) through a quantizer to significantly reduce the amount of floating-point operations (FLOPs) in matrix multiplication. Most of the existing methods only explore quantization algorithms in the classification task. In this paper, the proposed CA-SpaceNet is quantized by LSQ-Net [35] and deployed in a low-bit-width PIM accelerator.

III. METHOD

A. Overview

The network architecture is given in Fig. 2. The network’s input is two images: the first one is the raw image I with satellite, and the second is the image I_m with only the background after removing the satellite. Two DarkNet-53

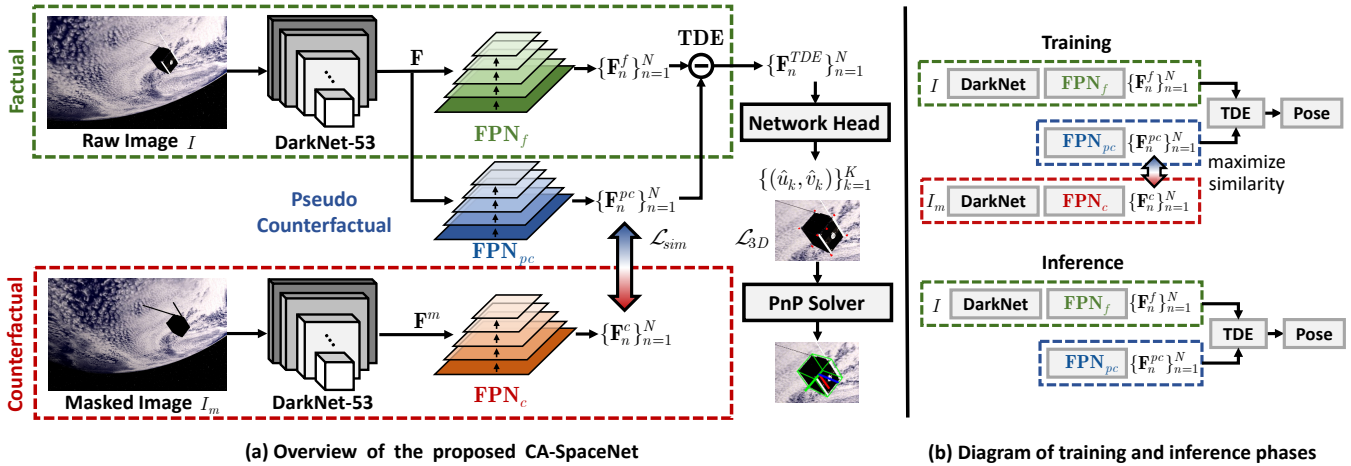


Fig. 2. (a) The CA-SpaceNet consists of five stages: 2) Three independent feature pyramid networks FPN_f , FPN_{pc} and FPN_c with the same structure complete the feature aggregation. For clarity, different colors are assigned to the three paths. 3) Unbiased feature $\{\mathbf{F}_n^{TDE}\}_{n=1}^N$ is obtained by counterfactual analysis. 4) The keypoint detector regresses the 2D projections of the 3D corners of the satellite’s cubic body. 5) Finally, the PnP solver is utilized to calculate the 6D pose of the target satellite through 2D-3D correspondences. (b) To clearly explain the differences between the training and inference phases, we ignore unnecessary feature arrows. In the training phase, FPN_{pc} will imitate FPN_c by maximizing the similarity between $\{\mathbf{F}_n^{pc}\}_{n=1}^N$ and $\{\mathbf{F}_n^c\}_{n=1}^N$, while the whole counterfactual path will be removed during inference.

networks with the same weights are adopted to perform features extraction of I and I_m , respectively, resulting in $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$ and $\mathbf{F}^m \in \mathbb{R}^{C \times H \times W}$. After feature extraction, three feature pyramid networks FPN_f , FPN_{pc} , and FPN_c are constructed to perform counterfactual analysis. These feature aggregation modules with the same network structure but different weights are the core components of **factual path**, **pseudo counterfactual path**, and **counterfactual path**. Through counterfactual analysis, the side effect can be decoupled and removed from the total effect to obtain the final TDE, *i.e.*, unbiased feature $\{\mathbf{F}_n^{TDE}\}_{n=1}^N$ after weakening background interference, where N denotes the number of layers in FPN. Finally, the unbiased feature will be sent to an anchor-based keypoint detector. A PnP solver is adopted to predict the final 6D pose of the target satellite.

The rest of this section is organized as follows: In subsections III-B, III-C, and III-D, the construction of the factual path, counterfactual path, and pseudo counterfactual path are described in detail, respectively. The network quantization method adopted in this paper is briefly reviewed in subsection III-E. The training and inference processes of the CA-SpaceNet are introduced in subsection III-F.

B. Factual Path

It can be seen from the ranking of SPEED competition that the methods based on PnP solver are much more stable than the methods of directly estimating 6D pose. Therefore, this paper adopts the strategy based on a PnP solver. In this strategy, the 6D pose estimation task is divided into two subtasks: 2D keypoint detection and PnP problem. Detailed compositions of the proposed framework are shown in Fig. 2. The factual path is the central part of the CA-SpaceNet, which refers to the structure of [4]. Hu *et al.* [4] explored the problem of huge changes in the depth range of space-borne objects. However, they directly adopted neural networks to

extract features of the whole image without considering the impact of complex background information on 6D pose estimation tasks. In some general computer vision tasks, *e.g.*, object detection and semantic segmentation, background interference will not cause a significant decrease in performance. While in some tasks requiring high precision, such as 6D pose estimation, the background interference will affect the accuracy of keypoint detection, which will significantly deteriorate the final performance.

The factual path is designed to simulate the phenomenon of background interference. In this path, FPN_f completes feature aggregation and generates features with the target satellite and irrelevant background:

$$\mathcal{F}^f = \text{FPN}_f(\mathbf{F}), \quad (1)$$

where \mathcal{F}^f denotes the factual feature set $\mathcal{F}^f = \{\mathbf{F}_n^f\}$, ($n = 1, 2, \dots, N$) generated by different layers of FPN_f . This feature is regarded as the total effect (TE) in counterfactual analysis. The total direct effect (TDE) in CA-SpaceNet is replaced with TE when analyzing the factual path separately. \mathcal{F}^f will be directly sent to the network head to perform keypoint detection, resulting in $\{(\hat{u}_k, \hat{v}_k)\}$, ($k = 1, 2, \dots, K$), where K denotes the number of the corners of the satellite’s cubic body. 3D loss \mathcal{L}_{3D} and object class loss \mathcal{L}_{cls} are adopted in this paper. Please refer to [4] for more details about these loss functions.

C. Counterfactual Path

The idea of counterfactual analysis is to imagine a non-existent path, that is, to study the effect under the *What If* scenario. In space scenes, the complex satellite-earth relationship and harsh illumination conditions will cause significant changes in the background. These factors will negatively impact the feature extraction stage and eventually

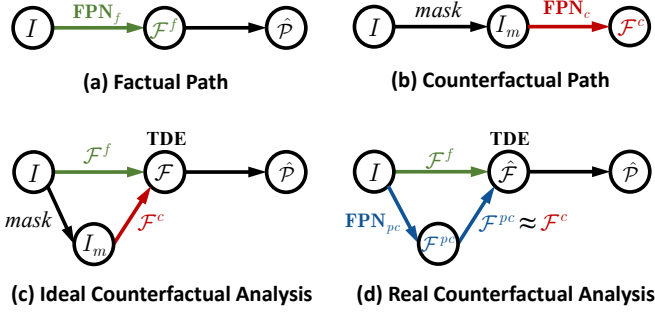


Fig. 3. Simplified causal graphs of CA-SpaceNet in four situations. These causal graphs consist of four types of nodes: image node, feature node, TDE node, and pose results node. Consistent with Fig. 2, different colors are assigned to different elements. The causal graphs of the factual and counterfactual paths are shown in (a) and (b). The difference between the ideal (c) and the real (d) situation is caused by the unavailable masks during the inference phase. For clarity, \mathcal{F} refers to \mathcal{F}^{TDE} .

lead to suboptimal results. Therefore, we imagine what features will be generated through "what if there is no target?" in the counterfactual path. A path composed of a DarkNet-53 and an \mathbf{FPN}_c is constructed to realize this assumption. The path in red box of Fig. 2 shows more details. The input of the counterfactual path is I_m with only background information after erasing the target through the ground-truth mask. Due to the absence of the target, the generated feature map only contains background information:

$$\mathcal{F}^c = \mathbf{FPN}_c(\mathbf{F}^m), \quad (2)$$

where \mathcal{F}^c denotes the counterfactual feature set $\mathcal{F}^c = \{\mathbf{F}_n^c\}, (n = 1, 2, \dots, N)$.

Simplified causal graphs of CA-SpaceNet are shown in Fig. 3. A causal graph is a directed acyclic graph (DAG) that consists of nodes and directed edges. The nodes denote variables, and the directed edges denote cause-effects between nodes [7], [6]. Counterfactual analysis in CA-SpaceNet aims to disentangle the pure object features from the mixed features. In Fig. 3, the factual path (a) and counterfactual path (b) constitute the ideal counterfactual analysis (c). The total direct effect feature \mathcal{F} in (c) is obtained by subtracting the side effect feature \mathcal{F}^c from the total effect feature \mathcal{F}^f :

$$\mathcal{F} = \mathcal{F}^f - \mathcal{F}^c. \quad (3)$$

It should be noted that two DarkNet-53 networks in the counterfactual path and factual path share the same weights and will be frozen during training of the CA-SpaceNet. The reason for adopting this strategy is to only equip the FPN modules with the ability to distinguish the background and the target to avoid the backbone becoming a confounder.

D. Pseudo Counterfactual Path

Although the unbiased feature \mathcal{F} can be directly calculated by Eq. 3 theoretically, the segmentation information of the target is not available in application. Therefore, we elaborately design a pseudo counterfactual path to imitate the counterfactual feature \mathcal{F}^c , which is colored in blue in Fig. 2(a). As its name implies, *pseudo* means that this path

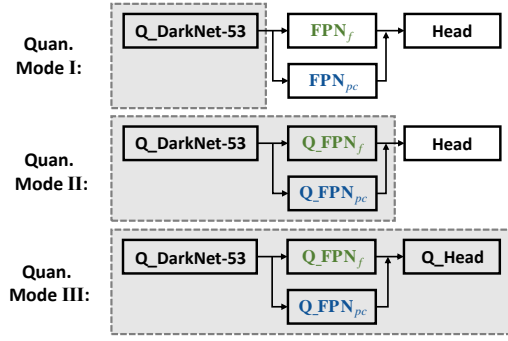


Fig. 4. Three different quantization modes. In order to explore the influence of quantizing different parts on the performance of CA-SpaceNet, three quantization modes are set up: only quantizing the backbone, quantizing the backbone and FPN, and quantizing all modules. Note that all counterfactual paths are removed for clarity.

is a fake path, which aims to imitate the counterfactual path. \mathbf{FPN}_{pc} is the main component of the pseudo counterfactual path. It takes the factual feature \mathbf{F} as input and generates imitation feature:

$$\mathcal{F}^{pc} = \mathbf{FPN}_{pc}(\mathbf{F}), \quad (4)$$

where \mathcal{F}^{pc} denotes the pseudo counterfactual feature set $\mathcal{F}^{pc} = \{\mathbf{F}_n^{pc}\}, (n = 1, 2, \dots, N)$.

In order to make \mathcal{F}^{pc} and \mathcal{F}^c as similar as possible, *i.e.*, $\mathbf{F}_n^{pc} \approx \mathbf{F}_n^c$, a smoothed L1 norm loss function $sl_1(\cdot, \cdot)$ is adopted to measure the discrepancy between them:

$$\mathcal{L}_{sim} = \sum_{n=1}^N sl_1(\mathbf{F}_n^{pc}, \mathbf{F}_n^c). \quad (5)$$

With minimizing \mathcal{L}_{sim} , \mathbf{FPN}_{pc} can learn the ability to disentangle the pure background feature from mixed feature \mathbf{F} . As shown in Fig. 3(d), the counterfactual feature \mathcal{F}^c will be replaced by the pseudo counterfactual feature \mathcal{F}^{pc} in real counterfactual analysis. The final approximate total direct effect feature is calculated by

$$\hat{\mathcal{F}} = \mathcal{F}^f - \mathcal{F}^{pc}. \quad (6)$$

This strategy skillfully solves the problem that the ground-truth mask is lacking during inference of the CA-SpaceNet.

E. Network Quantization

In the space scene, special working conditions and limited computing resources put forward higher requirements for the power consumption and latency of the algorithm. This paper adopts the classical LSQ-Net [35] to quantize the proposed CA-SpaceNet to a low-bit-width model and then deploys a 3-bit convolutional layer into a PIM architecture on FPGA. As shown in Fig. 4, three quantization modes are proposed. The quantization range is gradually expanded to explore the impact of quantization on CA-SpaceNet finely. This is the first work that applies network quantization methods to the 6D pose estimation task of space-borne targets.

Fig. 5(a) demonstrates the detailed quantization and fusion process. Next, we will elaborate on a convolutional layer with

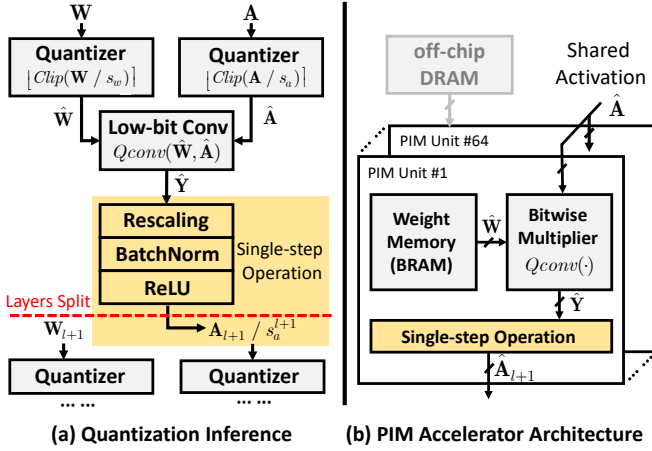


Fig. 5. (a) Detailed quantization inference process. The red dotted line is the separation line of two Conv.-BN-Activ. layers. All computation modules in the yellow box are fused into a single-step operation. (b) Simplified PIM accelerator architecture adopted by this paper. Since the quantized weights have been loaded into BRAMs in 64 separate PIM units, there is no need to read data from the off-chip DRAM, which will cause low efficiency and high power consumption.

the kernel size of 3×3 and the input size of $H' \times W'$ as an example. Before convolution, the weights $\mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in} \times 3 \times 3}$ and activations $\mathbf{A} \in \mathbb{R}^{C_{in} \times H' \times W'}$ in FP32 are quantized into low-bit features $\hat{\mathbf{W}} = \lfloor \text{clip}(\mathbf{W}/s_w) \rfloor$ and $\hat{\mathbf{A}} = \lfloor \text{clip}(\mathbf{A}/s_a) \rfloor$, where $\text{clip}(\cdot)$ returns values within quantization limits, and $\lfloor \cdot \rfloor$ returns every element in $\hat{\mathbf{W}}$ or $\hat{\mathbf{A}}$ to its nearest integer. After the quantization convolution

$$\hat{\mathbf{Y}} = Q_{conv}(\hat{\mathbf{W}}, \hat{\mathbf{A}}), \quad (7)$$

dequantization is performed to recover the activation through rescaling factors:

$$\mathbf{Y} = \hat{\mathbf{Y}} * s_w * s_a, \quad (8)$$

where $\hat{\mathbf{Y}}$ denotes the quantization feature and \mathbf{Y} denotes the recovery feature. Then the convolutional layer and its adjacent BatchNorm layer are fused to a single operation by equivalence relation:

$$\begin{aligned} \mathbf{Y}_{(i, :, :, :)}^{bn} &= \frac{\mathbf{Y}_{(i, :, :, :)} - \mu_i}{\sigma_i} \gamma_i + \beta_i \\ &= \frac{\gamma_i}{\sigma_i} * s_w * s_a * \hat{\mathbf{Y}}_{(i, :, :, :)} - \frac{\mu_i \gamma_i}{\sigma_i} + \beta_i \end{aligned}, \quad (9)$$

where the subscript i denotes the index of the output channel C_{out} . μ , σ , γ , and β are four types of parameters in the BatchNorm layer. Following this, the ReLU layer and scaling step in the next layer are also absorbed by a single operation.

The overall weight memory occupied by the CA-SpaceNet is greatly reduced through quantization and layer fusion, which is suitable for PIM chips with the merits of energy efficiency and avoidance of the memory wall. In PIM architecture, the quantized network is pre-loaded into BRAM, and intermediate data accessed from/to the off-chip DRAM access is entirely eliminated during inference. This paper implements a part of the CA-SpaceNet into the PIM accelerator proposed by Jiao *et al.* [37], which is demonstrated in

TABLE I
COMPARISON WITH STATE-OF-THE-ARTS ON SWISSCUBE.

| Method | Near \uparrow | Medium \uparrow | Far \uparrow | All \uparrow |
|------------------|-----------------|-------------------|----------------|----------------|
| SegDriven [39] | 41.1 | 22.9 | 7.1 | 21.8 |
| SegDriven-Z [39] | 52.6 | 45.4 | 29.4 | 43.2 |
| DLR [5] | 63.8 | 47.8 | 28.9 | 46.8 |
| WDR [4] | 65.2 | 48.7 | 31.9 | 47.9 |
| WDR* [4] | 92.37 | 84.16 | 61.27 | 78.78 |
| CA-SpaceNet | 91.01 | 86.32 | 61.72 | 79.39 |

TABLE II
COMPARISONS OF THE RE-TRAINING WDR* MODEL AND THE CA-SPACE NET.

| Method | Near \uparrow | Medium \uparrow | Far \uparrow | All \uparrow |
|--------------------|------------------|-------------------------|-------------------------|-------------------------|
| WDR* [4] | 92.37 | 84.16 | 61.27 | 78.78 |
| WDR* [4] w. 30-Ep. | 89.93 (-2.44) | 82.09 (-2.07) | 56.50 (-4.77) | 75.76 (-3.02) |
| CA-SpaceNet | 91.01 (-1.36) | 86.32 (+2.16) | 61.72 (+0.45) | 79.39 (+0.61) |

Fig. 5(b). The feasibility of the deployment is confirmed on FPGA.

F. Training and Inference

As shown in Fig. 2(b), the training and inference phases are separated in CA-SpaceNet. All three paths are activated during training. While learning to detect keypoints, the CA-SpaceNet is supposed to minimize the discrepancy between features generated by FPN_{pc} and FPN_c . Therefore, the network is trained with a weighted combination of the loss terms:

$$\mathcal{L} = \lambda_{3D} \mathcal{L}_{3D} + \lambda_c \mathcal{L}_{cls} + \lambda_s \mathcal{L}_{sim}, \quad (10)$$

where the loss weights λ_{3D} , λ_c , and λ_s are set to 1, 1, and 0.25, respectively.

It should be noted that the ground-truth mask is available during training, while unavailable during inference. The pseudo-counterfactual path is constructed to replace the function of the counterfactual path to address this issue. Therefore, the counterfactual path (*i.e.*, FPN_c) will be deleted during inference.

IV. EXPERIMENTS

Comprehensive experiments on SwissCube [4] and SPEED [3] are conducted to evaluate the proposed CA-SpaceNet. The influence of low-bit-width quantization on 6D pose estimation performance is explored. In the end, we deploy a layer of the quantized CA-SpaceNet to a PIM accelerator equipped on an FPGA SoC platform and evaluate the efficiency of the software and hardware co-design system.

A. Datasets and Evaluation Metrics

SwissCube [4]. The SwissCube dataset is a high fidelity dataset for 6D object pose estimation in space scenes. Accurate 3D meshes and physically-modeled astronomical objects are included in this dataset. It contains 500 scenes,

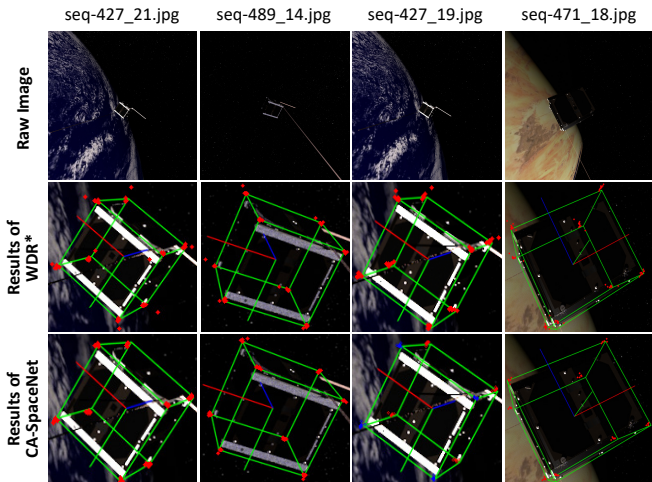


Fig. 6. Visualization of the prediction keypoints generated by the WDR* and the CA-SpaceNet. The ground-truth boxes (in green) and three axes are drawn for clarity. All prediction points are marked in red. The CA-SpaceNet can significantly reduce background interference and generate robust pose estimation results.

of which each scene has 100 image sequences, resulting in 50K images in total. Consistent with [4], 40K images are used for training, and the remaining 10K ones are used for testing.

SPEED [3]. The *Spacecraft Pose Estimation Dataset* (SPEED) was firstly released on the Kelvins Satellite Pose Estimation Challenge in 2019. It contains a large number of synthetic images and a small number of real satellite images. The ground-truth labels of the testing set are not available because the competition is not ongoing. Therefore, we divide the training set into two parts at random, 10K images for training and the remaining 2K ones for testing.

Evaluation metrics. Standard ADI-0.1d [40], [4] accuracy is adopted as the evaluation metric in SwissCube, which represents the percentage of samples whose 3D reconstruction error is less than 10% of the object diameter. The metric $\mathbf{e}_q + \mathbf{e}_t$ is adopted in SPEED, where \mathbf{e}_q is the quaternion error and \mathbf{e}_t is the normalized translation error.

B. Implementation Details

The proposed CA-SpaceNet is built on the PyTorch [41] and implemented on a system with the Intel(R) Core(TM) i7-9700K CPU @ 3.60GHz. All methods are trained on a single Nvidia Titan GPU. For all CA-SpaceNet frameworks in this paper, the DarkNet-53 pretrained on ImageNet is chosen as the backbone. Stochastic gradient descent (SGD) optimizer is adopted for network optimization with initial learning rate $1e-3$, momentum 0.9, and weight decay $1e-4$. The training epoch is set to 30. Online data augmentation strategies such as random shift, scale, and rotation are performed during training. Different experimental settings are adopted because of different complexity. For the SwissCube dataset, the number of minibatch is set to 8, and the resolution is set to 512×512 . For the SPEED dataset, the number of minibatch is set to 4, and the resolution is set to 960×960 .

TABLE III
COMPARISON WITH STATE-OF-THE-ARTS ON SPEED

| Method | $\mathbf{e}_q + \mathbf{e}_t \downarrow$ |
|----------------------|--|
| SLAB Baseline [3] | 0.0626 |
| pedro-fairspace [42] | 0.0571 |
| WDR [4] | 0.0180 |
| WDR* [4] | 0.0400 |
| CA-SpaceNet | 0.0385 |

An Ultra96v2 FPGA board is implemented to deploy the quantization convolutional layer of the CA-SpaceNet.

C. Results on the SwissCube Dataset

The SwissCube dataset is the largest released dataset in 6D pose estimation of space-borne targets. We choose this dataset as the main benchmark. Both quantitative and qualitative experiments are carried out on this dataset.

Quantitative Results. Tab. I compares the CA-SpaceNet with existing methods on the SwissCube dataset. Note that there are two results of WDR model [4]. The results of WDR is obtained by the original paper [4] while the results of WDR* is from our reproduction. The improvement of performance lies in sufficient training epochs. The CA-SpaceNet is obtained by introducing counterfactual analysis strategy to WDR* and performing another 30 training epochs. Therefore, WDR* is regarded as the main competitor of the CA-SpaceNet. In Tab. I, the proposed framework achieves state-of-the-art results on ADI-0.1d (86.32 in *Medium*, 61.72 in *Far* and 79.39 in *All*). Under the *Medium* and *Far* setting, the satellite area in the image is much smaller than the background area, which usually causes suboptimal results. The proposed CA-SpaceNet can eliminate the interference of background through counterfactual analysis, so as to achieve better results. This is also confirmed by subsequent qualitative experiments. Under the *Near* setting, the performance degradation is mainly caused by large masks in these scenes: the large area of the mask leads to the loss of background information, and the counterfactual path can't provide effective background features, resulting in performance degradation.

In order to verify that the performance improvement is brought by the counterfactual analysis strategy rather than the additional 30 training epochs, we conducted another model named WDR* w. 30-Ep.. Experimental results in Tab. II show that the additional 30 training epochs lead to over-fitting, while the additional 30 training epochs with counterfactual analysis (CA-SpaceNet) achieves better performance, which confirms the superiority of the proposed framework.

Qualitative Results. The prediction results of the WDR* and CA-SpaceNet are visualized and compared in Fig. 6. It can be seen from the results of WDR* that the background interference makes the predicted points largely offset from the ground-truth corners. These imprecise corners will lead to large pose estimation errors in the PnP stage. However, with the help of the causal inference method, the CA-SpaceNet successfully handles these complex situations. High quality

TABLE IV

RESULTS ON THREE DIFFERENT QUANTIZATION MODES OF 8-BIT AND 3-BIT CA-SPACENET ON SWISSCUBE

| #Bits | Quan. Mode | ADI-0.1d \uparrow | OPs & FLOPs | Perc.(%) |
|-------|------------|---------------------|------------------------------|----------|
| 8 | I | 76.21 | 36.91 GOPs + 33.79 GFLOPs | 52.21 |
| | II | 75.04 | 44.51 GOPs + 26.19 GFLOPs | 62.96 |
| | III | 74.65 | 70.47 GOPs + 0.23 GFLOPs | 99.67 |
| 3 | I | 75.10 | 36.91 GOPs + 33.79 GFLOPs | 52.21 |
| | II | 74.47 | 44.51 GOPs + 26.19 GFLOPs | 62.96 |
| | III | 68.68 | 70.47 GOPs + 0.23 GFLOPs | 99.67 |

prediction points show that the counterfactual analysis strategy is able to weaken the adverse impact of background interference to the final results.

D. Results on the SPEED Dataset

Tab. III compares our method with several top-performing solutions in the Kelvins Satellite Pose Estimation Challenge. Considering the limitation of computing resources, we didn't adopt multiple model strategies and refinement to the final 6D pose results as done in [5]. There is no official masks of the satellite in SPEED, so we utilize the *cv2.convexHull* function in OpenCV to connect 11 corners to generate approximate masks. As the original paper of WDR [4] did not release the code on the SPEED dataset, we reproduce this model and obtain the results of WDR*. Due to some unknown tricks, the results reproduced by us is slightly worse than WDR [4], but in the same order of magnitude (0.018 for WDR and 0.040 for WDR*). Note that this does not prevent us from evaluating the proposed framework, because the CA-SpaceNet is developed from the WDR*. Compared with the WDR* reproduced by us, the CA-SpaceNet achieves lower error on $\mathbf{e}_q + \mathbf{e}_t$ (0.0385 for CA-SpaceNet and 0.04 for WDR*). The decrease of error proves that the CA-SpaceNet is robust to complex backgrounds interference in 6D pose estimation of space-borne targets.

E. Network Quantization and Deployment

After confirming the effectiveness of the proposed framework, we quantize the CA-SpaceNet into a low-bit-width model and deploy a part of the quantized model into a real hardware accelerator.

Quantization Results. We evaluate the performance of 8-bit and 3-bit CA-SpaceNet on the SwissCube dataset. The performance and operation statistics of these quantized models are listed in Tab. IV. The *Perc.* refers to how many FLOPs in matrix multiplication are converted to low-bit-width OPs. Details of three quantization modes are demonstrated in Fig. 4. Following the common setting of quantization methods, the first layer of DarkNet-53 is kept in FP32. So even in mode III, there are still some floating-point operations that

TABLE V

SUMMARY OF PARAMETER STORAGE SIZE

| Format | #Para. | Model Size | Stor. Saving (%) \uparrow |
|--------|---------|------------|-----------------------------|
| FP32 | 51.29 M | 205.17 MB | 0.00 |
| 8-bit | 51.29 M | 51.29 MB | 75.00 |
| 3-bit | 51.29 M | 19.23 MB | 90.63 |

TABLE VI

MEASURED LATENCY ON DIFFERENT HARDWARE

| Device | Latency (ms) \downarrow |
|-------------------------------------|---------------------------|
| ARM v8.2 64-bit CPU (Nvidia Xavier) | 26.16 |
| Intel Core i7-8700K CPU | 10.25 |
| PIM Arch. on Ultra96v2 FPGA | 5.99 |

cannot be avoided. Tab. IV shows that the performance of 8-bit model and 3-bit model decreases with the increase of quantization range, which is consistent with intuition. Under the 8-bit mode I setting, quantifying the DarkNet-53 can save half of the FLOPs with only reducing ADI-0.1d by 3.18 (79.39 to 76.21). Same setting in 3-bit will reduce ADI-0.1d by 4.29 (79.39 to 75.10). This shows that the quantization strategy can save a large amount of computation without significantly reducing the performance. Under mode III setting, 8-bit and 3-bit models reduce the performance of 4.74 (79.39 to 74.65) and 10.17 (79.39 to 68.68) respectively, which shows that 3-bit quantization will have a great negative impact on the *Network Head* module. In DNNs, deeper layers represent more high-level features. Therefore, the quantization of these modules should be firstly avoided.

Through quantization, floating-point weights in the network are transformed into low-bit-width values for storage, which greatly reduce the size of the network. It is more easier for the network to be deployed to devices with limited memory. Tab. V gives an occupation summary of 8-bit and 3-bit networks under mode III setting. The size of the network is reduced by 75% and 90.63% in 8-bit and 3-bit settings, respectively. The minimal memory occupation provides the basis for the real deployment on chips.

Deployment on the PIM Chip. Due to the limitation of hardware resources of FPGA, we only deploy a single 3-bit quantized convolutional layer of the CA-SpaceNet with the feature map size of $128 \times 128 \times 64$ and kernel size of $128 \times 64 \times 3 \times 3$. In Tab. VI, the latency of PIM architecture [37] on FPGA is compared with the latency on ARM v8.2 CPU and Intel Core-i7 CPU. Note that GPU is not listed because of the low-power consumption setting. The results show that the PIM accelerator achieves the lowest latency (5.99ms) at a clock rate of 100MHz. Our deployment achieves 4.4x speedup compared with ARM v8.2 CPU and 1.7x speedup compared with Intel Core-i7 CPU. Lower latency proves the high efficiency of the low-bit-width quantization and actual deployment.

V. CONCLUSIONS

In order to address the issue that 6D pose estimation in space is vulnerable to background interference, this paper proposes CA-SpaceNet based on counterfactual analysis to weaken the interference of background features from the mixed features. Experimental results show that the proposed framework achieves robust performance. Further, we quantize the CA-SpaceNet into 3-bit and 8-bit and deploy part of the quantized network to a neural network accelerator on FPGA. We believe that our exploration can bring new contributions to the computer vision and space technology community. In the future, we will deploy the entire quantized network to PIM chips to better meet the demands of real space missions.

REFERENCES

- [1] A. Cropp and P. Palmer, "Pose estimation and relative orbit determination of a nearby target microsatellite using passive imagery," in *Proc. of the 5th Cranfield Conf. Dyn. Control Syst. Struct. Space*, 2002, p. 389–395.
- [2] S. D'Amico, M. Benn, and J. L. Jorgensen, "Pose estimation of an uncooperative spacecraft from actual space imagery," *International Journal of Space Science & Engineering*, vol. 2, no. 2, pp. 171–188, 2014.
- [3] M. Kisantal, S. Sharma, T. H. Park, D. Izzo, and S. D. Amico, "Satellite pose estimation challenge: Dataset, competition design and results," *IEEE Trans. Aero. Elec. Sys.*, vol. 56, no. 5, pp. 4083–4098, 2020.
- [4] Y. Hu, S. Speierer, W. Jakob, P. Fua, and M. Salzmann, "Wide-depth-range 6d object pose estimation in space," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021.
- [5] B. Chen, J. Cao, A. Parra, and T. J. Chin, "Satellite pose estimation with deep landmark regression and nonlinear pose refinement," in *Proc. of the IEEE Int. Conf. Comput. Vis. Workshops (ICCVW)*, 2019.
- [6] J. Pearl, M. Glymour, and N. P. Jewell, "Causal inference in statistics: A primer," 2016.
- [7] J. Pearl and D. Mackenzie, "The book of why : the new science of cause and effect," vol. 361, no. 6405, pp. 855.2–855, 2018.
- [8] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab, "Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again," in *Proc. of the IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017.
- [9] M. Rad and V. Lepetit, "Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth," in *Proc. of the IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017.
- [10] B. Tekin, S. N. Sinha, and P. Fua, "Real-time seamless single shot 6d object pose prediction," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018.
- [11] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," in *Robotics: Science and Systems Conference*, 2018.
- [12] C. P. Lu, G. D. Hager, and E. Mjølness, "Fast and globally convergent pose estimation from video images," *IEEE Trans. Pattern Anal. Mach. Intell. (TPAMI)*, vol. 22, no. 6, pp. 610–622, 2000.
- [13] X. Yu, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," in *Robotics: Science and Systems*, 2018.
- [14] P. F. Proenca and G. Yang, "Deep learning for spacecraft pose estimation from photorealistic rendering," in *Proc. of the IEEE Int. Conf. Robot. Autom. (ICRA)*, 2020.
- [15] B. Wen, C. Mitash, B. Ren, and K. E. Bekris, "se(3)-tracknet: Data-driven 6d pose tracking by calibrating image residuals in synthetic domains," in *Proc. of IEEE/RSJ Int. Conf. Intell. Robot. Syst. (IROS)*, 2020.
- [16] S. Sharma and S. D'Amico, "Reduced-dynamics pose estimation for non-cooperative spacecraft rendezvous using monocular vision," in *Proc. of the AAS Guidance, Navigation and Control Conf.*, 2017.
- [17] S.-G. Kim, J. Crassidis, C. Y., A. Fosbury, and J. Junkins, "Kalman filtering for relative spacecraft attitude and position estimation," *Journal of Guidance, Control, and Dynamics*, 2007.
- [18] S. Zhang and X. Cao, "Closed-form solution of monocular vision-based relative pose determination for rvd spacecrafts," *Aircraft Engineering & Aerospace Technology*, vol. 77, no. 3, pp. 192–198, 2005.
- [19] A. Petit, E. Marchand, and K. Kanani, "Vision-based space autonomous rendezvous : A case study," in *Proc. of IEEE/RSJ Int. Conf. Intell. Robot. Syst. (IROS)*, 2011.
- [20] S. Sumant, V. Jacopo, and D. Simone, "Robust model-based monocular pose initialization for noncooperative spacecraft rendezvous," *Journal of Spacecraft and Rockets*, vol. 55, no. 6, pp. 1414–1429, 2018.
- [21] S. D'Amico, P. Bodin, M. Delpéch, and R. Noteborn, "Prisma," in *Distributed Space Missions for Earth System Monitoring*, 2013.
- [22] S. Sharma, S. D'Amico, and C. Beierle, "Pose estimation for non-cooperative spacecraft rendezvous using convolutional neural networks," in *IEEE Aerospace Conference*, 2018.
- [23] B. G. King, "A political mediation model of corporate response to social movement activism," *Administrative Science Quarterly*, vol. 53, no. 3, pp. 395–421, 2008.
- [24] K. Luke, "The statistics of causal inference: A view from political methodology," *Political Analysis*, vol. 23, no. 3, pp. 313–335, 2015.
- [25] L. Richiardi, R. Bellocco, and D. Zugna, "Mediation analysis in epidemiology: methods, interpretation and bias," *International Journal of Epidemiology*, vol. 42, no. 5, pp. 1511–1519, 2013.
- [26] V. Chernozhukov, I. Fernández-Val, and B. Melly, "Inference on counterfactual distributions," *Econometrica*, vol. 81, no. 6, p. 2205–2268, 2013.
- [27] K. Tang, J. Huang, and H. Zhang, "Long-tailed classification by keeping the good and removing the bad momentum causal effect," in *Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2020.
- [28] T. Zhang, W. Min, J. Yang, T. Liu, S. Jiang, and Y. Rui, "What if we could not see? counterfactual analysis for egocentric action anticipation," in *Proc. of the Int. Joint Conf. Artif. Intell. (IJCAI)*, 2021.
- [29] K. Tang, Y. Niu, J. Huang, J. Shi, and H. Zhang, "Unbiased scene graph generation from biased training," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020.
- [30] K. Tang, J. Huang, and H. Zhang, "Counterfactual vqa: A cause-effect look at language bias," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021.
- [31] S. Gao, F. Huang, W. Cai, and H. Huang, "Network pruning via performance maximization," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021.
- [32] Z. Wang, C. Li, and X. Wang, "Convolutional neural network pruning with structural redundancy reduction," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021.
- [33] Y. Zhu and Y. Wang, "Student customized knowledge distillation: Bridging the gap between student and teacher," in *Proc. of the IEEE Int. Conf. Comput. Vis. (ICCV)*, 2021.
- [34] D. Y. Park, M.-H. Cha, C. Jeong, D. Kim, and B. Han, "Learning student-friendly teacher networks for knowledge distillation," in *Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2021.
- [35] S. K. Esser, J. L. McKinstry, D. Bablani, R. Appuswamy, and D. S. Modha, "Learned step size quantization," in *Proc. of the Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [36] K. Yamamoto, "Learnable companding quantization for accurate low-bit neural networks," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021.
- [37] B. Jiao, J. Zhang, Y. Xie, S. Wang, H. Zhu, X. Kang, Z. Dong, L. Zhang, and C. Chen, "A 0.57-gops/dsp object detection pim accelerator on fpga," in *26th Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2021.
- [38] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proc. of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [39] Y. Hu, J. Hugonot, P. Fua, and M. Salzmann, "Segmentation-driven 6d object pose estimation," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019.
- [40] Y. Hu, P. Fua, W. Wang, and M. Salzmann, "Single-stage 6d object pose estimation," in *Proc. of the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020.
- [41] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. Devito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *Adv. Neural Inf. Process. Syst. (NIPS)*, 2017.
- [42] P. F. Proenca and G. Yang, "Deep learning for spacecraft pose estimation from photorealistic rendering," in *Proc. of the IEEE Int. Conf. Robot. Autom. (ICRA)*, 2020.