

Data-Efficient and Hardware Decentralized Visual SLAM

Jincheng Yu¹ and Feng Gao¹

Abstract—Decentralized simultaneous localization and mapping (DSLAM) is essential to a multi-robot system, especially in environments lacking absolute positioning equipments like GPS. Visual based SLAM is a widely adopted solution in industry for its low cost and high flexibility. There are two essential components need to be efficiently deployed on each agent: 1) Visual Odometry (VO) and 2) Place Recognition. However, both of these components require intensive computation and storage on embedded system. The place recognition task is usually done with CNN based methods. We adopt CNN as the VO to provide 6-D pose between different frames, for both intra-robot or inter-robot. Thus we can use the CNN accelerator based on FPGA to execute these two components. In this work, we propose a hardware-software co-design DSLAM framework and use embedded FPGA to accelerator these two components.

We evaluate our framework on the hardware platform Xilinx ZU9 SoC and we can perform DSLAM in real time on each agent. We also evaluate our system on publicly available dataset.

I. Introduction

In recent years, with the development of the hardware and algorithms, the capabilities of a single agent have been greatly improved. To further expand the capabilities of intelligent robots, using several robots can accelerate many tasks, such as localization, exploration, and mapping. As simultaneous localization and mapping (SLAM) is an essential component in many tasks, it is important to do SLAM across different robots in many multi-agent applications. The camera is a widely used sensor in SLAM for its rich information and low cost. However, in many scenarios, communication is limited, so that there is no server or an agent can stably collect all of the visual data from each robot.

Therefore, to reduce communication requirements, the previous work [?] proposes a data-efficient decentralized SLAM (DSLAM) system. The DSLAM frame in [?] is illustrated in fig 1(a). It makes improvements in three typical components in the DSLAM system: 1) Using ORB-SLAM [?] in stereo configuration as the visual odometry (VO) algorithm which provides basic intra-robot pose estimation. 2) Using NetVLAD [?] algorithm to do place recognition which relates the current observation to previous scenes and other robots. 3) Using distributed Gauss-Seidel algorithm [?] as the optimization back-end which optimizes the intra-robot position and fuses the inter-robot locations and maps. Each agent executes the ORB-SLAM which contains three steps for each input frame: feature extraction, feature matching

and RANSAC. The NetVLAD method can encode the camera frame to a short vector which can be transformed to the server or a central agent with low communication cost.

However, both ORB-SLAM and NetVLAD require tremendous computation and storage resources, and thus, the deployment of DSLAM on embedded system is challenged by the limited resources and power supply.

Though NetVLAD consumes huge computation, with the development of FPGA accelerators, we use the embedded CNN accelerator on FPGA [?] to perform NetVLAD for each frame. We also notice that there are also some previous works regression the 6-D pose directly from the input stereo camera [?] or monocular camera [?], [?], With the development of CNN. We adopt Depth-VO-Feat [?] in DSLAM system to estimate the pose from the input monocular camera. Because Depth-VO-Feat is trained with stereo input frames and inferred with monocular camera, the CNN method can provide absolute scale from monocular camera, and also be accelerated with our CNN accelerator. Thus we do not need to execute ORB-SLAM on embedded CPUs.

The proposed DSLAM framework is illustrated in fig 1(b). To make the DSLAM system more energy efficient and hardware friendly, we propose a novel hardware-software co-design DSLAM framework with the following contributions:

- We implement NetVLAD on an embedded SoC platform with CPU cores and FPGA fabric.
- We use an end-to-end CNN based method to estimate the 6-DoF pose between intra-robot successive frames and matched scenes between different robots.
- We demonstrate that our proposed hardware-software co-design decentralized SLAM system can achieve a similar accuracy with the current state-of-the-art DSLAM system without increase of communication.

The rest part of this article is organized as follows. Section II will give the basic idea of CNN based methods and the hardware architecture of embedded FPGA. Section III will detail the implementation of our hardware-software co-design DSLAM system. The experiment result will be given in Section IV. Section VI will conclude this paper.

II. Background and Motivation

A. CNN based methods in DSLAM

As described before, there are two essential components on each agent: 1) Visual Odometry (VO) and 2)

¹Electronic Engineering Department, Tsinghua University, Beijing, China yjc16@mails.tsinghua.edu.cn

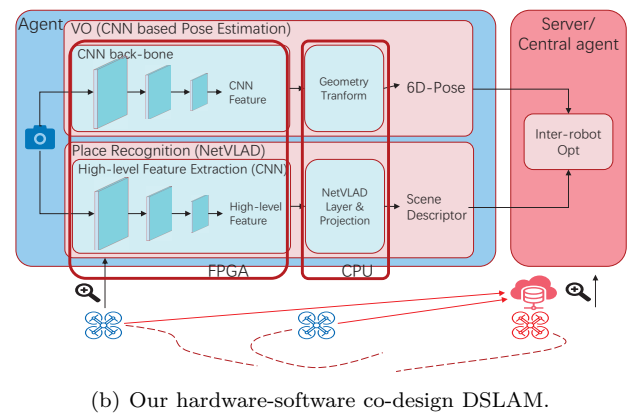
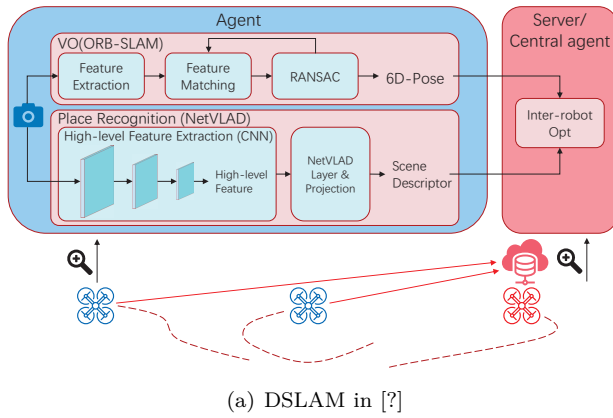


Fig. 1. Overview of the DSLAM in [?] and our hardware-software co-design DSLAM. Each agent (blue drones) will send the result of 6-D pose estimation and scene descriptor to a server or a central agent (red server or drone in figure) to do inter-robot place recognition and optimization. We use CNN instead of feature points to do pose estimation so that we can use CNN accelerator to speed up the whole process.

Place Recognition.

1) Visual Odometry (VO): Visual odometry estimation is the task to infer ego-motion from a sequence of images, and is an essential component in SLAM system. Some feature-based SLAM systems have enjoyed great success, like ORB-SLAM[?] and ORB-SLAM2[?]. Recently, several studies have shown that these feature-based SLAM systems require high computing resources. Fang et al.[?] shows that the feature extraction stage is the most computation-intensive, consuming over 50% of the CPU resources.

As FPGA is one of the most promising platforms as accelerator, the SLAM system on FPGAs has become a hot research topic. However, FPGA-accelerated feature extraction still consumes a lot of time and computing resources, which cannot be deployed simultaneously with an FPGA-accelerated neural network.

2) Place Recognition: The goal of place recognition is to calculate a given frame into a limited set of places. Each places can be encoded as a very short code which can be easy transformed with low communication cost. Traditional place recognition method usually translate the input frame as the aggregation of handcrafted feature point and local descriptors, like SIFT [?] or ORB [?], using vectorization techniques like bag-of-words (BoW) [?] or vector of locally aggregated descriptors(VLAD) [?].

Recent advances in the deep learning and the convolution neural network (CNN) enable powerful end-to-end mode for place recognition [?], [?], and the NetVLAD method is one of the most accurate method based on CNN. The NetVLAD algorithm based on VGG-16 model [?] consumes more than 80G operations for a single 300×300 input image (each operation means an addition or a multiplication). It is very challenging to deploy the NetVLAD on a traditional embedded hardware platform.

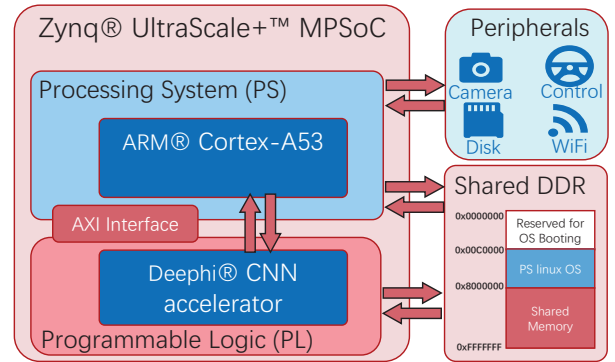


Fig. 2. Hardware architecture of Zynq SoC

B. Hardware architecture of Zynq MPSoC

The Xilinx Zynq MPSoC is a chip with ARM cores and FPGA fabric. The system is illustrated in fig 2. The ARM cores with an embedded Linux operation system are called Processing System (PS). The FPGA fabric is called Programmable Logic (PL). The peripherals like camera and communication unit (WiFi or others) are accessible with PS. The high-bandwidth on-chip AXI interface is used to communicate between PS and PL. PS and PL can also share the DDR to transfer large volume of data such as each frame of camera. Deepphi CNN accelerator [?] is one of the state-of-the-art accelerators and is famous for high energy efficiency on various of CNN structure. We deploy the accelerator on the PL side of Zynq SoC.

Though FPGA can greatly improve the performance and energy efficiency of CNN inference, FPGA cannot efficiently calculate float-point number and requires fixed-point parameters and intermediate data in CNN.

C. Motivation

Though previous work [?] proposes data-efficient DSLAM system, it is difficult to implement the two

essential components, VO and place recognition simultaneously on a communication-limited and energy-constrained embedded hardware platform on a real robot. We propose this hardware-software co-design DSLAM system to use Xilinx Zynq MPSoC and Deephi accelerator to execute these two components on real system.

III. Hardware-Software Co-design DSLAM

Our hardware-software co-design DSLAM system contains two essential improvements in the pose estimation and the place recognition tasks. As illustrate in fig 1(b), both of these two components are divide into two stages: 1) CNN front end to extract features which is deployed to the CNN acclerator on PL and 2) geometric operations to present final results which is deployed on the PS ARM core. To make full use of the Zynq MPSoC (illustrated in fig 2), we optimize the data follow for both of these components.

A. Pose Estimation

We adopt Depth-VO-Feat [?] in DSLAM system to estimate the pose from the input monocular camera. Monocular visual SLAM is a key issue in the field of robotics, while there are two challenging problems: 1) it's difficult and expensive to obtain accurate labeled data. 2) the methods that use monocular sequences in training always suffer from the scale-ambiguity problem, i.e. the actual scale of translations is missing and only direction is learned. In Depth-VO-Feat [?], we use image reconstruction loss as a self-supervised signal to train the convolutional neural networks, and jointly train two networks for depth and odometry estimation without external supervision, which can be used independently in testing phase. Besides, to fix scale-ambiguity issue, we use stereo sequences in training phase and monocular sequences in testing phase. With the known spatial relationship between the left and right cameras, our neural networks can learn the real world scale. Feature reconstruction loss is an additional supervision signal, used to improve the robustness of this framework. And we use depth smoothness loss to encourage the predicted depth to be smooth, which demonstrated success in prior works. Then the final loss becomes

$$L = \lambda_{ir}L_{ir} + \lambda_{fr}L_{fr} + \lambda_{ds}L_{ds}$$

, where L_{ir} , L_{fr} and L_{ds} are image reconstruction loss, feature reconstruction loss and depth smoothness loss respectively, λ_{ir} , λ_{fr} and λ_{ds} are the loss weightings for each loss term. The training framework is illustrated in section III-A.

In order to run efficiently on the FPGA platform, we use fixed-point arithmetic units in the hardware to replace the floating-point number format in GPU and CPU. Many previous works have shown that 8-bit quantization for weights and featuremaps can make the networks run faster on FGPA. Here we adopt the

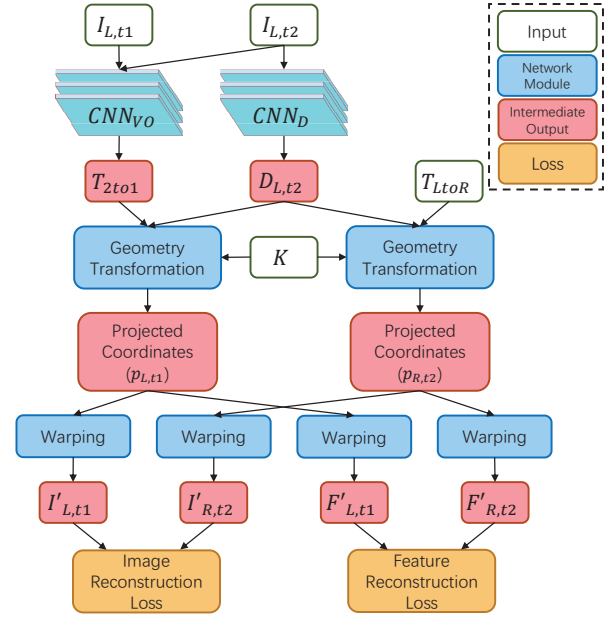


Fig. 3. Illustration of Depth-VO-Feat framework in training phase, where T_{LtoR} is the relative camera pose transformations between right and left views, and K denotes the known camera intrinsic matrix. CNN_{VO} and CNN_D are convolutional nets for visual odometry and depth estimation respectively, which can be used independently in testing phase.

fixed-point finetune method in [?], in that we use the fixed-point number representation in the feed forward phase and keep floating-point number representation for backpropagation, and both weights and data will be re-quantized after each backpropagation. As the fixed-point method will lead to the accuracy loss of the model, We attempt several different quantization strategies to balance speed and accuracy, which will be shown in detail in section IV.

B. Place Recognition

The place recognition method provide the encoded vector transformed to the central agent for inter-robot place matching. As described in section II, CNN has achieved great improvements in place recognition tasks, and NetVLAD [?] is one of the most impressive methods. The computation flow of NetVLAD is illustrated in fig 4. The CNN-based place recognition methods give the global descriptor of a camera frame in a two-step manner: 1) Firstly, a CNN encoder fetches the high-level feature map. 2) A vectorization component that aggregates the feature map into a shot global descriptor. The VLAD layer [?] is a recently proposed plug-and-play operation that greatly improve the performance of place recognition. In the original work with the VLAD layer [?], the feature extraction encoder is a typical CNN VGG-16 [?]. The output dimonsion of original NetVLAD is usually tens of thousands, which is very difficult to stored on embedded system, not to mention in the communication-constrained environment. The PCA

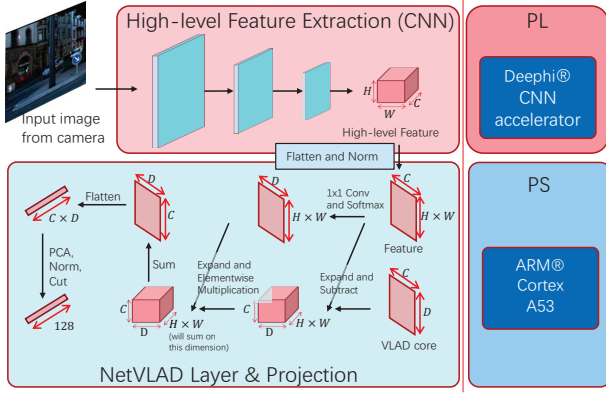


Fig. 4. Process of NetVLAD. The CNN encoder is running at the CNN accelerator on PL side, and the VLAD layer as well as the PCA is running at the ARM core at PS side.

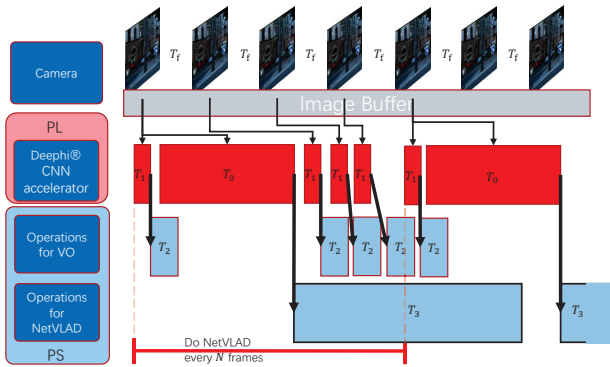


Fig. 5. Scheduling pipeline. There are 4 threads: Camera read, Deepphi core at PL, Operations for VO, and Operations for NetVLAD.

and the projection method can drastically reduce the output dimension. The previous works[?] show that 128 dimension is plenty for DSLAM.

Unlike the fixed-point finetune method used for pose estimation. The training procedure with huge non-public datasets is very complex, and also we cannot finetune the NetVLAD model because of the lack of training data. We simply analyse dynamic range of the weight and intermediate feature map of each CNN layer, and figure out the optimal decimal point position for each layer respectively to minimize the truncation error of each layer. This method is proposed in [?] and is used in many tasks such as image classification and image detection.

C. Parallel Scheduling for VO and NetVLAD

The time consumption of NetVLAD and VO is imbalance. We do pipeline optimization to schedule the two components on Zynq MPSoC. The pipeline is illustrated in fig 5. The interval time for reading camera is T_f . The CNN time for NetVLAD and VO is T_0 and T_1 . The computation time cost on PS for VO and NetVLAD is T_2 and T_3 . We do VO every input frame and do NetVLAD every N frames.

Considering the thread on PL, the time constrain is given as eq. (1).

$$N \times T_f > T_0 + N \times T_1 \quad (1)$$

The thread for VO on PS constrains the NetVLAD frequency as eq. (2).

$$N \times T_f > T_0 + T_1 + (N - 1) \times T_2 \quad (2)$$

The PS part of NetVLAD should finish before computing the PS part of next NetVLAD frame. This constrain can be written as eq. (3).

$$N \times T_f > T_3 \quad (3)$$

The execution time of our design will be given in section IV.

IV. Experiments

The speed the accuracy of our proposed hardware-software co-design DSLAM will be evaluated in this section.

V.

VI. Conclusion

Acknowledgment

This work is not supported by any fund.