

A 400MHz NPU with 7.8TOPS²/W High-Performance-Guaranteed Efficiency in 55nm for Multi-Mode Pruning and Diverse Quantization Using Pattern-Kernel Encoding and Reconfigurable MAC Units

Zhanhong Tan¹, Sia-Huat Tan¹, Jan-Henrik Lambrechts¹, Yannian Zhang², Yifu Wu², Kaisheng Ma¹

¹Tsinghua University, Beijing, China

²IIISCT, Xi'an, China

Deep neural networks present a promising future in applications, ranging from face ID on mobile phones to self-driving cars. Weight pruning and quantization act as valuable solutions to release the burden of computation and memory. Figure 1 shows the family of weight pruning, including the fine-grained and several structural pruning methods. With similar compression rates, coarse-grained pruning results in more accuracy drop. A new structural solution called pattern pruning [5] achieves excellent precision with uniform sparsity rates among kernels, which is friendly to hardware. Kernels are encoded into non-zero values with sparse pattern masks (SPM). This work adopts 16 types of patterns with 4b SPM for the 3x3 convolution, which gains up to 8x compression for eight-zero kernels. As for quantization, the optimal choice generally depends on models.

To support various quantization and pruning methods, including the new pattern granularity, we clarify the following challenges: 1) For pruning, the coarse-grained channel- and filter-level can be easily implemented by the layer width reconfiguration. However, the kernel-level leads to an unbalanced number of kernels in different filters, thus reducing resource utilization. Pattern pruning is a brand-new way that requires an efficient hardware design with a particular data-encoding and computation design that is also compatible with other pruning. However, prior works only focused on a single condition like unstructured sparsity [3] or relatively coarse-grained pruning [4] but did not cover multiple modes. 2) For quantization, resources are expected to be reused as much as possible for linear and nonlinear quantization. UNPU [2] studied various bit-width of weights but not for diversity. Consequently, there lacks a unified architecture facing to optimize diverse types of pruning and quantization.

We introduce a 55nm 400MHz 8.1TOPS/W versatile NPU with the following features: 1) Efficient processing engines (PE) and the input channel scheduler (ICS) supporting diverse pruning via pattern-kernel encoding; 2) An algorithm-enabled data layout for kernel pruning achieving better computation utilization; 3) Unified MAC units supporting 4b weight with linear and nonlinear quantization in a multiplier-free fashion that saves 30.7% power and 3.5% area.

Figure 2 shows the overall NPU architecture and processing flow for pattern and kernel pruning. The proposed architecture mainly consists of 64 PEs to process 64 output channels in parallel. Each PE comprises a sparsity interface (SI) and a pattern-based calculation core (PCC) that exploits the weight and activation sparsity. Each PE is interconnected with an 8-entry input feature map (IFM) FIFO and a 2.3KB weight buffer bank. ICS issues and multicasts input channels to PEs with order indices and masks, only requiring 6.8% extra memory of the 75KB activation buffer. These scheduling orders are generated from a software flow following the group-wise optimized kernel pruning. For the weight pre-processing shown at the bottom-right, the 4b weights are encoded in a pattern-kernel fashion: for pattern pruning, weights within a kernel are compressed in pattern-based format with a 4b SPM and non-zero values; for kernel pruning, kernels in each bank are rearranged to match the input channel order. Two pruning methods are compatible since the pattern-kernel encoding is decoupled by inter/intra-kernel.

Figure 3 depicts the high-performance and pattern-based sparsity-aware PE module with full pipeline design. The kernel in pattern-kernel format is decoded through the pattern decoder and restored to a full 9-weight kernel with a 9-bit mask. The mapping relationship between SPM and weight sparsity mask is stored in a 16-entry LUT reconfigured for each layer. An activation tile is pipelined into the local IFM tile register file, and simultaneously, a sparsity mask for activations is also produced. With sparsity masks, pair-masks for MAC units are generated by the AND logic, followed by a pointer generator to locate all the non-zero activation-weight pairs with the bisection-method. The sparsity rates across kernels are uniform in

pattern pruning. The proposed group-wise kernel pruning attains balanced workloads for multicast. With a more uniform workload distribution profit from the optimized pruning on our hardware, the PE idle time is reduced by 27.3%~35.9%.

Figure 4 presents a high PE-utilization kernel pruning optimization method. The number of pruned kernels in different filters varies randomly for the

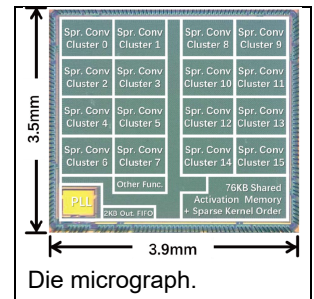
regular kernel pruning, leading to unbalanced workloads. To tackle this problem, we propose a group-wise manner that each group of output channels shares the same set of input channels, and the number of connected input channels of all groups is restricted as the same. Evaluations under three models with 50% group-wise kernel pruning on CIFAR10 show an accuracy drop of less than 2%. To further improve PE utilization, we rearrange scheduling using the flow shown in figure 2, aiming to prevent PEs from being idle for too long. Results show that this policy efficiently skips ineffectual kernels, achieving up to 2.83x performance improvement.

Figure 5 presents the unified MAC unit to maximize resource reuse in different quantization manners. For the linear mode (LIN), we implement the multiplication via accumulating three parts ($1 \times A$, $2 \times A$, $4 \times A$). Two small shifters calculate $2 \times A$ and $4 \times A$. Three weight bits serve as control signals to decide the effectual parts. For the single power-2 quantization (POW), only one shifter is activated to generate $2^0 \sim 2^6$, and we utilize $W[2:0]$ of 3'b111 as the code of zero. For the mixed power-2 quantization (MIX), two small adders leverage (α , β) as adjustment factors to fine-tune two power-2 levels. Both shifters generate results of two parts. Three modes reuse the same 28b accumulator to complete the final results. In summary, this unified MAC unit for 4b weight calculation replaces the traditional decoder and arithmetic multiplier with a few small adders and two shifters, which saves 30.7% power and 3.5% area without performance loss.

Figure 6 shows the measurement results of our NPU implemented in UMC CMOS 55nm technology, operating at 0.75-to-1.00V with 26-to-400MHz frequency. VGG-16 and ResNet-18, which employed different quantization and pruning methods with guaranteed accuracy, are evaluated by the NPU at 1.0V and 300MHz. The average performance and energy efficiency of convolution layers are 104.7-to-792.9GOPS and 1.07-to-3.64TOPS/W. We evaluate the peak energy efficiency using different sparsity scaling and show that this work achieves up to 8.1TOPS/W efficiency. Utilizing 4b weight quantization combined with pattern and kernel pruning, our NPU saves up to 91.7x weight storage. Compared to prior works that usually compromised performance for the peak efficiency [1~4], our NPU achieves the highest performance-guaranteed efficiency of 8.5TOPS²/W (peak-efficiency x the-corresponding-performance). This newly defined metric is profound, aiming to guarantee that the NPU can achieve high performance and high efficiency at the same time. Besides, this work is the first one supporting inter/intra-kernel sparsity and linear/nonlinear quantization in a unified architecture with high performance and efficiency.

References:

- [1] D. Shin, *et al.*, "DNPU: An 8.1 TOPS/W Reconfigurable CNN-RNN Processor for General-Purpose Deep Neural Networks," ISSCC, Feb. 2017.
- [2] J. Lee *et al.*, "UNPU: A 50.6TOPS/W Unified Deep Neural Network Accelerator with 1b-To-16b Fully-Variable Weight Bit-Precision," ISSCC, Feb. 2018.
- [3] C. Lin *et al.*, "A 3.4-to-13.3TOPS/W 3.6TOPS Dual-Core Deep-Learning Accelerator for Versatile AI Applications in 7nm 5G Smartphone SoC," ISSCC, Feb. 2020.
- [4] J. Yue *et al.*, "A 65nm Computing-in-Memory-Based CNN Processor with 2.9-to-35.8TOPS/W System Energy Efficiency Using Dynamic-Sparsity Performance-Scaling Architecture and Energy-Efficient Inter/Intra-Macro Data Reuse," ISSCC, Feb. 2020.
- [5] X. Ma *et al.*, "PCONV: The Missing but Desirable Sparsity in DNN Weight Pruning for Real-Time Execution on Mobile Devices," AAAI, Feb. 2020.



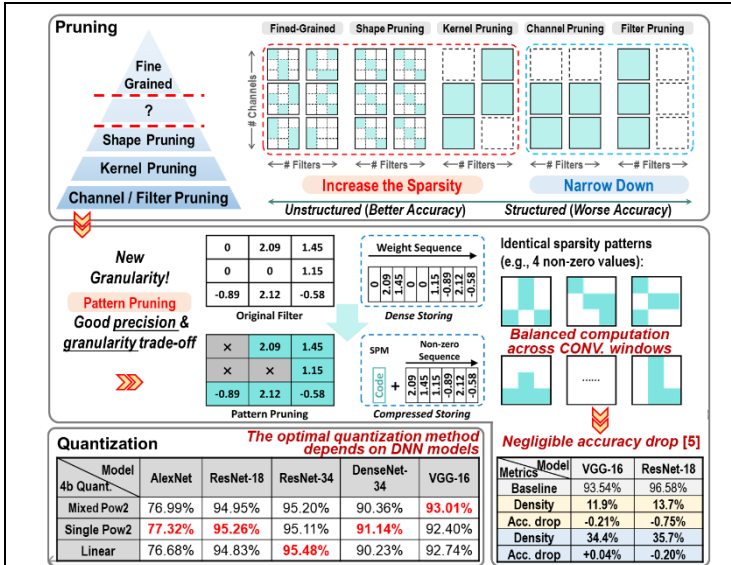


Fig. 1. The pruning family with a new promising pattern pruning and different optimal quantization across models.

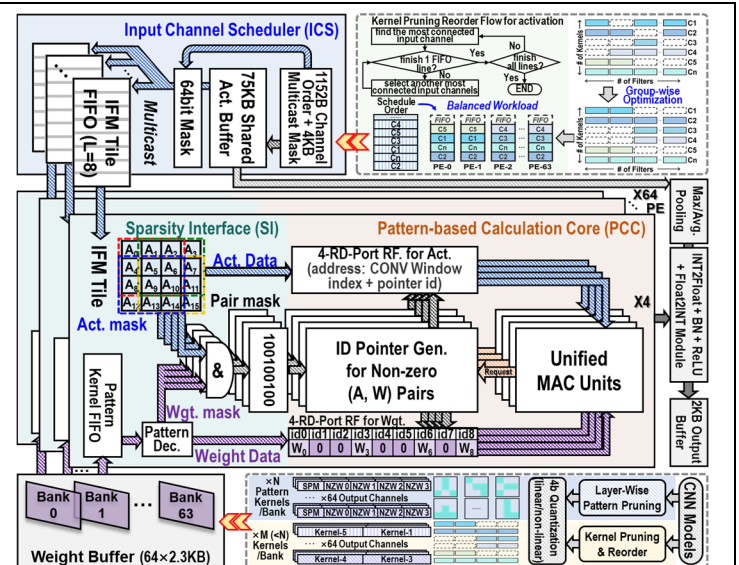


Fig. 2. Overall architecture and weight (bottom-right) & activation (top-right) pre-processing flow for pattern and kernel pruning.

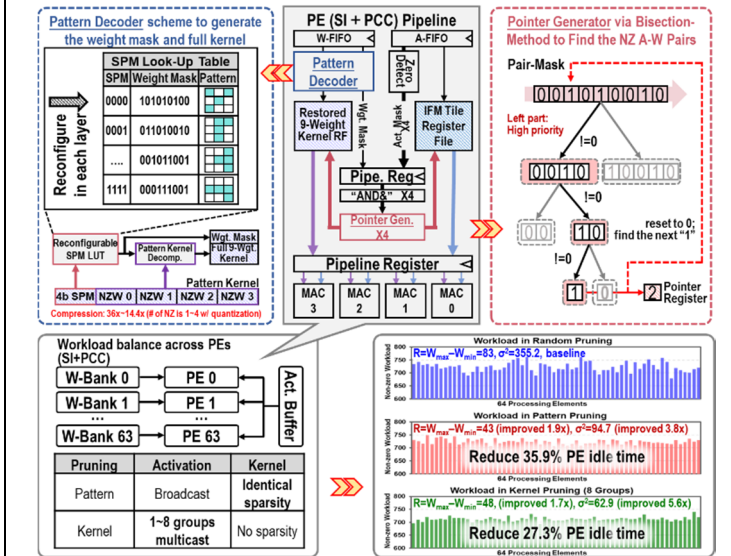


Fig. 3. Sparsity interface (SI) and pattern-based calculation core (PCC) for sparse computation.

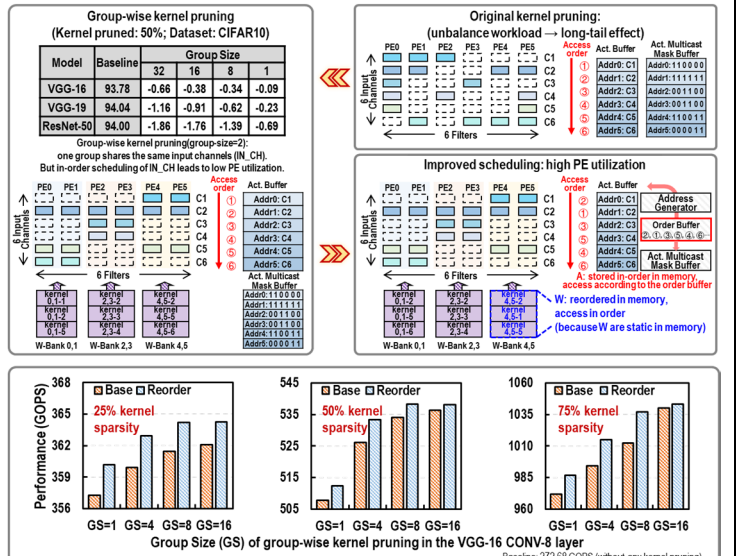


Fig. 4. The input channel scheduler (ICS) scheme for PE utilization optimization via group-wise kernel pruning and reordering.

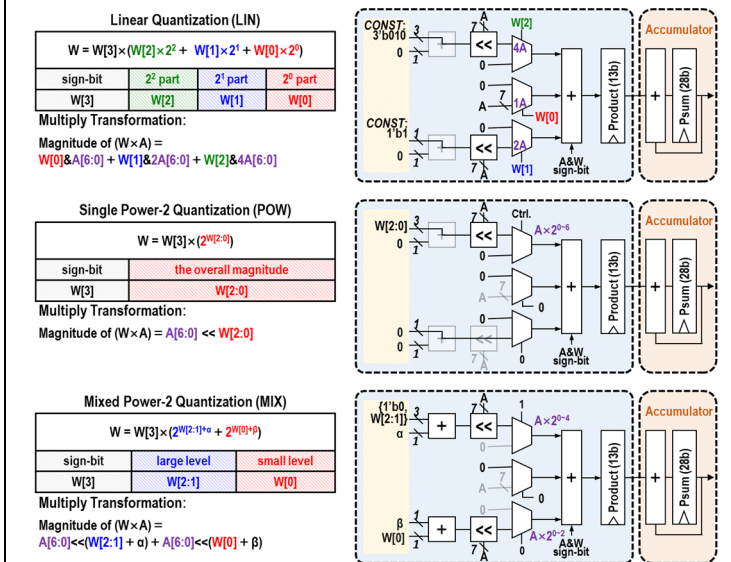


Fig. 5. Unified MAC unit for 4bit linear (LIN), single power-2 (POW), mixed power-2 (MIX) quantization for weights.

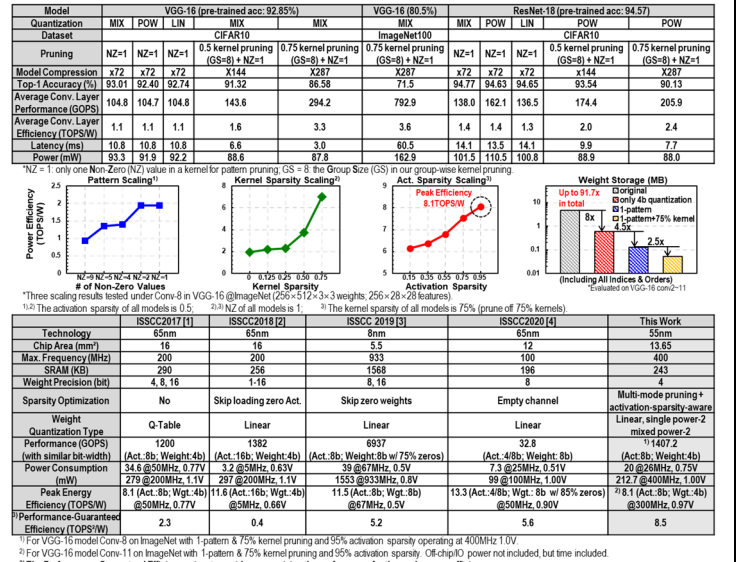


Fig. 6. Measurement results and the comparison table.