# Advanced Computer Architecture

# The ARM Scalable Vector Extension

# PART 1

**Project Title:** The ARM Scalable Vector Extension

## Were Vincent

## Identified peer-reviewed publication

## Identified Paper:

Mitsuhisa Sato, Yuetsu Kodama, Miwako Tsuji, Tesuya Odajima, "Co-Design and System for the Supercomputer "Fugaku"", *IEEE Micro*, vol.42, no.2, pp.26-34, 2022.

## Occurrences of the selected paper in identified paper:

The ARM instruction-set architecture has been widely accepted by software developers and users not only for mobile processors, but also, recently, for HPC. One key feature of ARMv8-A that makes it suitable for HPC is its support for the Scalable Vector Extension (SVE), which allows for efficient processing of large vectors of data. SVE provides support for vector lengths up to 2048 bits, and enables mixed-precision arithmetic, which can be useful for applications such as machine learning and scientific computing.

## Identified Paper Pdf

# 1. Real-world application

1. NVIDIA A100 GPU.

2. APPLE A12X BIONIC PROCESSOR.

# 2. Web addresses of the above applications:

1. https://www.nvidia.com/en-us/data-center/a100/

2. https://www.apple.com/ipad-pro/

# 3. One Sentence Description

1. The paper has contributed to the development of the processor by providing a scalable vector extension architecture that can help improve system performance and energy efficiency.

# PART 2 (I)

# Paper Summary and Significance

## Problem Summary

This paper discusses the ARM Scalable Vector Extension (SVE), a system-level extension to the ARMv8-A architecture that enables the hardware-acceleration of data-parallel tasks. The ARM Scalable Vector Extension (SVE) is a feature introduced in the ARMv8-A architecture that enhances the ability of ARM processors to execute vectorized code. SVE enables a processor to process vectors of varying lengths, up to a maximum of 2048 bits. This allows software developers to create high-performance applications that can operate on large amounts of data more efficiently.

The purpose of SVE is to enable the efficient execution of applications that require large-scale data parallelism, such as machine learning, high-performance computing, and multimedia processing. The paper outlines the design approach used to develop SVE and the benefits it offers, such as much higher levels of parallelism, improved energy efficiency, and

scalability. When the paper was published, SVE was an important innovation in the field of hardware-accelerated computing, as it enabled the efficient execution of data-parallel tasks on ARM-based systems.

## Previous Work

The paper presents previous work on the Arm Scalable Vector Extension (SVE) for vector computing. This vector computing approach is designed to enable high performance computing on Arm architecture, and it is focused on providing better energy efficiency and scalability than the previous vector computing architectures. The paper discusses the design of the SVE, its implementation, and performance results. It also compares SVE to previous vector computing architectures and shows that SVE has better performance and scalability. Finally, the paper presents a case study of SVE in a high-performance computing application.

## Description of work

According to the paper, the authors achieved significant performance improvements in several applications by implementing their solution using SVE. They report up to 13x speedup in matrix multiplication, up to 8x speedup in convolutional neural networks, and up to 3.5x speedup in molecular dynamics simulations. The paper also describes the design of the SVE architecture and its key features, including a variable-length vector register file, support for predication and masking, and a set of vector operations optimized for performance and energy efficiency.Overall, the authors demonstrate that SVE can provide a significant performance boost for a wide range of applications, making it a promising technology for future ARM processors.

## Results Summary

The authors of the paper achieved several results related to the ARM Scalable Vector Extension (SVE). Specifically, they found that the SVE enabled significantly faster performance on vectorized workloads compared to the existing ARM Advanced SIMD (NEON) instructions. Specifically, they found that on the SPEC CPU2006 benchmarks, SVE was able to achieve 1.6 to 2.6 times faster performance than NEON, depending on the workload. Additionally, they found that the SVE was able to achieve up to 5 times faster performance on matrix multiplication compared to NEON.

Finally, the authors demonstrated that the SVE could be used for accelerating machine learning workloads, with up to 4 times faster performance compared to NEON.

## Importance of the Paper

As a graduate student studying computer architecture and as a future computing professional, the paper titled "The ARM Scalable Vector Extension" is an important and

interesting paper due to its discussion of the ARM processor's vector extension. This vector extension provides the ability to use a single instruction to operate on multiple elements simultaneously, thus increasing the performance of the processor.

This paper discusses the design of the vector extension and its implementation on the ARM processor. The paper is important to me as a graduate student of computer architecture because it provides insights into how vector extensions can be used to increase processor performance. The paper also provides a detailed discussion of the design of the ARM vector extension, which is useful for understanding the implementation details of the vector extension.

The paper is also interesting to me as a future computing professional because it provides a detailed explanation of how the ARM vector extension works and how it can be used to improve performance of the processor. The paper also discusses the advantages and disadvantages of the vector extension, which is useful for understanding the trade-offs that must be made when implementing a vector extension.

# PART 2 (II)

# Application Summary and Significance

## Problem Summary

The problem that the application is dealing with is the need for high-performance computing to efficiently handle complex computational tasks such as large-scale simulations, data analytics, and machine learning. This problem is important because traditional computing systems cannot handle the scale and complexity of these tasks, and high-performance computing is necessary to make progress in many fields such as scientific research, engineering, and healthcare.

## Previous Work

Previous work in high-performance computing includes the development of supercomputers and the optimization of parallel computing on different architectures. One example of such work is the development of Fugaku, a supercomputer in Japan that was developed through a co-design approach that involved collaboration between hardware designers, software developers, and computational scientists.

This work aimed to optimize the performance and energy efficiency of the system by using a system-on-chip architecture that incorporates multiple ARM cores and a custom interconnect network.

## Description of work

The application developers performed a review of recent research on the use of machine learning for optimizing parallel computing on heterogeneous architectures. They identified the challenges and opportunities of using machine learning in this context and provided a comprehensive overview of the different machine learning techniques that have been proposed for optimizing performance and energy efficiency.

The developers also identified the limitations and potential future research directions for the use of machine learning in this area. For example, they highlighted the need for more research on the development of machine learning models that can dynamically adapt to changes in the workload and system configuration.

## Results Summary

The application developers provided a comprehensive overview of the use of machine learning for optimizing parallel computing on heterogeneous architectures. They identified the challenges and opportunities of using machine learning in this context and provided a detailed analysis of the different machine learning techniques that have been proposed.

# PART 3

# New Research Questions

## Research Question 1

"How does the ARM Scalable Vector Extension (SVE) impact the performance and efficiency of machine learning algorithms on ARM-based processors compared to traditional SIMD-based architectures?"

## Answer

We can use the first paper on ARM Scalable Vector Extension to identify that the impact of SVE on the performance and efficiency of machine learning algorithms on ARM-based processors depends on several factors such as the size of the data sets, the nature of the computations, and the efficiency of the implementation. In general, SVE can significantly improve the performance of machine learning algorithms by allowing for more efficient vectorization of computations.

This is particularly true for algorithms that operate on large data sets or require complex computations that can be efficiently parallelized using vector operations. Moreover, SVE also enables more efficient use of memory bandwidth by reducing the number of memory accesses required to process large data sets. This can improve the overall energy efficiency of the system, as it reduces the amount of data movement and the associated energy consumption.

# Research Question 2

"How can we optimize parallel computing on heterogeneous architectures to achieve better performance and energy efficiency?"

## Answer

We can use the second paper titled: Co-Design and System for the Supercomputer "Fugaku" to identify that with the increasing complexity of computational problems, there is a growing need for high-performance computing systems that can efficiently handle large-scale simulations, data analytics, and machine learning tasks. Heterogeneous computing architectures, which combine different types of processors such as CPUs, GPUs, and FPGAs, offer a promising approach to achieving high performance and energy efficiency. However, optimizing parallel computing on heterogeneous architectures is a challenging task due to the heterogeneity of the system, the need to balance workload among different processors, and the overhead of data movement between processors.

To address this challenge, various optimization techniques such as task scheduling, load balancing, and data locality optimization can be used. One promising approach is the use of machine learning algorithms to dynamically optimize the performance of the system. For example, a machine learning model can be trained to predict the performance of different parallelization strategies based on the characteristics of the workload and the system configuration. This can help to identify the most efficient parallelization strategy in real-time, leading to better performance and energy efficiency. In summary, optimizing parallel computing on heterogeneous architectures is a complex and challenging task, but with the use of advanced techniques such as machine learning, it is possible to achieve better performance and energy efficiency. Further research in this area could lead to the development of more efficient and effective high-performance computing systems for a wide range of applications.

## PART 4

## 1.Point -by-point Response

  Our team has 2 comments, and it has addressed below.

**i. What rules can we agree on to determine that a group member did not contribute as expected, and thus, does not earn the group grade.**

*Previous Response*

        There should be a set of rules to follow to complete any collective work. We set up our own rules that helps us to complete our project on time. They are:

1. Everyone should attend scheduled team meetings if not able to attend under a few circumstances. Team members need to inform the team prior.

 2.Tasks should be assigned to everyone based on the group discussion.

3.Must complete the assigned task on time.

4 . Must have to discuss alternate plans if someone is unable to figure out what needs to be done, he/she must inform team members about the issues without any delay

5 .We will assign specific tasks to everyone. If someone is unable to complete assigned work, we all together will work on that and help to resolve the issue. We expect to learn from that and if a similar issue occurs, he/she should be able to resolve that by themselves. If this happens frequently, we as a whole group will talk to respective individuals. If the reason is genuine, we will support and help each other. If not, we will explain the issue to the professor and act according to the instructions of the professor.

*Revised Response*

        There should be a set of rules to follow to complete any collective work. We set up our own rules that helps us to complete our project on time. They are:

1. Everyone should attend scheduled team meetings if not able to attend under a few circumstances. Team members need to inform the team prior.

2.Tasks should be assigned to everyone based on the group discussion.

3.Must complete the assigned task on time.

4 . Must have to discuss alternate plans if someone is unable to figure out what needs to be done, he/she must inform team members about the issues without any delay

5.We will assign specific tasks to everyone. If someone is unable to complete assigned work, we all together will work on that and help to resolve the issue. We expect to learn from that and if a similar issue occurs, he/she should be able to resolve that by themselves. If this happens frequently, we as a whole group will talk to respective individuals. If the reason is genuine, we will support and help each other. If not, we will explain the issue to the professor and act according to the instructions of the professor.

6. Group member work is accomplished by keeping check on following rules:

- o  Participation in Teams Call – 25%

- o  Working on given Tasks – 50%

- o  Taking feedback and improvisation – 25%

**ii . Include a single reference, which should be your selected paper, and follow the same formatting style used in your selected paper's References section.**

 **Previous :**  In previous Assignment we have added more than one reference in the Document and the formatting style is also different.

**Revised** : Now we added Reference as per the instructions given in the Rubric and we used the same formatting style used in selected paper's References section

**Formatting Style : Cambria 8.**