

## **Bit-Level Parallelism**

Loic Niragire

School of Technology, Northcentral University

TIM 8121: Distributes Algorithms and Parallel Computing

Dr. Carol Cusano

April 9<sup>th</sup>, 2023.

The rapid evolution of computing technology has led to an ever-growing demand for higher-performance processors capable of handling complex tasks in various fields, such as artificial intelligence, big data processing, scientific simulations, and real-time applications. Bit-level parallelism, which refers to the simultaneous processing of multiple bits of data within a single operation, has significantly enhanced computing performance by increasing the efficiency of data manipulation and transfer within processors. As a result, understanding the latest developments in bit-level parallelism and their implications for computing systems is essential for researchers, engineers, and industry professionals working in this area.

The present study examines recent trends in bit-level parallelism, encompassing the progression toward larger word sizes, the utilization of specialized processors like GPUs and TPUs, and novel processor architectures, and evaluates their impact on overall computing performance in different application domains. The paper also addresses these trends' challenges and potential drawbacks, including increased complexity, power consumption, and fabrication challenges.

### **Background and Historical Context**

Bit-level parallelism improves throughput by increasing the number of bits a processor can operate on in a single operation. A 32-bit processor, for instance, can operate on 32 bits of data in a single operation, which has several implications for the design of the processor. From a data processing perspective, the arithmetic and logic unit (ALU) must be designed to perform arithmetic, logic, and bitwise operations on 32-bit operands. Moreover, the processor's registers must store, manipulate, and transfer 32-bit data values efficiently. Meanwhile, the processor's instruction set architecture (ISA) must be designed to work with

32-bit data values and memory addresses<sup>1</sup>. However, the actual throughput depends on other factors, such as processor architecture, clock speed, and memory hierarchy. Processing more data simultaneously allows the processor to quickly complete tasks (Patterson & Hennessy, 2013).

The performance gains associated with increased bit-level parallelism can be calculated using metrics such as instructions per second (IPS) or floating-point operations per second (FLOPS). For example, if a 32-bit processor can execute 1 billion instructions per second (1 GIPS) and processes 32 bits of data per instruction, it can process 32 billion bits per second (32 GBPS). In contrast, a 64-bit processor with the same GIPS can process 64 billion bits per second (64 GBPS), doubling the data throughput. Table 1 highlights processors' historical progression in terms of word length.

Table 1 Processor's historical progression in terms of their word length

| Processor     | Year introduced | Word length |
|---------------|-----------------|-------------|
| Intel 4004    | 1971            | 4-bit       |
| Intel 8080    | 1974            | 8-bit       |
| Intel 8086    | 1978            | 16-bit      |
| Intel 80386   | 1985            | 32-bit      |
| Intel Itanium | 2001            | 64-bit      |
| AMD64         | 2003            | 64-bit      |

<sup>1</sup> A 32-bit processor can generate 2<sup>32</sup> unique memory addresses, corresponding to an addressable memory space of 4GB.

## Benefits of Bit-Level Parallelism in Computing Performance

### Arithmetic and Logic Operations

Bit-Level parallelism directly affects the efficiency of arithmetic and logic operations in processors. As word size increases, the processor can handle larger data widths in a single operation, reducing execution time for specific tasks. For example, a 64-bit processor can perform arithmetic operations on 64-bit integers in a single instruction, whereas a 32-bit processor would require multiple instructions to handle the same data (Hennessy & Patterson, 2011). The instructions per second (IPS) metric can quantify this performance improvement. If a 32-bit processor requires  $n$  instructions to operate on a 64-bit integer, and a 64-bit processor requires  $m$  instructions, then the performance improvement is determined by Equation (1).

$$\frac{(n - m)}{n} * 100\% \quad (1)$$

### Memory Addressing

Increasing word size also expands the addressable memory space, allowing processors to handle larger data sets and more complex applications (Patterson & Hennessy, 2013). For instance, a 32-bit processor can address up to 4GB of memory, while a 64-bit processor can address up to 16 exabytes. This increase in addressable memory space can be quantified as  $(2^{64} - 2^{32})/2^{32} * 100\%$ , representing a substantial expansion of over 4 billion times. Furthermore, processors can perform tasks more efficiently with more addressable memory, reducing the need for time-consuming disk swapping and other memory management techniques (Crawford, 1998).

## **Recent Trends in Bit-Level Parallelism**

### **Energy-Efficient Processor Design**

Energy efficiency has become a critical factor in processor design, with researchers exploring ways to maximize performance while minimizing power consumption (Esmaeilzadeh, Blem, Amant, Sankaralinam, & Burger, 2011). One approach is to selectively reduce the bit-width of data in specific operations, a technique known as approximate computing (Chippa, 2013). For example, Chippa et al. (2013) demonstrated that by reducing the bit-width of data by 50% in certain arithmetic operations, performance could be improved by up to 2x with minimal impact on application quality. This trend highlights the potential of bit-level parallelism to enhance computing performance while addressing energy constraints.

### **Neuromorphic Computing**

Neuromorphic computing aims to mimic the human brain's structure and function by leveraging highly parallel, lower-power, and fault-tolerant architectures (Merolla, Arthur, & Alvarez-Icaza, 2014). Neuromorphic processors, such as IBM's TrueNorth chip, utilize bit-level parallelism in novel ways. Each neuron processes multiple input spikes concurrently and updates its internal state based on these inputs (Merolla et al., 2014). In their study, Merolla et al. (2014) reported that the TrueNorth chip achieved a 3-orders-of-magnitude improvement in energy efficiency compared to conventional processors while maintaining comparable performance for cognitive computing tasks.

### **In-Memory Computing**

In-memory computing addresses performance bottlenecks associated with data movement between the processor and memory by performing computations directly within memory cells (Zha & Li, 2018). In addition, in-memory computing architectures have been

shown to exploit bit-level parallelism by enabling the simultaneous processing of multiple bits within a single memory cell, leading to performance improvements and reduced energy consumption (Li et al., 2018). For instance, Li et al. (2018) demonstrated that their in-memory computing architecture could achieve a 13.8x speedup compared to conventional processors for matrix-vector multiplication.

### **RISC-V and Customizable Instruction Set Architecture**

The RISC-V architecture, an open-source instruction set architecture (ISA) that supports various word sizes, including 32-bit, 64-bit, and 128-bit, has gained significant interest due to its flexibility, modularity, and extensibility (Waterman, 2023). For example, research by Lu et al. (2020) showed that RISC-V-based processors could achieve competitive performance and energy efficiency compared to commercial ARM and x86 processors. In addition, the customizable nature of RISC-V architecture enables researchers and engineers to design processors with specific word sizes and bit-level parallelism capabilities tailored to their applications' requirements.

## **Challenges and Potential Drawbacks**

### **Power Consumption and Thermal Challenges**

As bit-level parallelism increases, so does the power consumption of processors, leading to thermal challenges that need to be addressed (Borkar & Chien, 2011). In addition, power density increases as more transistors are packed into processors to support higher bit-level parallelism, leading to overheating and reduced reliability (Borkar & Chien, 2011). Consequently, researchers and engineers must find innovative solutions to manage power consumption and thermal dissipation without compromising performance.

## **Diminishing Returns in Performance**

Although increasing bit-level parallelism can enhance performance, there are diminishing returns beyond a certain point (Hennessy & Patterson, 2011). Furthermore, as word size increases, the potential performance gains from bit-level parallelism may be offset by other factors, such as increased complexity, longer instruction pipelines, and greater control overhead (Hennessy & Patterson, 2011). This suggests that increasing bit-level parallelism may only sometimes yield significant performance improvements, and a balanced approach that considers another architectural and system-level optimization is necessary.

## **Software and Compatibility Challenges**

Implementing higher bit-level parallelism often requires software and hardware modifications, which can lead to compatibility issues (Crawford, 1998). For instance, transitioning from a 32-bit to a 64-bit processor necessitates operating systems, compilers, and application software changes to fully exploit the benefits of an increased bit-level parallelism (Crawford, 1998). Addressing these compatibility challenges requires considerable software development and updates investment and can slow the adoption of new bit-level parallelism advancements.

## **Complexity and Cost**

Increasing bit-level parallelism often results in more complex processor designs, increasing manufacturing costs and design challenges (Hennessy & Patterson, 2011). As word size and bit-level parallelism increase, so does the complexity of the underlying circuitry and the required number of transistors. This increased complexity can result in higher manufacturing costs and lower yields, which may offset the potential performance gains from increased bit-level parallelism (Hennessy & Patterson, 2011).

Bit-level parallelism has significantly driven computing systems' evolution and performance improvements over the past few decades. This paper explored the concept of bit-level parallelism, its benefits on computing performance, recent trends, and the challenges and potential drawbacks associated with its implementation. The analysis revealed that bit-level parallelism had played a critical role in enhancing arithmetic and logic operations, increasing addressable memory space, and enabling the development of specialized processors such as GPUs and TPUs.

Recent trends in bit-level parallelism, including energy-efficient processor design and neuromorphic computing, demonstrate the ongoing efforts to leverage this concept for further performance gains, efficiency improvements, and scalability. However, several challenges and potential drawbacks must be addressed, such as power consumption and thermal issues, diminishing returns in performance, software and compatibility challenges, and increased complexity and cost.

Bit-level parallelism remains a vital aspect of computing technology, shaping the development of processors and computing systems that meet the ever-growing demands for performance, energy efficiency, and scalability. Although developing a 128-bit processor is technically feasible, the benefits may not outweigh the challenges and potential drawbacks now. For example, current 64-bit processors can already address a vast memory space (16 exabytes), which is sufficient for most applications. Furthermore, the increased complexity and cost associated with a 128-bit processor might not be justifiable, given the diminishing returns in performance improvement. Nevertheless, as computing demands continue to grow and technological breakthroughs enable more efficient designs, developing a 128-bit processor might become more viable.



## References

- Borkar, S., & Chien, A. A. (2011). The future of microprocessors. *Communications of the ACM* Vol. 54 No.5, 67-77.
- Chippa, V. K. (2013). Analysis and characterization of inherent application resilience for approximate computing. *Proceedings of the 50th Annual Design Automation Conference* (pp. 1-9). ACM.
- Crawford, J. (1998). The development of the Intel IA-64 architecture. *Intel Technology Journal*, 10-20.
- Esmailzadeh, H., Blem, E., Amant, R., Sankaralingam, K., & Burger, D. (2011). Dark Silicon and the End of Multicore Scaling. *The proceedings of the 38th International Symposium on Computer Architecture (ISCA '11)* (pp. 365-376). San Jose, California, USA: ACM.
- Merolla, P., Arthur, J. V., & Alvarez-Icaza, R. (2014). A Million spiking-neuron integrated circuit with scalable communication network and interface. *Science* Vol. 345, 668-673.
- Patterson, D. A., & Hennessy, J. L. (2013). *Computer Organization and Design*. Morgan Kaufmann.
- Waterman, A. S. (2023, April 8). *people.eecs.berkeley.edu*. Retrieved from <https://people.eecs.berkeley.edu/~krste/papers/EECS-2016-1.pdf>
- Zha, Y., & Li, J. (2018). Liquid Silicon: A Data-Centric Reconfigurable Architecture Enabled by RRAM Technology. *Proceedings of the 45th Annual International Symposium on Computer Architecture* (pp. 51-60). ACM.