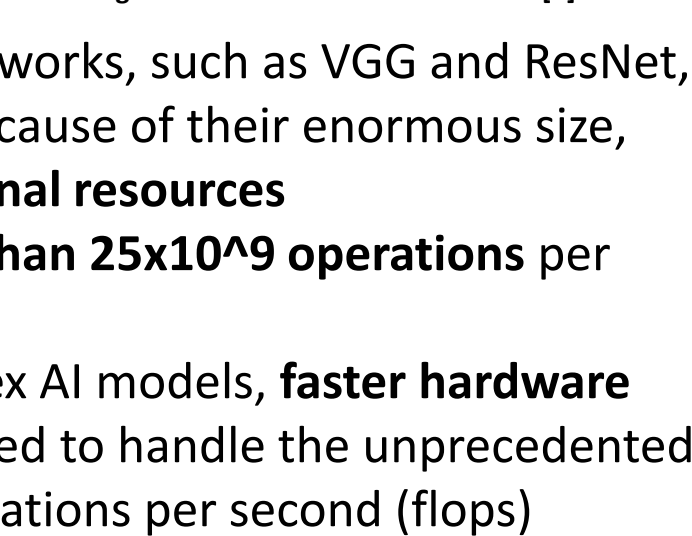
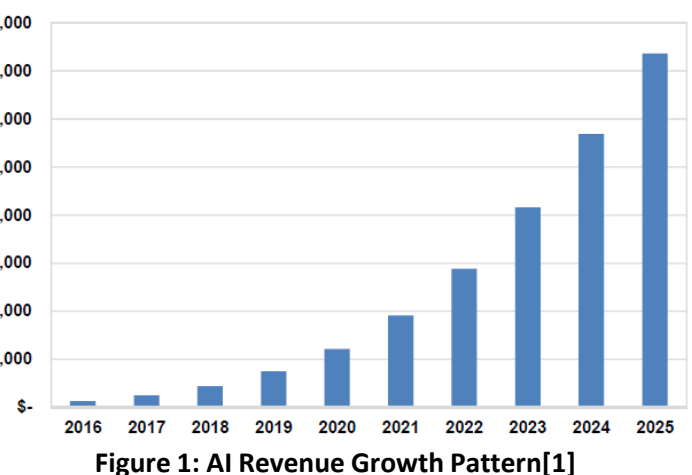


Optimizing Artificial Intelligence Performance with Morphable Parallel Computing Architectures

Introduction and Motivation

- Artificial Intelligence (AI), the field associated with modeling human intelligence for computation, has skyrocketed over the past few years
- AI has an unmatched **ability to make deep connections** between arbitrary inputs and outputs, often in ways researchers can't understand
- The most popular AI Structure is the **Neural Network (NN)**
 - As seen in Figure 2, Neural Networks, such as VGG and ResNet, are extremely accurate, but because of their enormous size, require **enormous computational resources**
 - Each of these require **greater than 25x10⁹ operations** per forward pass of data
 - To develop increasingly complex AI models, **faster hardware architectures** must be developed to handle the unprecedented demand for floating point operations per second (flops)

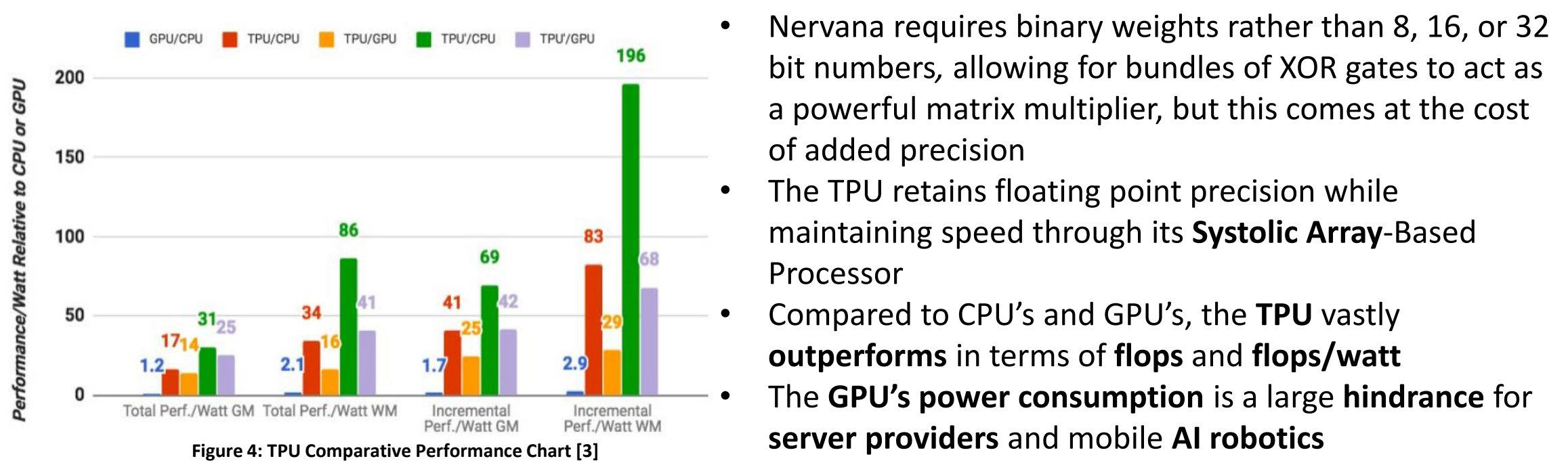


State of AI Hardware

The propagation of inputs through the weights of a Neural Networks is computed through **matrix multiplication**.

Here lies the computational overhead of Neural Networks, because operations (multiplications and additions) per second increases linearly with added nodes and added layers, making this a quadratic problem.

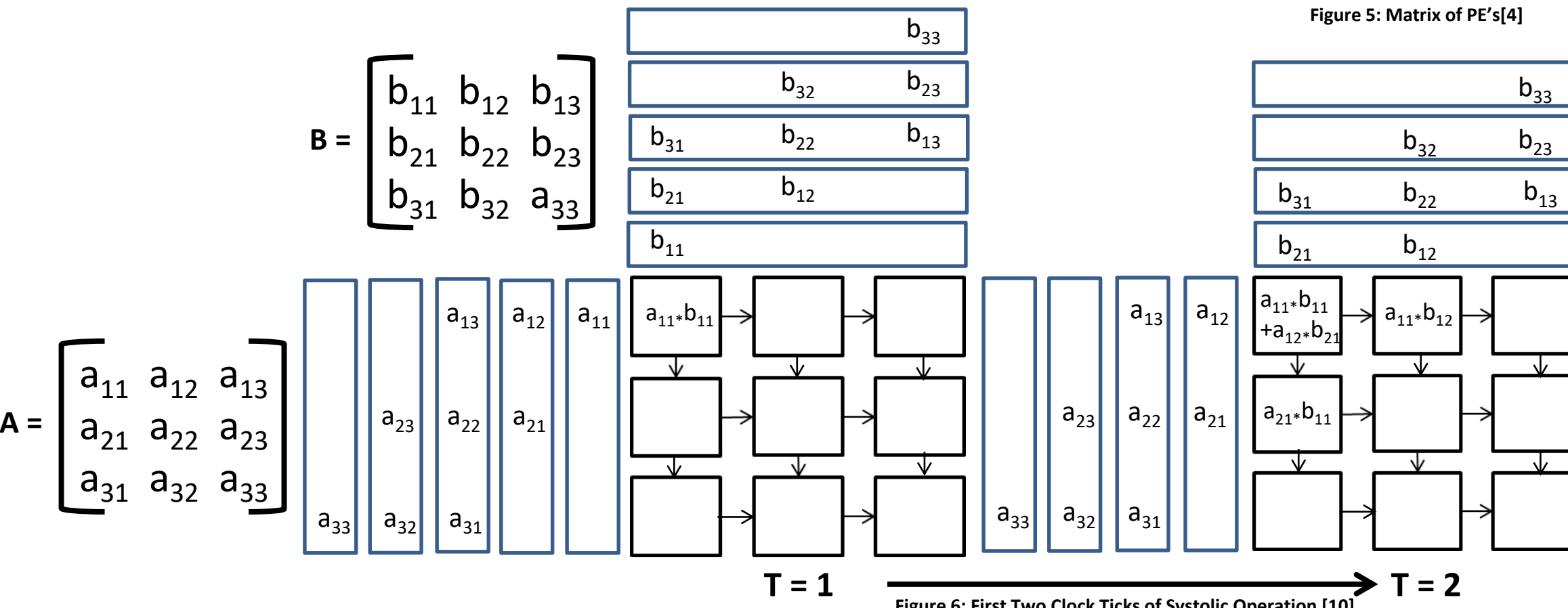
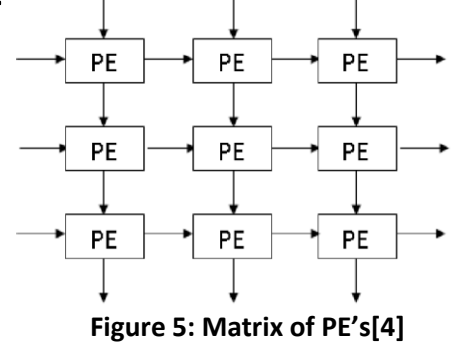
- Therefore, AI Hardware research is focused on **optimizing large-scale matrix multiplication**
- The most popular processors are graphics processing units (GPU)'s
- However, recently there have been several innovations in the realm with **Google's Tensor Processing Unit (TPU)** and **Intel's Nervana** Chip



- Nervana requires binary weights rather than 8, 16, or 32 bit numbers, allowing for bundles of XOR gates to act as a powerful matrix multiplier, but this comes at the cost of added precision
- The TPU retains floating point precision while maintaining speed through its **Systolic Array-Based** Processor
- Compared to CPU's and GPU's, the **TPU** vastly **outperforms** in terms of **flops and flops/watt**
- The **GPU's power consumption** is a large hindrance for **server providers** and **mobile AI robotics**

Systolic Array Background

- A Systolic Array is a **highly parallelized** computer architecture that uses a number of processing elements (PE) which multiply and accumulate (MAC) the dot products of the end matrix
- In the case of the TPU, the PE's were arranged in **rectangular formation**
- Consider the example, $C = AB$, where A and B are 3x3 matrices to be multiplied

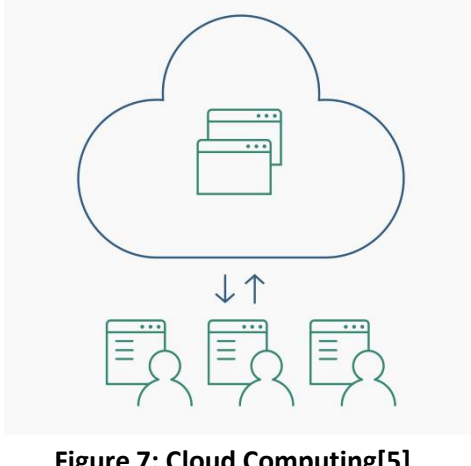


- This process continues until $T = 7$, or arbitrarily when all elements of the systolic array are filled and the input data has been emptied.
- Compared to a normal CPU, which would require $T=27$ to multiply the matrices, this system's concurrency drastically improves total latency

Scalability Issues and Focus

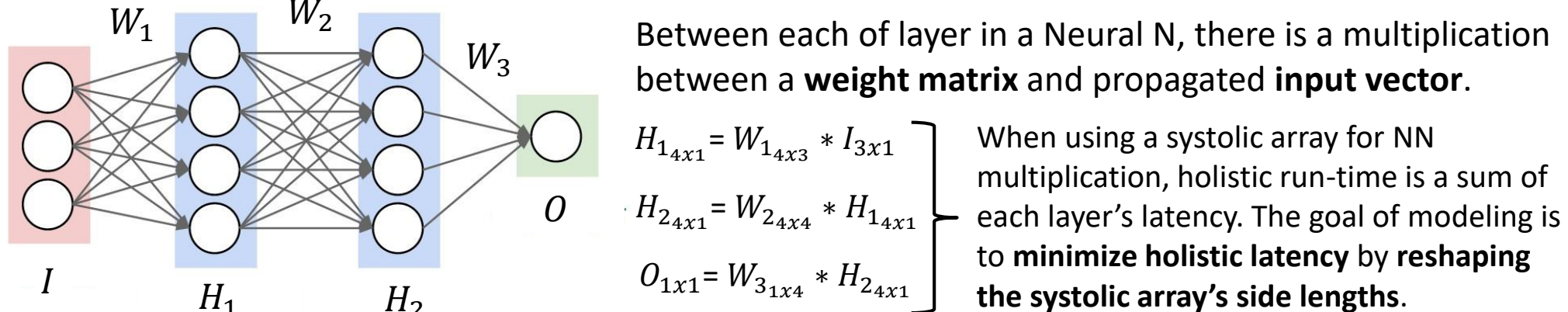
While the systolic array is state of the art, there are several perform hindrances at scale:

- Repetition** – For matrices which exceed the size of a systolic array's row and column size, the systolic array is forced to **repetitively cycle through the larger matrices** and perform block matrix multiplication. On Δ^2 problem
 - Memory** – To simply enlarge the systolic array for larger matrices is a faulted strategy because memory loadup latency **quadratically increases** with the side length of the matrix
 - Clock Cycles** – The **time required** for a systolic array to finish computation is **dependent on the size of the systolic array** itself, so the larger the systolic, the longer it takes to complete
- *Developing a sizing algorithm which minimizes these three functions is the central focus of this research**
- It is **imperative** that the **algorithm** not create additional latencies, but minimizes the 3 stated **hindrances**
 - The **Systolic Array** must be of 1-D Size to handle **matrix-vector** multiplications
 - Neural Network Training often contains **thousands of large-scale matrix** multiplications, since training sets often contain >1000 examples and >10 epochs
 - This is why Cloud Computing is a rising industry, but **power costs** for server providers and **hourly compute costs** for consumers is an issue
 - The development of a faster systolic array, which has been proven to be more power-efficient than GPUs, could help both server providers and consumers by lessening server-time usage.

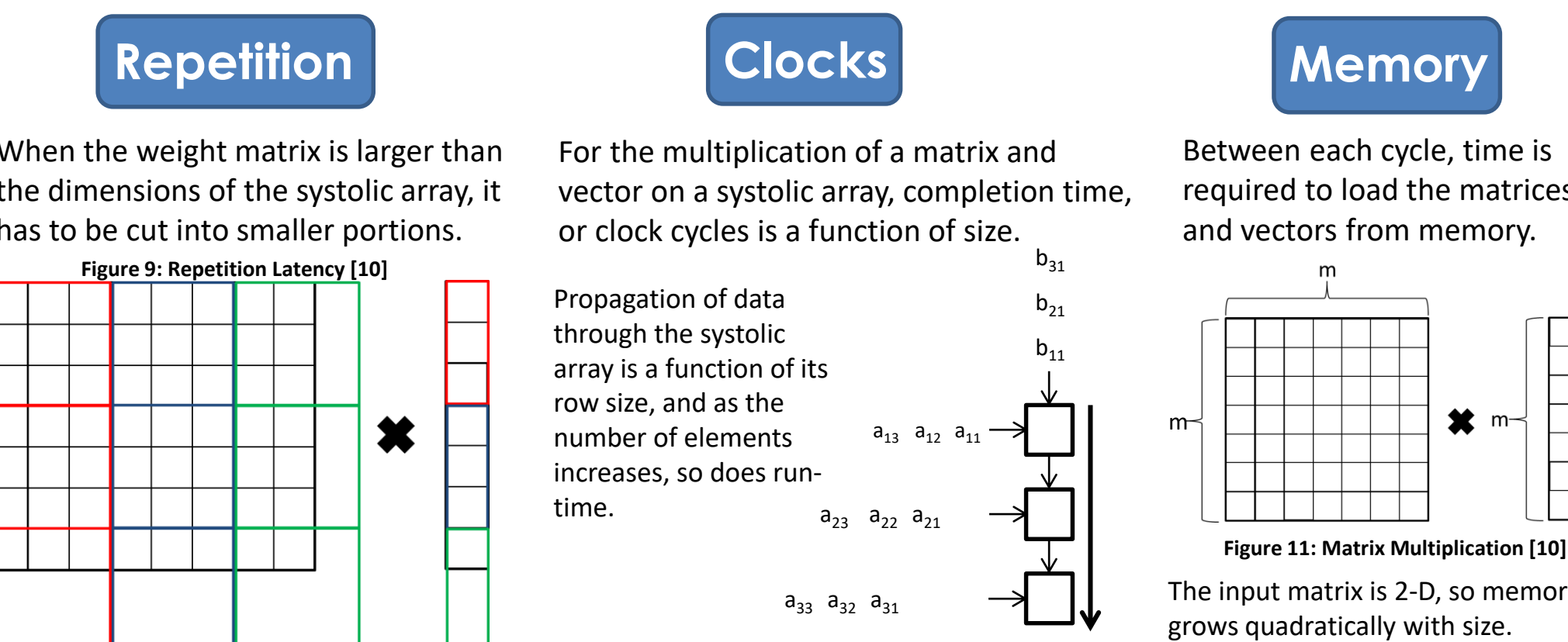


System Modeling

Latency Modeling



In the below equations, m = **systolic array row size**, A = **weight matrix**, B = **input vector**



When the weight matrix is larger than the dimensions of the systolic array, it has to be cut into smaller portions.

For the multiplication of a matrix and vector on a systolic array, completion time, or clock cycles is a function of size.

Between each cycle, time is required to load the matrices and vectors from memory.

Figure 10: Systolic Matrix-Vector Multiply [10]

$$\text{Clocks}(m) = 2 * m - 1$$

Figure 11: Matrix Multiplication [10]

The input matrix is 2-D, so memory grows quadratically with size.

$$\text{Memory}(m) = m^2 + m$$

Latency for Layer-wise Multiplication can be approximated by:

$$\text{LayerLatency}(m) = \text{Repetitions}(m) * (\text{Memory}(m) + \text{Clocks}(m))$$

The goal is now to find to **find the size (m)**, which will **minimize this function**. Rather than take a derivative and find x-axis intersections, a **1-D search algorithm** finds the smallest latency and its corresponding systolic size.

In the side case, the optimal size for multiplying a 512x512 matrix and a 512x1 vector is 170 Processing Elements.

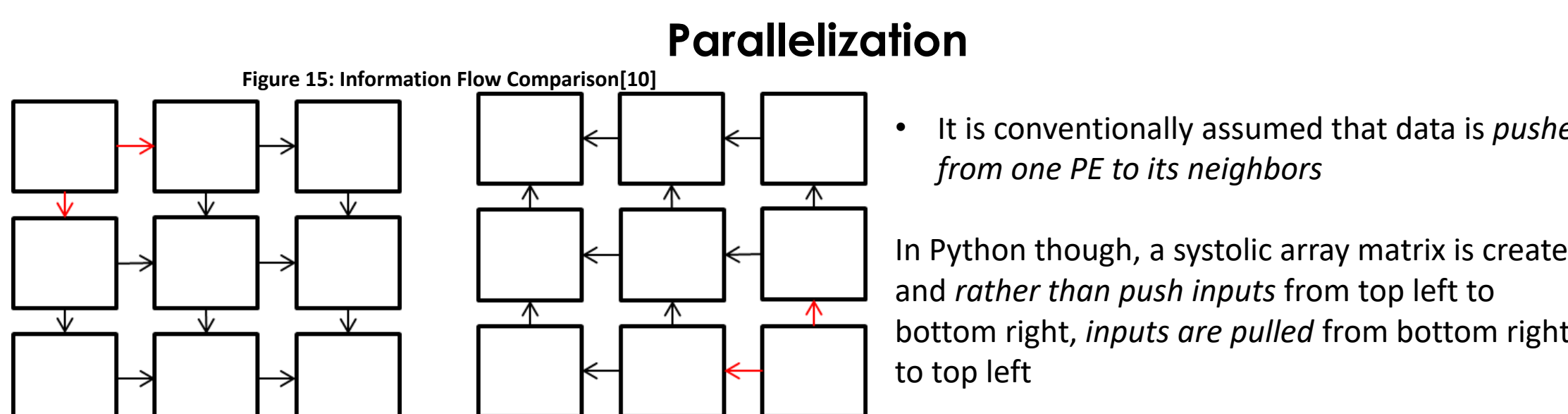
Figure 12: Theoretical Latency[10]

Simulator Design and Testing

Design Challenges and Goals

- In order to test the efficacy of the proposed algorithm and Circular-Flow Systolic Structure, a high-level simulator was sought to be designed.
- Previous works have attempted such a task [5], but not with the programmability incurred with a higher-level programming language like Python, which is the suggest language.

Challenges	Requirements
Parallelization	Systolic Array's are effective because of their parallel nature, and while all higher-level languages are executed asynchronously and line-by-line, the proposed system will have to mask any iterative nature incurred by trying to model the systolic array.
Memory Access	In ASIC and FPGA applications, it takes time to lead floating-point numbers from memory. In fact, it was stated that compared to a 512x512 size, the 256x256 systolic array performed much faster, simply because less memory had to be retrieved per clock cycle [8].
Processing Element Complexity	Processing Elements have two functions: Multiply and Accumulate or Pass on Inputs to Neighbors. Therefore, to model simplicity, the Processing Element Objects must be written concisely.



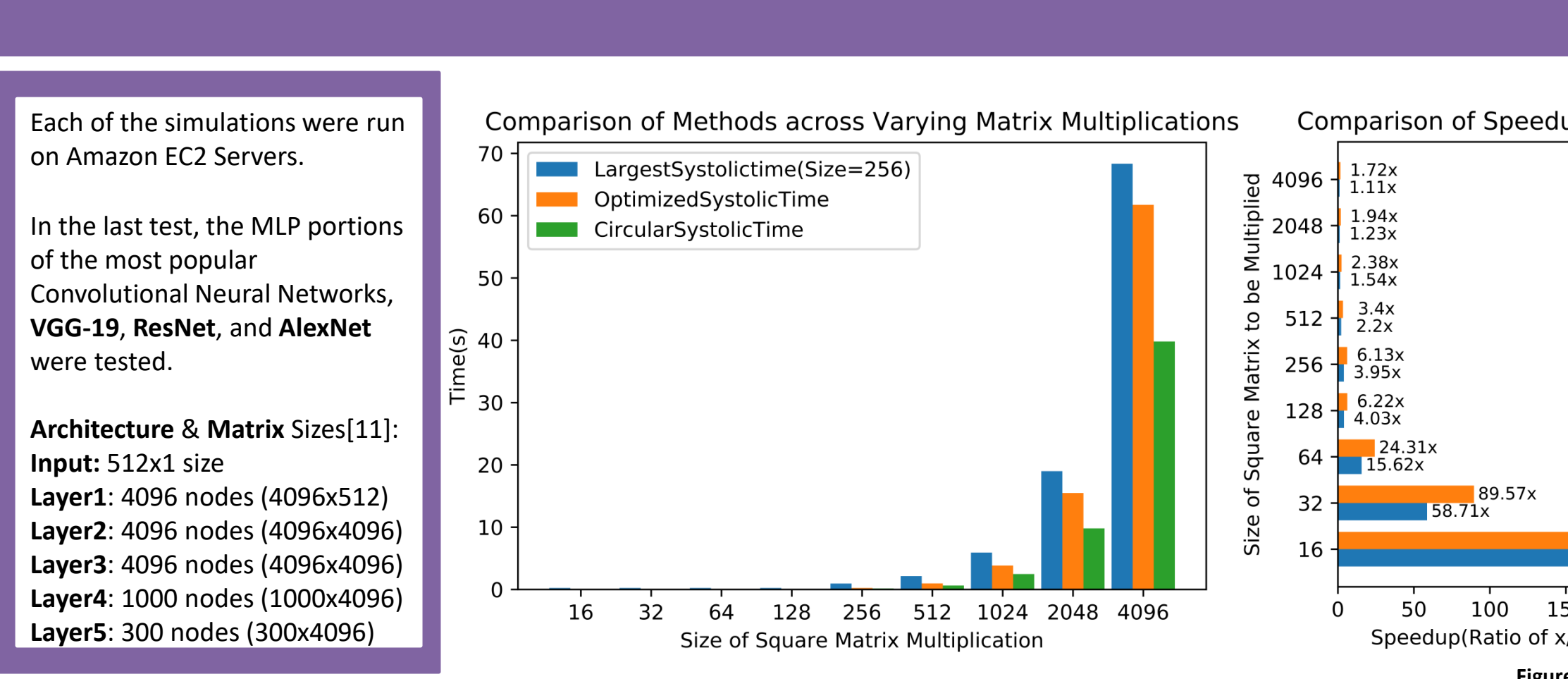
- It is conventionally assumed that data is **pushed from one PE to its neighbors**
- In Python though, a systolic array matrix is created and **rather than push inputs** from top left to bottom right, **inputs are pulled** from bottom right to top left

In Figure 16, the information within the systolic matrix can be modeled by a set of discrete time steps. A while loop cycles through these steps until multiplication is finished.

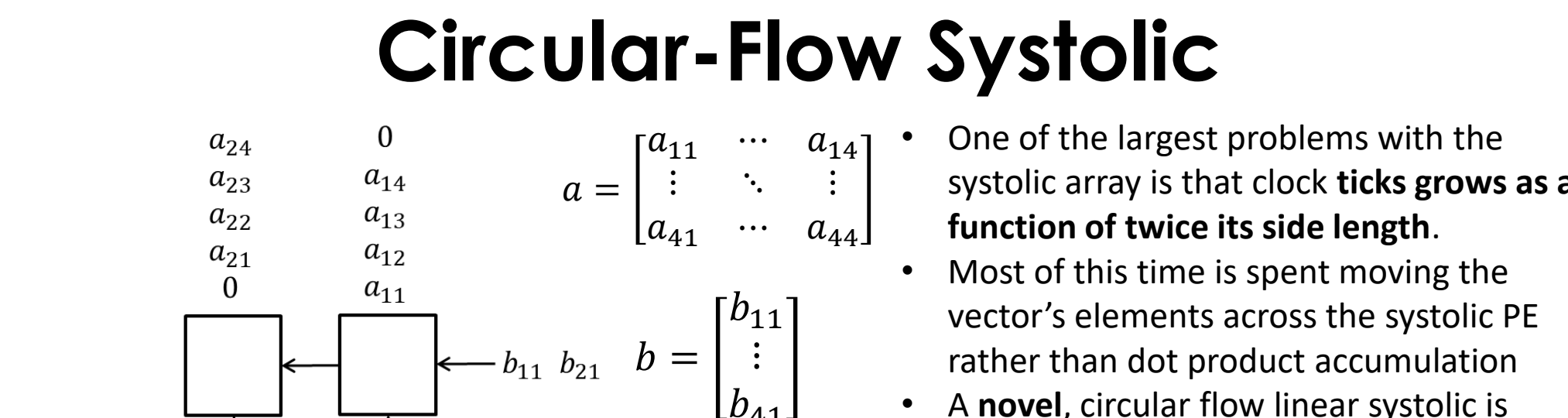
To reduce PE complexity, each PE evaluates 4 rules:

- If(NorthValue and WestValue aren't None): **perform a MAC operation**
- If(NorthValue is None, but NorthNeighbor's NorthValue isn't None): **NorthValue = NorthNeighbor's NorthValue**
- If(WestValue is None, but WestNeighbor's WestValue isn't None): **WestValue = WestNeighbor's WestValue**
- If(!MAC's for all PE's = col_size of A): **stop systolic array**

Figure 16: Discrete Time Increments of Systolic[10]



Circular-Flow Systolic



The minimum number of clock cycles required to multiply the two matrices is:

$$\text{Clocks}(m) = 1.5m - 1$$

Theoretically, in the case of a 256x256 and 256x1 matrix multiply, the repetition speedup with a circular-flow systolic is around 1.333x

Similar to the Morphable Systolic System, the optimum shape is first calculated based on the sizes of matrices in the instantaneous layer.

However, holistically, in order to find the optimum size for the entire network, the network optimum size is a weighted average of individual layer latency and size.

Figure 14: Clock Cycle Comparison[10]

Figure 18 shows the CircularSystolic's true speed enhancement compared to a Normal Systolic.

Figure 18: Statistical Comparison of Systolic's[10]

Figure 19: Systolic Performance Analysis [10]

Figure 20: Approach of Optimal to Max [10]

Figure 21: Relative Speedup of Circular [10]

Figure 22: Mojo FPGA [7]

Figure 23: PE at FPGA Scale [8]

Figure 24: Convolutional Kernel Operation[9]

Figure 25: Convolutional Kernel Operation[9]

Figure 26: Convolutional Kernel Operation[9]

Figure 27: Convolutional Kernel Operation[9]

Figure 28: Convolutional Kernel Operation[9]

Figure 29: Convolutional Kernel Operation[9]

Figure 30: Convolutional Kernel Operation[9]

Figure 31: Convolutional Kernel Operation[9]

Figure 32: Convolutional Kernel Operation[9]

Figure 33: Convolutional Kernel Operation[9]

Figure 34: Convolutional Kernel Operation[9]

Figure 35: Convolutional Kernel Operation[9]

Figure 36: Convolutional Kernel Operation[9]

Figure 37: Convolutional Kernel Operation[9]

Data Analysis

- As seen in Figure 20, simply **optimizing the size of a Systolic Array** with the proposed algorithm leads to **speedups of 1.54x** when multiplying matrices and vector's which **quadruple the size of the array**
- However, when approaching the largest of matrices, the algorithm's benefits begins to decrease, since at a **4096x4096 matrix multiply**, the speedup is only around **1.11x**
- This is because when specifying a max # of PE's, which is necessary since **hardware resources are finite**, memory approaches a threshold, while repetitions and clock cycles continue to increase
- As seen in Fig. 20, the **optimum size quickly approaches the max # of PE's**
- However, the Circular Systolic circumvents this issue with a **constantly smaller number of clock-ticks compared to a Normal Systolic Array**
- As seen in Figure 21, for a 512x512 multiplication, regardless of spatial dimension, the **speedup** for a Circular Systolic is **constant** compared to a Normal Systolic
- The **true test** came with Figure 19, and the comparison of the various Systolic Structures versus the Multi-Layer-Perception (MLP) Layers of the largest, most popular Convolutional Neural Networks
- Each CNN model uses the **same MLP architecture**, and across each layer, **both the Optimized Systolic and Circular Systolic showed speedups over the traditional, maximum size Systolic**
- This was a unique test of the Systolic's Performance since unlike the tests, the matrices weren't always square, but rectangular.
- The algorithm accounted for this and produced optimal sizes dissimilar to that of square matrices.

Summary and Conclusion

- Artificial Intelligence** has been **growing at an exponential rate** over the past few years, which has warranted the research of faster, more computationally-capable hardware
- AI is reliant upon **fast matrix multiplication**, which is where the **GPU** has excelled, but research at **Google** has shown that the **Systolic Array is a promising option for its flops and flops/watt performance**
- The Systolic Array is highly parallelized, but when bound in the spatial dimension, can significantly suffer with large-scale matrix multiplication
- The intent of the research was to **develop a cost function for any systolic array**, based on a given matrix-vector operation, and **output the optimum systolic array shape**
- A mathematical model was created and an algorithm in Python performs a 1-D search to find the optimal size
- Additionally, a unique Systolic Structure called the **Circular Flow Systolic** is proposed to **minimize** the Systolic's **clock cycles** while keeping other factors constant
- Due to hardware limitations, a **simulator was designed** in Python to evaluate the proposed algorithm and Circular Systolic Structure
- It was shown that **across various square matrix multiplications** and rectangular multiplications as found in popular CNN Models, both the **algorithm** and the **Circular Systolic provided significant speed enhancements**
- Applications of this technology** might find itself in **mobile robotics**, where matrix multiplication latency is required to be minimized or in **Cloud Computing**
- Companies such as Amazon and Google provide hardware for developers to train their AI models, and with faster, **more electricity efficient systems**, the **providers would spend less money** hosting the servers, and **consumers would pay less** for hourly consumption as training time would be lessened

Preliminary FPGA Development

- While the focus of this project was to develop an algorithm for enhancing matrix-vector multiplication, a true test of its efficacy is to test on a **Field Programmable Gate Array (FPGA)**
- FPGA's are programmable chips that are commonly used for prototyping, but in this scenario, its **morphability** makes it an ideal candidate for the proposed algorithm
- FPGA's are also significantly **more power efficient than GPU's**, but lower-end FPGA's drastically suffer with memory access
- Unlike in Python, where commands are executed sequentially, FPGA's execute lines instantaneously
- To the Left, is an example of a PE on an FPGA, simply consisting of an adder, multiplier, RAM, and registers for the outputs
- A simplistic example with **multiplying 2, 3x3 matrices** was conducted, and the execution time was **46ns** at a clock speed of **128 MHz**
- A CPU executing **C++ multiplication**, chosen for the languages speed, takes **0.054s**.

Figure 22: Mojo FPGA [7]

Figure 23: PE at FPGA Scale [8]

Figure 24: Convolutional Kernel Operation[9]

Figure 25: Convolutional Kernel Operation[9]

Figure 26: Convolutional Kernel Operation[9]

Figure 27: Convolutional Kernel Operation[9]

Figure 28: Convolutional Kernel Operation[9]

Figure 29: Convolutional Kernel Operation[9]

Figure 30: Convolutional Kernel Operation[9]

Figure 31: Convolutional Kernel Operation[9]

Figure 32: Convolutional Kernel Operation[9]

Figure 33: Convolutional Kernel Operation[9]

Figure 34: Convolutional Kernel Operation[9]

Figure 35: Convolutional Kernel Operation[9]

Figure 36: Convolutional Kernel Operation[9]

Figure 37: Convolutional Kernel Operation[9]

Figure 38: Convolutional Kernel Operation[9]

Figure 39: Convolutional Kernel Operation[9]

Figure 40: Convolutional Kernel Operation[9]

Figure 41: Convolutional Kernel Operation[9]

Figure 42: Convolutional Kernel Operation[9]

Figure 43: Convolutional Kernel Operation[9]

Figure 44: Convolutional Kernel Operation[9]

Figure 45: Convolutional Kernel Operation[9]

Figure 46: Convolutional Kernel Operation[9]

Figure 47: Convolutional Kernel Operation[9]

Figure 48: Convolutional Kernel Operation[9]

Figure 49: Convolutional Kernel Operation[9]

Figure 50: Convolutional Kernel Operation[9]

Figure 51: Convolutional Kernel Operation[9]

Figure 52: Convolutional Kernel Operation[9]

Figure 53: Convolutional Kernel Operation[9]

Figure 54: Convolutional Kernel Operation[9]

Figure 55: Convolutional Kernel Operation[9]

Figure 56: Convolutional Kernel Operation[9]

Figure 57: Convolutional Kernel Operation[9]

Figure 58: Convolutional Kernel Operation[9]

Figure 59: Convolutional Kernel Operation[9]

Figure 60: Convolutional Kernel Operation[9]

Figure 61: Convolutional Kernel Operation[9]

Figure 62: Convolutional Kernel Operation[9]

Figure 63: Convolutional Kernel Operation[9]

Figure 64: Convolutional Kernel Operation[9]

Figure 65: Convolutional Kernel Operation[9]

Figure 66: Convolutional Kernel Operation[9]

Figure 67: Convolutional Kernel Operation[9]

Figure 68: Convolutional Kernel Operation[9]

Figure 69: Convolutional Kernel Operation[9]

Figure 70: Convolutional Kernel Operation[9]

Figure 71: Convolutional Kernel Operation[9]

Figure 72: Convolutional Kernel Operation[9]

Figure 73: Convolutional Kernel Operation[9]

Figure 74: Convolutional Kernel Operation[9]

Figure 75: Convolutional Kernel Operation[9]

Figure 76: Convolutional Kernel Operation[9]

Figure 77: Convolutional Kernel Operation[9]

Figure 78: Convolutional Kernel Operation[9]

Figure 79: Convolutional Kernel Operation[9]

Figure 80: Convolutional Kernel Operation[9]

Figure 81: Convolutional Kernel Operation[9]

Figure 82: Convolutional Kernel Operation[9]

Figure 83: Convolutional Kernel Operation[9]

Figure 84: Convolutional Kernel Operation[9]

Figure 85: Convolutional Kernel Operation[9]

Figure 86: Convolutional Kernel Operation[9]

Figure 87: Convolutional Kernel Operation[9]

Figure 88: Convolutional Kernel Operation[9]

Figure 89: Convolutional Kernel Operation[9]

Figure 90: Convolutional Kernel Operation[9]

Figure 91: Convolutional Kernel Operation[9]

Figure 92: Convolutional Kernel Operation[9]

Figure 93: Convolutional Kernel Operation[9]

Figure 94: Convolutional Kernel Operation[9]

Figure 95: Convolutional Kernel Operation[9]

Figure 96: Convolutional Kernel Operation[9]

Figure 97: Convolutional Kernel Operation[9]

Figure 98: Convolutional Kernel Operation[9]

Figure 99: Convolutional Kernel Operation[9]

Figure 100: Convolutional Kernel Operation[9]