"What Intensivate is doing is taking accelerator technology in a totally new direction. With their architecture, we'll be able to provide a new Cloud offering."

Michael McGrath
Technologist
Microsoft Azure

## INTRODUCTION

In the past, servers were primarily used for business functions such as billing and CRM. These servers and hardware platforms were developed to provide maximum performance when running this type of software. But within the last few years, the exponential expansion of human and machine generated data has drastically changed the way business is done.  With the change in business has come a change in the software used to collect, process, analyze and leverage information.

## NEW GENERATION DATA PROCESSING SOFTWARE

A slew of new software platforms has evolved over the last decade to support the varied ways that very large amounts of data are put to use inside a business.  Data is used in not just a single way.  Rather, a plethora of approaches has arisen, which includes traditional analytics (but on much larger data sets), real time analytics, advanced analytics (which often involves some form of Machine Learning), and the increasing use  of webs of processing performed on streams of data.

These new software platforms have one thing in common: they scale to hundreds or thousands of servers.  The platform provides a convenient way for software developers to write relatively straight forward code.  The platform then handles making that code scale out to many servers.  This has enabled an explosion in the number, variety, and nature of applications that extract value from the vast amount of data available, and manage that data within the Enterprise environment.

Names of popular platforms include: Spark, Hadoop, Kafka, and various NoSQL databases.

Applications written using these new software platforms are called **scale-out applications**, because they automatically scale to hundreds or thousands of servers.

## HOW CPU IMPLEMENTATION RELATES TO SCALE OUT APPLICATIONS

CPUs are designed to run traditional business applications.  Such software was written in a style called "single threaded" or sometimes "multi-threaded".  It may run on a few servers, but it does not easily scale to a large number of servers. Applications like this force a CPU to provide high single threaded performance, which in turn requires a massive amount of circuitry inside the CPU.  The extra circuitry makes CPUs expensive and consumes a large amount of power.

Scale out applications have reached a critical mass, now consuming more than 20% of all servers in existence. Meanwhile, CPUs are forced to continue supporting traditional business applications, and so must continue to include the massive circuitry that makes them expensive and power hungry.

---

1  IDC "The Evolution of Data to Life Critical," 2017. https://www.seagate.com/www-content/our-story/trends/files/Seagate-WP-DataAge2025-March-2017.pdf

"The trend I'm seeing…is to move into stream processing and real-time analytics. This kind of software doesn't run on current accelerators. The Intensivate accelerator card is the only technology I've come across that addresses this growing problem."

Piet Loubser
SVP Paxata
Recent VP Product Marketing
Hortonworks

## ACCELERATOR TECHNOLOGIES

Today's IT system architects live in a world of "accelerator" implementations, where hardware is optimized to a workload having specific patterns. Each accelerator has made choices down inside the silicon about the pattern of wires placed on the chip. This pattern is fixed when the chip is manufactured.  CPUs choose a middle of the road pattern that performs acceptably on everything. A sort of jack of all trades — exceptional on nothing but also terrible on nothing.  In contrast, an accelerator makes choices that make is astoundingly good at a few things, but as a consquence terrible on others.  This is how an accelerator achieves an advantage over a CPU.  The narrower the set of applications, the higher the advantage over a CPU it can achieve on those specific applications that match the accelerator's internal patterns.

There are four main accelerator categories, based on the underlying pattern of the chip, namely:

1. GPU – Effective for Deep Learning training, and matrix multiplication based algorithms
2. AI accelerators — Highly specialized to Deep Learning
3. FPGA – Serves specific functions with certain characteristics, which are rarely found
4. Scale Out Accelerator – Introduced by Intensivate, for applications that scale to hundreds or thousands of nodes — within a node, broadly applies to any pattern

These underlying technologies are deployed as PCIe accelerator add-in cards, in standardized server platforms (Intel Xeon or AMD based) or embedded on system mother boards

## RISC–V

Today, the standardized server platforms are overwhelmingly x86 based, with one company monopolizing the market, which imposes massive premiums on the CPUs that power these servers. The x86 instruction set is proprietary and closed. Opposition has formed, and is making headway in the guise of ARM based servers, but adoption of ARM in data centers is still very limited.

Recently, a new open source Instruction Set Architecture named RISC–V has appeared. Fueled by these market forces, RISC–V currently has over 100 industry leading companies within its consortium (RISCV.org). Due to its open source nature and the availability of a complete, open source, high quality chip implementation called Rocket Chip (RISC–V based SOC generator), there has been a rapid build up of a wide ecosystem of support.  It promises large-scale industry-wide adoption and success, similar to the Linux OS domination seen in the server software platform space.

"I'm excited about Intensivate's accelerator card and see strong opportunity for this technology to be adopted by leading server vendors like Dell, Lenovo, and HPE."

Jai Menon
Former CTO, IBM Systems Business
Former CTO, Dell Enterprise Solutions

## INTENSIVATE SOLUTION

Intensivate has adopted RISC-V as the interface to its accelerator technology, thereby capitalizing on the support afforded by the community, and so lowering the cost to bring this new kind of accelerator to market. RISC-V further enables a seamless transition, to running existing scale out applications on RISC-V powered accelerator cards, with low effort and low friction.

Intensivate's accelerator chip has these key features:

a)  High performance — a clock speed of 3.0GHz with 1.8 instructions completed per cycle
b)  Very low power – 550mW/core
c)  Very small foot print (0.9mm2)
d)  Multiple cores per accelerator chip (16 cores)
e)  Multiple threads per core (16 threads)
f)  Superior memory behavior, sustaining 8 outstanding memory misses while maintaining the 1.8 instructions completed per cycle, making it ideal for typical scale out applications

Any instruction set could have been chosen as the interface to the accelerator core, but RISC-V comes with the full toolchain, OS, and applications in place, resulting in dramatic reduction in the cost to deploy and making adoption by customers seamless. Further, the open source Rocket Chip generator was fully utilized to reduce engineering cost and risk of producing a chip. The portions of Rocket Chip that don't affect power or performance are maintained, thereby saving many person-years of effort, while the performance critical core was ripped up and replaced with Intensivate's patent pending technology. The envelope around the core was maintained, allowing IntenCore to drop directly into Rocket Chip without the rest of the chip being able to tell the difference. This strategy serves the purpose of developing an accelerator processor at a fraction of the cost vs. developing from the ground up, while using mature, proven, components of Rocket Chip.

**intensivate**

## System Configuration Overview

### PCIe Accelerator Highlights

---

### Network Bandwidth: Gbps/Core

1 dual 100Gb card (200Gbps) dedicated to each Intensivate accelerator card

Each Intensivate card has 336 cores

Equates to 0.52Gbps/Core

---

### Memory Bandwidth: GB/s/Core

Each accelerator chip has 4 channels of LPDDR4 DRAM. Total 52 GB/s memory BW

Each Intensivate chip has 16 cores

Equates to 3.25 GB/s/Core

---

### For Comparison

Network Comparison

Intensivate : **.52Gbps/Core**

Typical AWS EC2
Machine Learning  Config: **.40Gbps/Core**

Memory Comparison

Intensivate: **3.25GB/s/Core**

Typical AWS EC2
Machine Learning Config: **2.8 GB/s/Core**

---

Accelerator card can be configured with 1U, 2U or 4U standard rackmount server configurations

- Each PCIe card has 21 Intensivate accelerator chips, each with its own 32GB of LP DDR4 providing a total of 672GB of DRAM on the card

- Scale out applications see each accelerator chip as an independent node with its own IP address

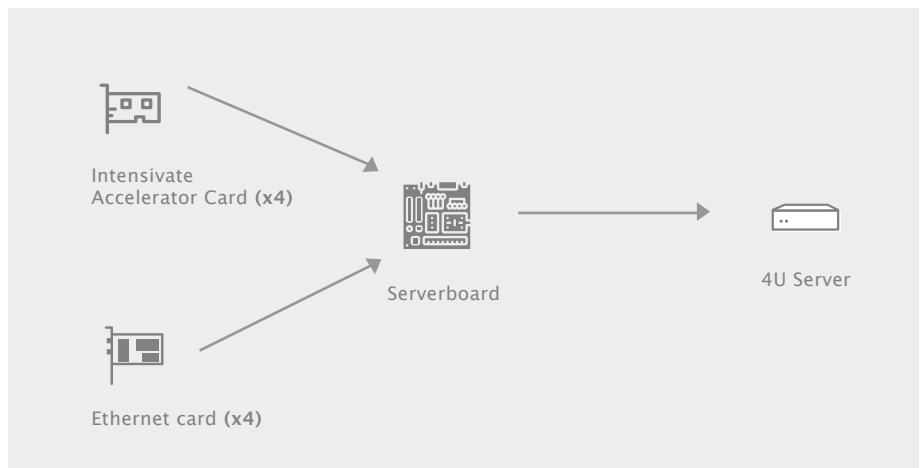- High speed 120 Gbps chip to chip mesh network supports network intensive scale out applications

Intensivate
Accelerator Card **(x4)**

Serverboard

4U Server

Ethernet card **(x4)**

Figure 1. example of typical system configuration

### Component Details (see figure 1)

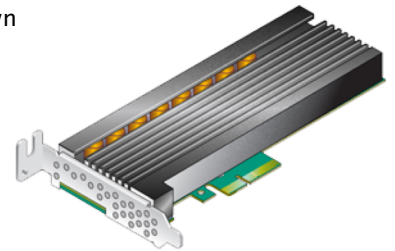| Part | Source | SKU | Quantity |
|------|--------|-----|----------|
| 4U Server | Supermicro | 4029GP–TRT2 | 1 |
| Serverboard | Supermicro | X11DPG–OT–CPU | 1 |
| 100 Gb/s x2 Ethernet Adapter Card | Mellanox | MCX415A–CCAT | 4 |
| Intensivate Accelerator card | Intensivate | IntenScale | 4 |

## Accelerator Card Overview

### High Level Card Specs

PCIe card holds 21 Intensivate ASIC Chips

336 cores per card

Each chip paired with 32GB DRAM plus

50GB/sec chip to chip mesh network

672 GB memory per card

Card seen as sub-net of 21 independent

servers, each indistinguishable from a

rack mount server

## IntenScale PCIe SCALE OUT ACCELERATOR CARD

The IntenScale PCIe accelerator card holds 21 IntenScale accelerator chips. Each chip is paired with 32GB of its own memory. Intensivate's patent pending technology makes each such chip on the card appear to software as a complete, independent rackmount server with its own IP address and own disk. Hence each accelerator card contains the equivalent of 21 1U rackmount servers on it, at a fraction of the cost and fraction of the power. All 21 server equivalents on the card consume a total 350W combined, compared to 4,400W for the 21 rackmount servers. A high-speed mesh network connects all 21 accelerator chips on the card, providing a very high compute intensive platform for scale out applications. One may configure multiple IntenScale cards in a 1U, 2U or 4U standard rackmount system.

## Block Diagram