Search...

Q

DIY Big Data (https://diybigdata.net/)

One byte at a time ...

HOME

PROJECTS

ABOUT

LEGAL

Pnoject Kickoff – Personal Compute Cluster 2019

Edition

September 1, 2019 (https://diybigdata.net/2019/09/personal-compute-cluster-2019-project-kickoff-2/)

michael (https://diybigdata.net/author/michael/)

Three years ago, I worked through a project of creating a low cost, low power computer cluster into s://diybigdata.net/odroid-xu4-cluster/), with the primary goal of becoming more familiar with the inner works of Apache Spark. This project did accomplish that goal, but since this cluster was made up of 32-bit ARM processors and each had only 2 GB of RAM, the cluster was not too useful for getting meaningful work done. What it did excel at, however, was showing the user how to write Spark code efficiently. If you wanted to get anything done on this constrained system you had to be mindful of every inefficiency.

Jump ahead to 2019, and I have decided to give it another go. This time, I want to make a cluster that is moderately useful for data analysis and machine learning, but does not break the bank either. The first step is to document my system design requirements and general goals:

- The cluster's primary purpose will be to run Apache Spark to do data analysis. However, I would also like to experiment with ElasticSearch.
- The cluster should be usable for data sets up to 256 GB or more in size, and have have storage for a few terabytes of data.
- The cluster should use 64 bit x86 CPUs. My last cluster was based on 32 bit ARM CPUs. This created <u>some compatibility issues (http://ARM)</u> that were interesting to sort through, but I don't want to deal with those sort of problems this time.

- Building, configuring, and operating this cluster should facilitate the following learning objectives:
 - Understand the hardware tradeoffs for performance in distributed computing systems.
 - Advance my understanding and skill in Apache Spark
 - Learn how to run Apache Spark in Kubernetes.
 - Experiment with HDFS 3.x and compare it to QFS 2.x.
- My target budget is \$3,000, and I would prefer to come in under that but would go over if there is value in it.
- I want at least four physical nodes.

Hardware Selection

Probably the biggest driver in my hardware selection is my budget. I would love to get some of the new server-class motherboards. For example, SuperMicro's new motherboard basedon AMD's EPYC 3251 CPU (https://www.supermicro.com/Aplus/motherboard/Embedded/M11SDV-4C-LN4F.cfm) particularly caught my eye, but building out a single cluster node that reasonably leverages the features of the SuperMicro motherboard resulted in a per node code of well over \$2,000. That would be most of my budget for a single node. Given my goal of have a distributed hardware environment, I would have to scale back my expectation on per-node capabilities.

I die a lot of searching for this. I considered doing my own node builds, looked at low cost computers rom Walmart, and even looked at some of the more powerful single board computers, such as the UDOO Bolt (https://www.udoo.org/udoo-bolt/). Given my general project goals, the technical parameters I looked at were:

- **CPU Cores and Threads** The more threads the better for each node. This would allow a higher level of parallelism within the cluster. It would also allow better handling of running the distributed file system across the same nodes that spark is running on.
- **CPU Speed** Ultimately I want this cluster to be useful and more than just a learning vehicle like my last cluster was. So I considered the CPU's benchmark when differentiating between choices.
- Maximum RAM Capacity Spark works best with the more RAM it has available to it, so you want to have a high RAM to thread count ratio.
- NVMe SSD Even though Spark leverages RAM as much as possible to make calculations fast, it still needs to spill data to disk in order to manage operations. Ideally, the isa it spills data to is fast. Furthermore, the distributed file system's performance is related to the performance of its storage media. NVMe SSDs allow high performance at a relatively cheap price (1 TB NVMe SSDs can be found for about \$100).
- Small Form Factor I don't have a data center to put this cluster into, just my desk in my home
 office. I want to keep it's footprint low.

• **Future Expandability** – The two facets of expandability that are important here are more storage and faster networking. Ideally the node type I select would allow me to upgrade storage and networking if I wish to invest more money into the cluster.

At one point I was seriously considering the UDOO Bolt. It seemed to check all the qualities I was looking for in a node computer. Its eight thread AMD V1605B

(https://www.cpubenchmark.net/cpu.php?cpu=AMD+Ryzen+Embedded+V1605B&id=3331) CPU seemed to provide strong performance for relatively low cost and thermal design power. I even placed a preorder as UDOO was accepting preorders for delivery in July 2019. But July 2019 came and went and there was no solid date for delivery, so I canceled my preorder and started looking again. And I am glad I did.

I ended up coming across the EGLOBAL s200 Mini PC (http://s.click.aliexpress.com/e/3rIO2ICs). This computer is manufactured and distributed by a Chinese company, so you would have to order it from a site like AliExpress. What is interesting about this computer is that you can choose from four CPU opulons, including a Core i9 8950HK (https://www.cpubenchmark.net/cpu.php?cpu=Intel+Core+i9-8950HK+%40+2.90GHz&id=3246) or a Xeon E-2176M (https://www.cpubenchmark.net/cpu.php?cpu=Intel+Xeon+E-2176M+%40+2.70GHz&id=3242), both of which provide 12 threads and have henchmarks more than twice as fast that of the UDOO Bolt's AMD V1605B. Furthermore, the computer can accept up to 64 GB of DDR4 2666 MHz RAM, it has two M.2 NVMe SSD slots, ability to add a 2.5" SATA3 SSD, and comes in an extremely small form factor. Best of all, this S200 computer can be purchased at about the same price as the UDOO Bolt. So I bought four. One thing I will note on sourcing these computers is that they seem to be widely side by different resellers on AliExpress and even on Amazon (https://amzn.to/2MXhJrc) (affiliate link) though its rebranded there, and can be found with a wide range of prices. At the time of this writing, I found the Topton Computer Store (http://s.click.aliexpress.com/e/4NhNHqzi) (affiliate link) as having the cheapest prices.

I bought the computers without any RAM or SSD, allowing me to select precisely which ones I felt would be best for my cluster build. Though this computer can accept up to 64 GB of RAM, the manufacture doesn't offer a prebuilt version with that much RAM installed. The computer has two SO-DIMM sockets, so I purchased two 32 GB RAM modules for each computer. For storage I decided to start with a 1 TB SSB. Given the computer's expansion abilities, I could add more later if I desired. The SSD storage that the computer manufactured offers was slow, and relatively expensive. The EGLOBAL S200 computer is using PCIe Gen 3, so this meant I couldn't use the fastest of PCIe Gen 4 SSDs, but there are plenty fast SSDs for Gen 3. I decided to go with the <u>Sabrent Rocket line of M.2 NVMe SSD (https://www.sabrent.com/rocket-nvme/)</u>, chiefly due to it's very high speed and relatively low co. What I ended up purchasing for each node is as follows (note that all product links are affiliate links):

Item	Source	Price per Item	Item Count per Node
EGLOBAL S200 Computer with Xeon E-2176M CPU	AliExpress (http://s.click.aliexpress.com/e/3rIO2ICs)	\$437.73	1
32 GB RAM Module Samsung DDR4 2666MHz 260 Pin SODIMM, 1.2V	<u>Amazon (https://amzn.to/2NCNIy8)</u>	\$137.75	2
1TB SSD Sabrent Rocket NVMe PCIe M.2 2280 SSD	<u>Amazon (https://amzn.to/2L2zrsF)</u>	\$109.98	1

This puts my total per node cost at \$823.21. For four nodes, my total cluster cost already exceeds my \$30,00 price target, but given the number of threads this cluster will have and the 256 GB RAM pool, I'm very happy with the this cluster's specs. You could cut costs by going with a cheaper CPU option for the \$200, or by using less RAM or a smaller SSD. I would note that the price for the computer seems to fluctuate a lot on AliExpress. Your price might vary.

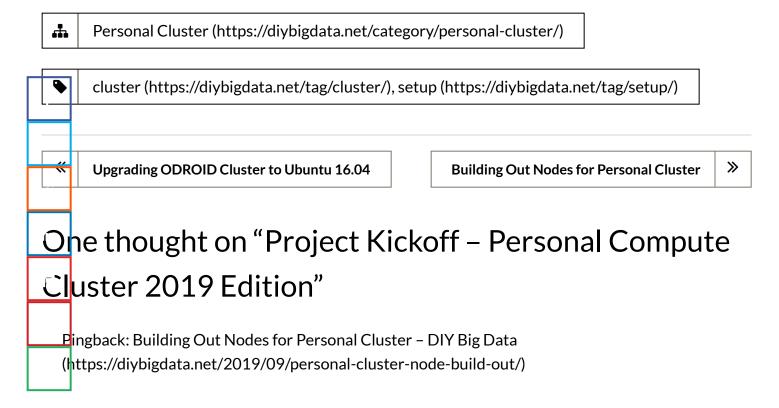
or der to complete the cluster, I needed to purchase some networking hardware. I plan to use the same network design as my last cluster (https://diybigdata.net/2016/06/network-design-for-the-low-cost-cluster/), so all I really needed to get is a switch, cables to connect the node to it, and a USB of the rnet transceiver to connect the cluster to the outside world. Altogether, my cluster build looked like this:

Item	Source	Price per Item	Item Count for Cluster
EGLOBAL S200 Computer Node built out as described above	See Above	\$823.21	4
5 Port Ethernet Switch	<u>Amazon (https://amzn.to/2LqQrte)</u>	\$14.99	1
1 ft Cat 6 cables	<u>Amazon (https://amzn.to/30pISqR)</u>	\$15.99	1
USB 3 Type C Ethernet Adapter	Amazon (https://amzn.to/2XBzuTp)	\$15.99	1
Power Strip	Amazon (https://amzn.to/2xiVAv8)	\$19.99	1

I used an existing ethernet cable I had to connect the cluster to my home network, but you don't have one, but sure to pick that up too. This brings my total cluster cost to \$3359.80. I feel that's not bad for a cluster with 48 thread, 256 GB of RAM, and 4 TB of disk space.

UPDATE – Since building out and using this cluster, I have noted that the CPU cooler that the S200 nodes came with is insufficient when running a node at 100% load across all CPU cores. I ended up upgrading the CPU cooler, <u>as documented here (https://diybigdata.net/2020/01/improving-cpu-cooling-of-eglobal-s200/)</u>. The upgraded coolers add a total of \$95.77 to the cost fo the cluster.

Next up will be the physical build out and initial setup of the cluster.



Leave a Reply

You must be logged in (https://diybigdata.net/wp-login.php? redirect_to=https%3A%2F%2Fdiybigdata.net%2F2019%2F09%2Fpersonal-compute-cluster-2019-project-kickoff-2%2F) to post a comment.



Improving Linux Kernel Network Configuration for Spark on High Performance Networks (https://diybigdata.net/2020/06/tweaks-for-spark-on-high-speed-ethernet-networks/)

Identifying Bot Commenters on Reddit using Benford's Law (https://diybigdata.net/2020/03/using-benfords-law-to-identify-bots-on-reddit/)

Upgrading the Compute Cluster to 2.5G Ethernet (https://diybigdata.net/2020/03/upgrading-cluster-to-2-5g-ethernet/)

Benchmarking Software for PySpark on Apache Spark Clusters (https://diybigdata.net/2020/01/pyspark-benchmark/)

Improving the cooling of the EGLOBAL S200 computer (https://diybigdata.net/2020/01/improving-cpu-cooling-of-eglobal-s200/)

Archives

June 2020 (https://diybigdata.net/2020/06/)

March 2020 (https://diybigdata.net/2020/03/)

<mark>⊭</mark>nuary 2020 (https://diybigdata.net/2020/01/)

December 2019 (https://diybigdata.net/2019/12/)

October 2019 (https://diybigdata.net/2019/10/)

September 2019 (https://diybigdata.net/2019/09/)

November 2017 (https://diybigdata.net/2017/11/)

January 2017 (https://diybigdata.net/2017/01/)

November 2016 (https://diybigdata.net/2016/11/)

October 2016 (https://diybigdata.net/2016/10/)

September 2016 (https://diybigdata.net/2016/09/)

August 2016 (https://diybigdata.net/2016/08/)

July 2016 (https://diybigdata.net/2016/07/)

June 2016 (https://diybigdata.net/2016/06/)



Categories

Computer Science (https://diybigdata.net/category/general/computer-science/)

Data Analysis (https://diybigdata.net/category/low-cost-cluster/data-analysis/)

Data Analysis (https://diybigdata.net/category/general/data-analysis-general/)

General (https://diybigdata.net/category/general/)

Hardware (https://diybigdata.net/category/low-cost-cluster/hardware/)

Low Cost Cluster (https://diybigdata.net/category/low-cost-cluster/)

Personal Cluster (https://diybigdata.net/category/personal-cluster/)

Set Up (https://diybigdata.net/category/low-cost-cluster/set-up/)

Spark Performance (https://diybigdata.net/category/general/spark-performance/)

Uncategorized (https://diybigdata.net/category/uncategorized/)

Meta

Log in (https://diybigdata.net/wp-login.php)

Entries feed (https://diybigdata.net/feed/)

Comments feed (https://diybigdata.net/comments/feed/)

WordPress.org (https://wordpress.org/)

Copyright © 2016 by Michael F. Kamprath. All rights reserved.

Proudly powered by WordPress (http://wordpress.org/) | WEN Business by WEN Themes (https://wenthemes.com/)

