

```
elif operation == "MIRROR_Y":
    mirror_mod.use_x = False
    mirror_mod.use_y = True
    mirror_mod.use_z = False
elif operation == "MIRROR_Z":
    mirror_mod.use_x = False
    mirror_mod.use_y = False
    mirror_mod.use_z = True

#selection at the end -add back the deselected mirror modifier object
mirror_ob.select= 1
modifier_ob.select=1
bpy.context.scene.objects.active = modifier_ob
print("Selected" + str(modifier_ob)) # modifier ob is the active ob
#mirror_ob.select = 0
name = bpy.context.selected_objects[0]
bpy.data.objects[name].select = 1
```

Apache Spark over Frontera cluster

Earlier Experiences

- Motivation
- Overview
- Frontera
- Apache Spark
- Benchmark
- Performance optimization
- Thanks

- Bring the popular Apache Spark processing engine to HPC
- Run the software without any cluster side setup requirements
- Bypass rootless execution environments limitations
- Provide an optimized Apache Spark configuration for Frontera cluster

- Apache Spark distributed computing engine:
 - task dispatching
 - scheduling
 - communication over compute cluster
- Spark cluster managers requiring dedicated clusters:
 - Hadoop YARN
 - Apache MESOS
 - Kubernetes
- Rootless execution context challenges:
 - software dependencies hard to keep over HPC environment
 - no guarantee that same software list is available on each cluster
 - restricted service deployment in operating system user space

Frontera HPC Cluster

The experience takes place over Frontera primary computing system:

- Provided by Dell EMC with 8008 nodes
- Intel Xeon Platinum 8280 “Cascade Lake”
- 56 cores/node - 28 cores/socket
- Clock rate 2.7Ghz (Base Frequency)
- Peak Node Performance 4.8TF, Double Precision
- Memory/Node 192GB DDR-4
- Local Disk with 480GB SSD drive
- Mellanox Infiniband, HDR-100, capable of 200Gbps



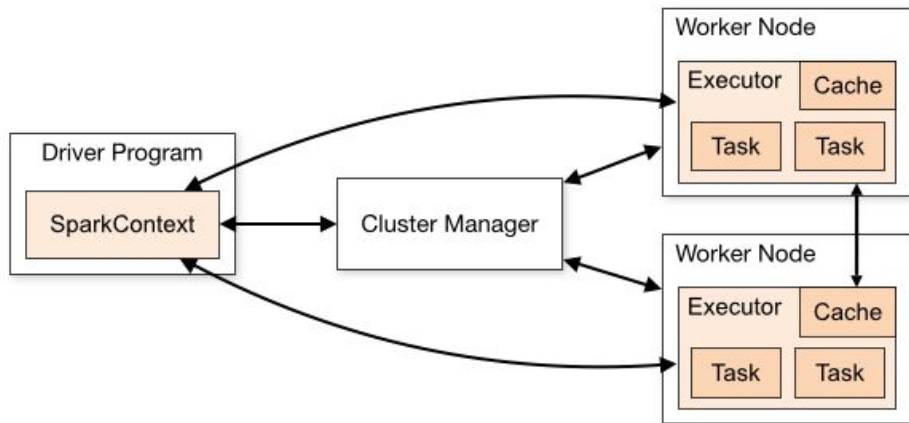
The experience takes place over Frontera primary computing system:

- LINPACK benchmark: 23.5PF
- Theoretical peak performance: 38.7PF
- Storage composed by LUSTRE distributed file system:
 - based storage of 60PB
 - fast flash storage of 3PB
- LUSTRE OSTs (object storage targets):
 - based storage: home, work and scratch
 - flash storage: flash
- home: code and collected logs
- scratch: application workspace

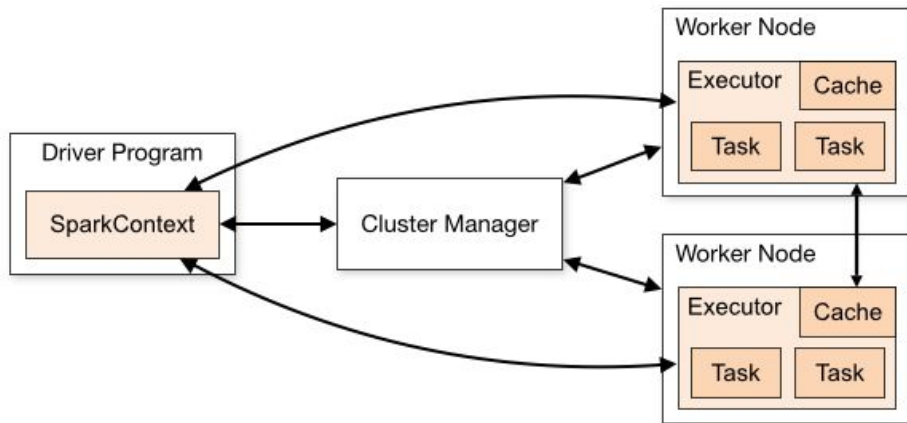


- Analytics engine for large-scale data processing
- High level API in multiple languages (our example uses python)
- MLlib provide tools for machine learning
- RDD: resilient distributed dataset
 - working set for distributed programs over distributed shared memory
 - bypass limitations in the MapReduce cluster computing paradigm such as forced linear dataflow structure
- Processing workflow managed as a directed acyclic graph (DAG)
 - nodes are the RDDs
 - edges are the applied operations over the RDDs
- Handle acyclic and cyclic graphs that allows representation of iterative methods
- DAG iterative methods representation bypass another MapReduce limitation, providing a solution to express iterative algorithms

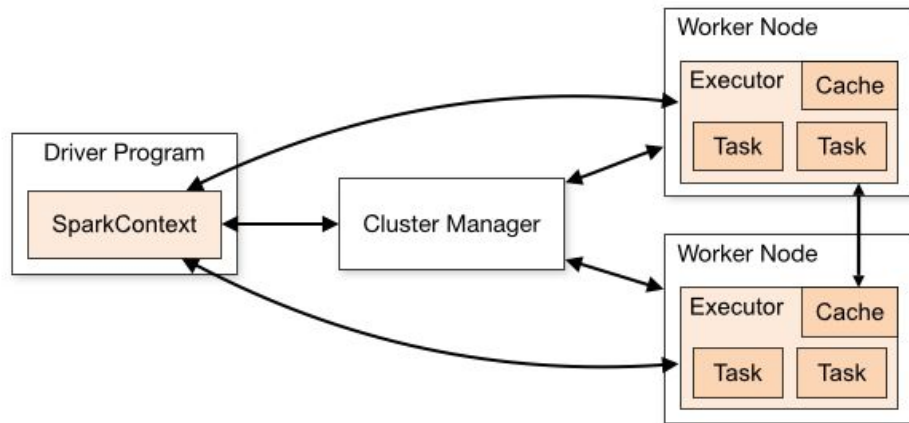
- Engine deployment requires
 - a cluster manager
 - distributed storage system
- Cluster Manager: standalone Spark cluster
 - allows better control over the cluster operation
 - bypass rootless environment limitations
 - allocate resources across applications
 - send application code to the executors
- SparkContext: Spark driver program
 - responsible to send application code and its tasks to the executors
- Executors: spawn on the nodes
 - run the computations
 - manage the application data



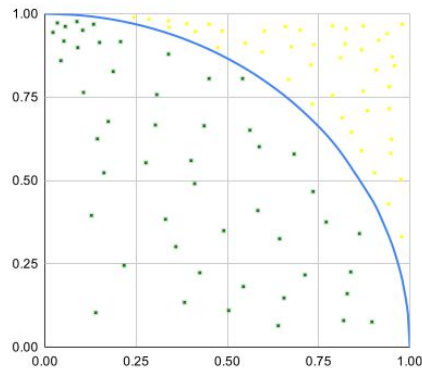
- Block Manager:
 - application data cache
 - key-value store for data blocks
 - acts as a local cache for the driver and executors on every node
 - provides interface for upload and fetch blocks locally and remotely using memory, disk and external block stores
- Client execution modes:
 - cluster mode: driver process is launched on a worker node
 - client mode: driver remains on the client node that submitted the application
 - local mode: all processes of the application run on a single machine



- Local execution mode is used in our base tests with a single node
- With implemented Cluster Manager support scripts for multi-node tests it is possible to achieve cluster mode, bypassing client mode limitation of standalone deployment
- Support scripts are developed using bash command language:
 - manage deployment jobs
 - always available in any site
 - lesser maintenance effort since command language assure regression testing for any release
- Use aptainer to deploy the required software



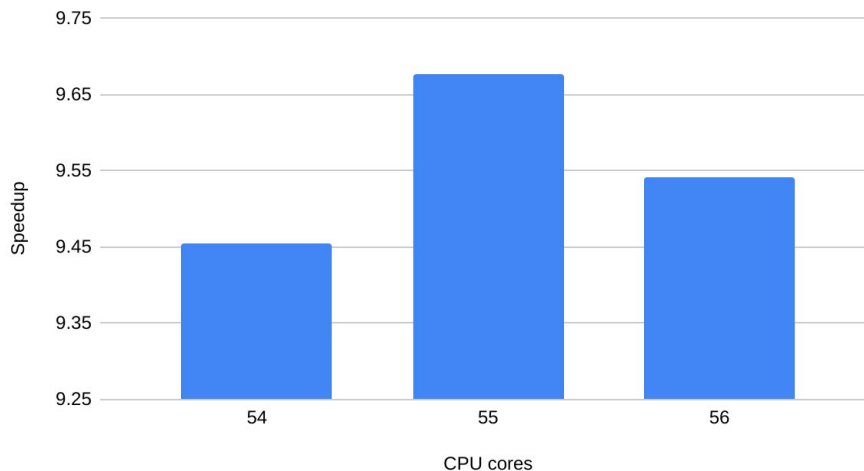
- Pi value estimation with Monte Carlo sampling methods
 - direct sampling
 - importance sampling
 - rejection sampling
- Best fit method for Pi value estimation is rejection sampling
 - allows to select samples within a region of the sampled distribution
 - simulate random points in a 2-D plane
 - domain is a square around a circle
 - square size equal the diameter of the circle
 - random number of points are generated inside the square
 - estimation is the ratio between points inside the circle and total number of generated points



Single Node

- Base test of 10^{10} samples
- This test uses all resources of a single node
- Algorithm tends to be more CPU intensive
- Best speedup when leaving one core available for the operating system and driver for this environment

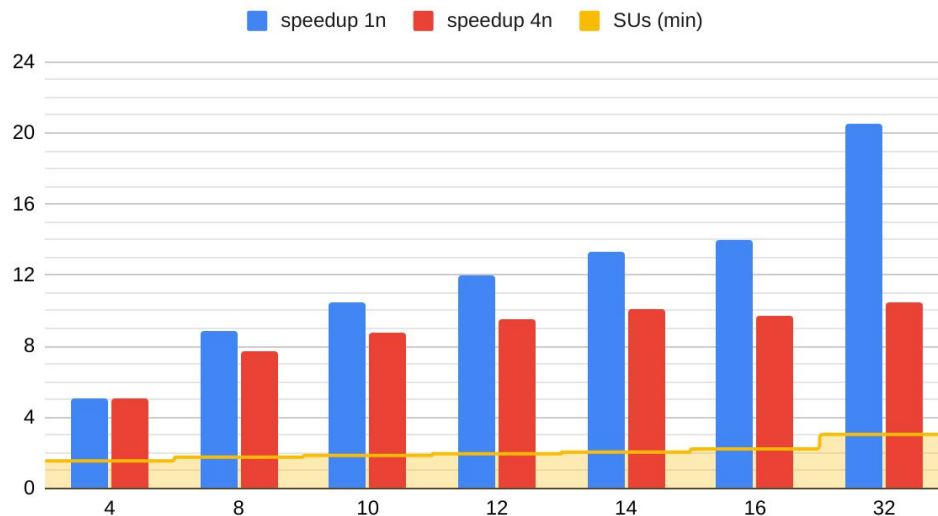
Speedup vs. CPU cores



Executors multi-node scaling

- We only use 140GB (at least 2.5GB per task) to keep executor memory equal in all nodes
- Start to scale out the single node test sample to different number of nodes
- We noted the superscalar effect of Intel processor caches that brought great speedups

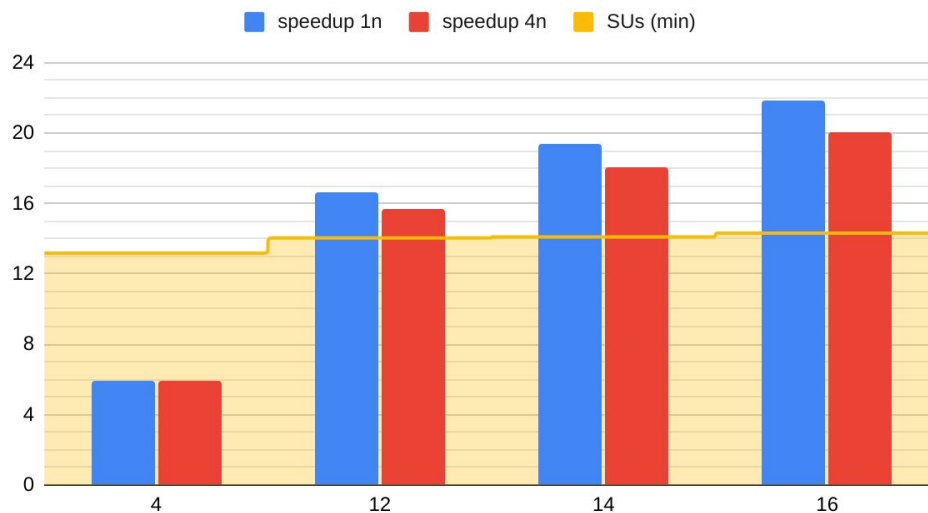
Speedup using Apache Spark Cluster for 10^{10} samples



Executors multi-node scaling

- Next step scaled up 10 times the sample
- This can still run over the 4 nodes available memory
- The performance with this sample have better results now with more nodes
- Communication takes lesser overhead with applied optimizations

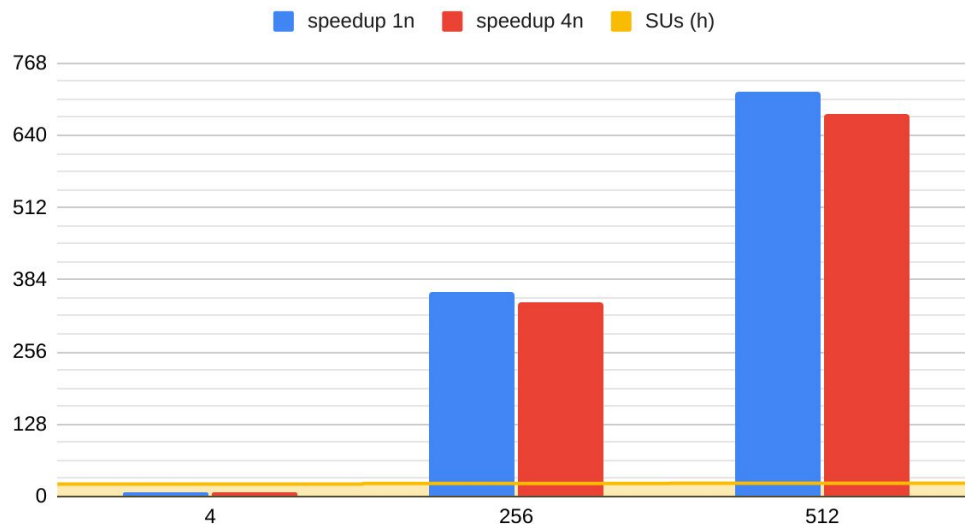
Speedup using Apache Spark for 10^{11} samples



Executors multi-node scaling

- Next step scaled up 1000 times the initial sample
- To run this simulation we need at least 256 nodes because of required memory size
- Here the superscalar effect overwhelms the previous results, turning evident the Frontera execution cost optimization for this job sizes

Speedup using Apache Spark for 10^{13} samples



- BigHPC consortium for supporting my research work on this case study
- To TACC team a special thanks with the research work around Frontera
- To LIP that is providing in Portugal the access to another HPC cluster to continue this research work and validate the optimizations with different hardware

BIG HPC HIGH PERFORMANCE COMPUTING

Partners:



Funding:

Co-financiado por:

