

Large-scale Parallel Computing

Outline

- Parallel program evaluation example
- Scaling-out methods
- High-performance cluster
- Big data cluster

Performance Evaluation of Matrix Multiplication

```
import numpy as np
import datetime
start = datetime.datetime.now()

size = 10000
a = np.random.random_sample((size, size))
b = np.random.random_sample((size, size))
n = np.dot(a,b)

end = datetime.datetime.now()
print(end - start, "h:mm:ss")
```

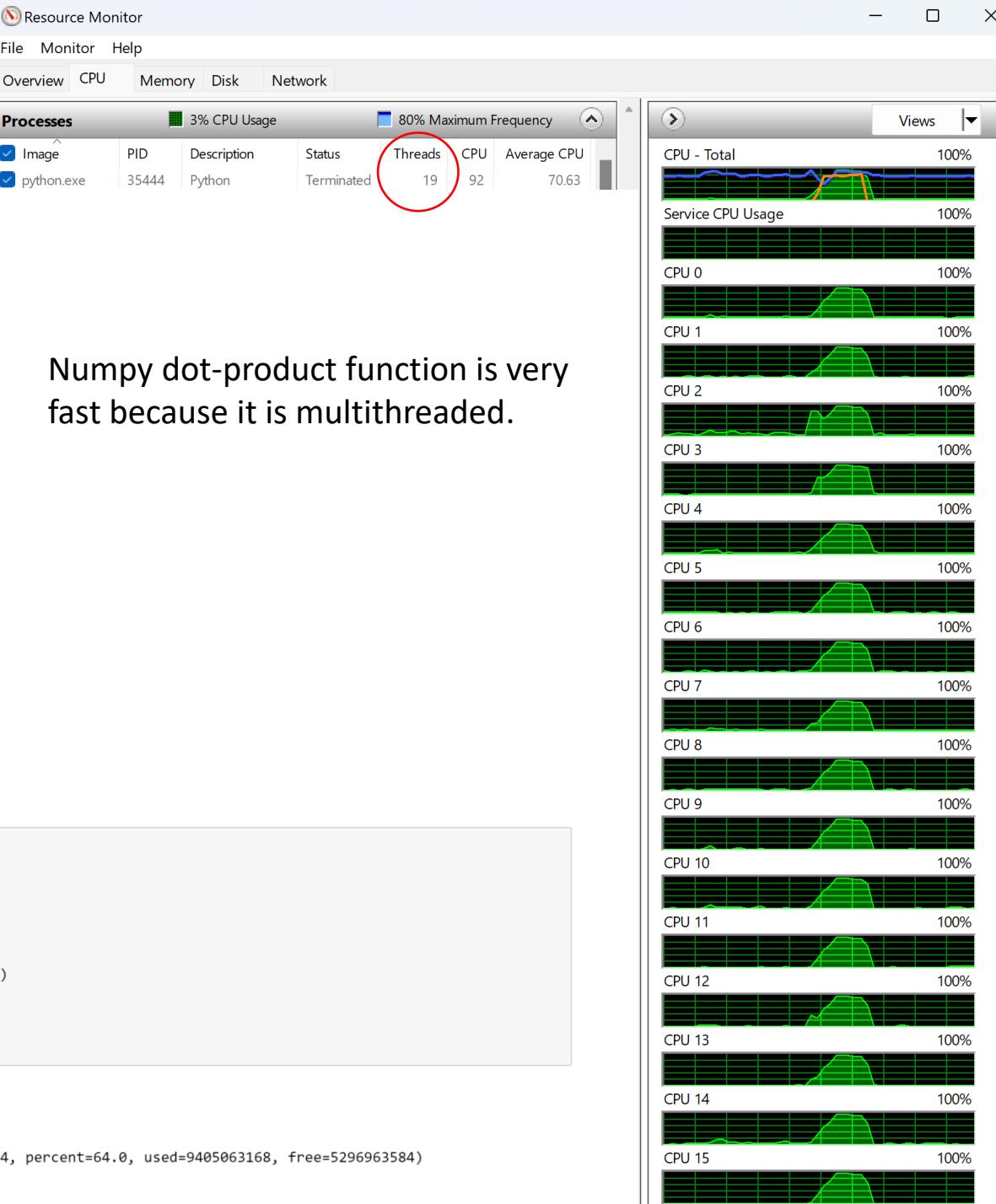
```
PS C:\Users\veera\Documents\2 Teaching\625
0:00:09.250277 h:mm:ss
```

Getting CPU information in Python

```
!pip install -q psutil
!pip install -q py-cpuinfo
import platform
import psutil
import cpuinfo

print('platform      :', platform.platform())
print('cpu model    :', cpuinfo.get_cpu_info()['brand_raw'])
print('physical cpu :', psutil.cpu_count(logical=False))
print('logical cpu  :', psutil.cpu_count(logical=True))
print('virtual memory:', psutil.virtual_memory())

platform      : Windows-10-10.0.22621-SP0
cpu model    : AMD Ryzen 7 6800HS Creator Edition
physical cpu : 8
logical cpu  : 16
virtual memory : svmem(total=14702026752, available=5296963584, percent=64.0, used=9405063168, free=5296963584)
```



Parallel Computing in Numpy

Numpy uses several libraries that are multithreaded and uses SIMD (vector) instructions.

```
import numpy as np
np.show_config()

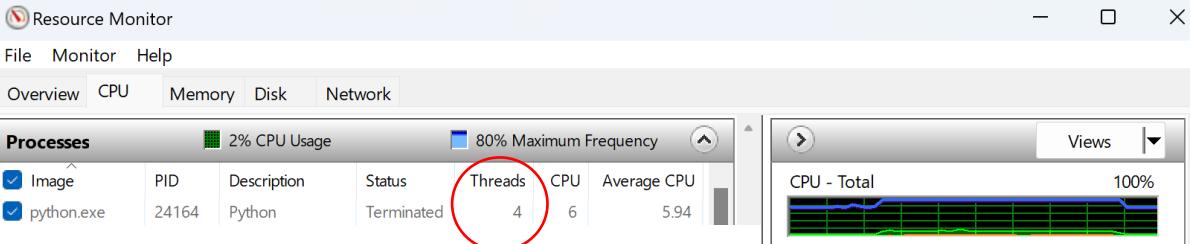
openblas64_info:
    library_dirs = ['D:\\a\\numpy\\numpy\\build\\openblas64_info']
    libraries = ['openblas64_info']
    language = f77
    define_macros = [('HAVE_CBLAS', None), ('BLAS_SYMBOL_SUFFIX', '64_'), ('HAVE_BLAS_ILP64', None)]
blas_ilp64_opt_info:
    library_dirs = ['D:\\a\\numpy\\numpy\\build\\openblas64_info']
    libraries = ['openblas64_info']
    language = f77
    define_macros = [('HAVE_CBLAS', None), ('BLAS_SYMBOL_SUFFIX', '64_'), ('HAVE_BLAS_ILP64', None)]
openblas64_lapack_info:
    library_dirs = ['D:\\a\\numpy\\numpy\\build\\openblas64_lapack_info']
    libraries = ['openblas64_lapack_info']
    language = f77
    define_macros = [('HAVE_CBLAS', None), ('BLAS_SYMBOL_SUFFIX', '64_'), ('HAVE_BLAS_ILP64', None), ('HAVE_LAPACKE', None)]
lapack_ilp64_opt_info:
    library_dirs = ['D:\\a\\numpy\\numpy\\build\\openblas64_lapack_info']
    libraries = ['openblas64_lapack_info']
    language = f77
    define_macros = [('HAVE_CBLAS', None), ('BLAS_SYMBOL_SUFFIX', '64_'), ('HAVE_BLAS_ILP64', None), ('HAVE_LAPACKE', None)]
Supported SIMD extensions in this NumPy install:
baseline = SSE,SSE2,SSE3
found = SSSE3,SSE41,POPCNT,SSE42,AVX,F16C,FMA3,AVX2
not found = AVX512F,AVX512CD,AVX512_SKX,AVX512_CLX,AVX512_CNL,AVX512_ICL
```

```
!pip install threadpoolctl

from threadpoolctl import threadpool_info
from pprint import pprint
import numpy
pprint(threadpool_info())
```

```
Requirement already satisfied: threadpoolctl in c:\\users\\veera\\appdata\\local\\programs\\python\\python39\\lib\\site-packages (3.1.0)
[{'architecture': 'Zen',
  'filepath': 'C:\\\\Users\\\\veera\\\\AppData\\\\Local\\\\Programs\\\\Python\\\\Python39\\\\Lib\\\\site-packages\\\\numpy\\\\.libs\\\\libopenblas.FB5A
E2TYXYH21JRDUGUGQ3XBKUJTF43H.gfortran-win_amd64.dll',
  'internal_api': 'openblas',
  'num_threads': 16,
  'prefix': 'libopenblas',
  'threading_layer': 'pthreads',
  'user_api': 'blas',
  'version': '0.3.20'}]
```

Performance Evaluation of Matrix Multiplication



Setting the number of threads before running Numpy code.

```
import os
os.environ["OMP_NUM_THREADS"] = "1"
os.environ["OPENBLAS_NUM_THREADS"] = "1"
os.environ["MKL_NUM_THREADS"] = "1"
os.environ["VECLIB_MAXIMUM_THREADS"] = "1"
os.environ["NUMEXPR_NUM_THREADS"] = "1"

import numpy as np

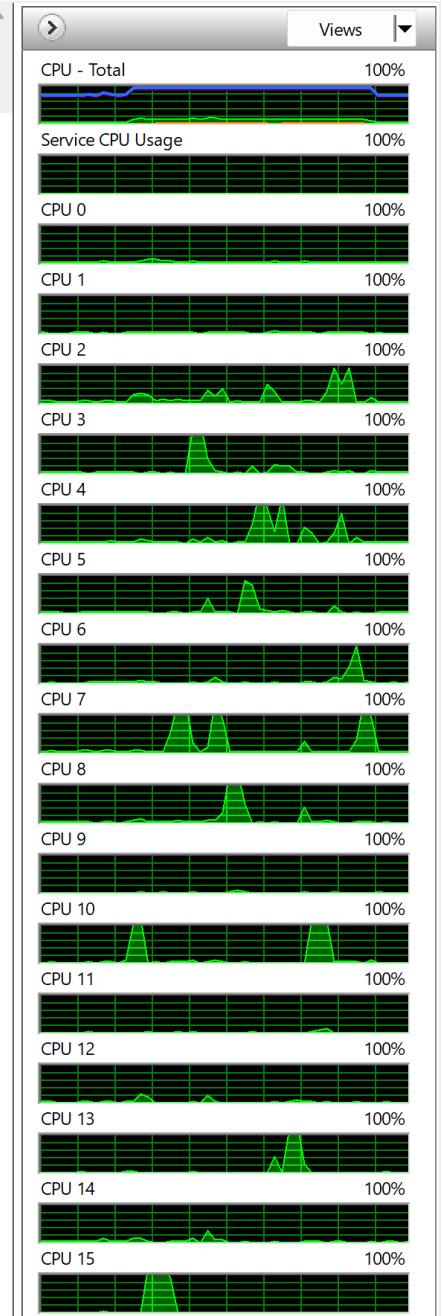
import datetime
start = datetime.datetime.now()

size = 10000
a = np.random.random_sample((size, size))
b = np.random.random_sample((size, size))
n = np.dot(a,b)

end = datetime.datetime.now()
print(end - start, "h:mm:ss")
```

```
PS C:\Users\veera\Documents\2 Teaching\625 Data Sci Arch\programs> python .\dotproduct.py
0:00:32.347729 h:mm:ss
```

```
# export OMP_NUM_THREADS=1
# export OPENBLAS_NUM_THREADS=1
# export MKL_NUM_THREADS=1
# export VECLIB_MAXIMUM_THREADS=1
# export NUMEXPR_NUM_THREADS=1
```



Performance Evaluation of Matrix Multiplication

Setting the number of threads before running Numpy code.

```
import os
os.environ["OMP_NUM_THREADS"] = "1"
os.environ["OPENBLAS_NUM_THREADS"] = "1"
os.environ["MKL_NUM_THREADS"] = "1"
os.environ["VECLIB_MAXIMUM_THREADS"] = "1"
os.environ["NUMEXPR_NUM_THREADS"] = "1"

import numpy as np

import datetime
start = datetime.datetime.now()

size = 10000
a = np.random.random_sample((size, size))
b = np.random.random_sample((size, size))
n = np.dot(a,b)

end = datetime.datetime.now()
print(end - start, "h:mm:ss")
```

```
PS C:\Users\veera\Documents\2 Teaching\625 Data Sci Arch\programs> python .\dotproduct.py
0:00:32.347729 h:mm:ss
```

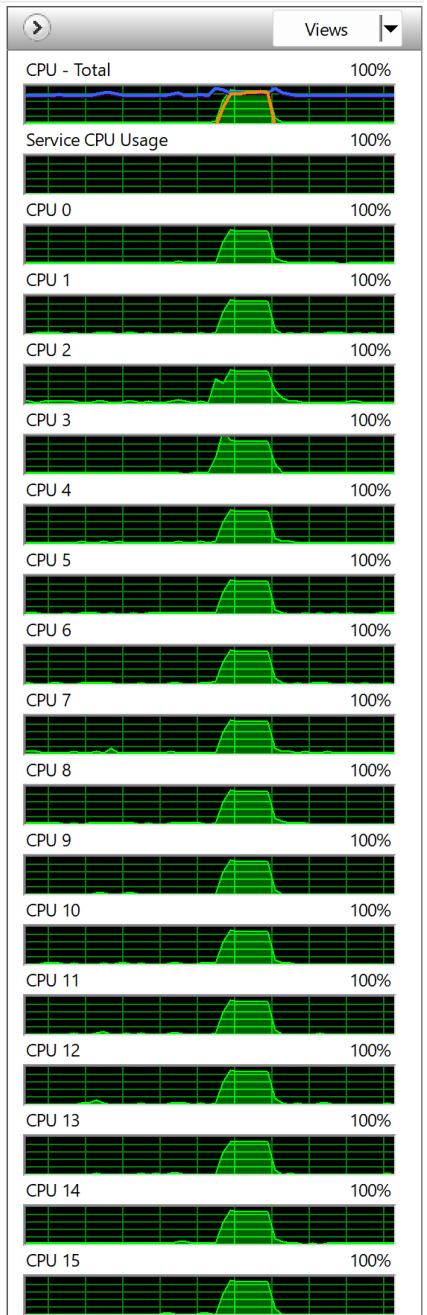
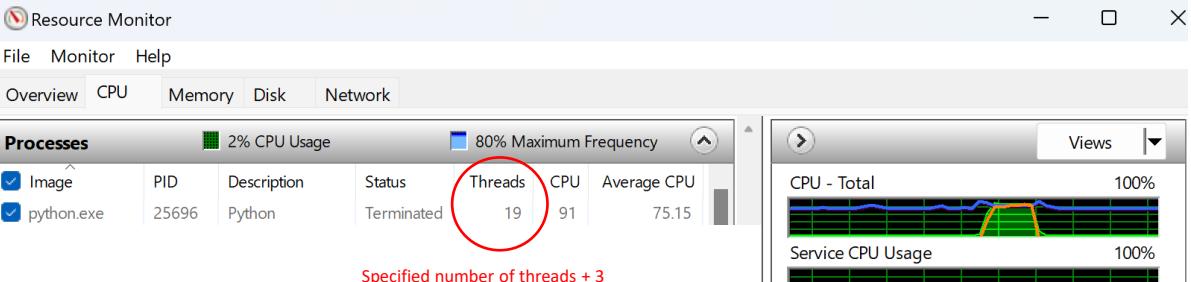
Change "1" to 2, 4, 8, 12, 16
and rerun the program

Testing on Jupyter Notebook needs restarting kernel.

```
import os
os.environ["OMP_NUM_THREADS"] = "16"
os.environ["OPENBLAS_NUM_THREADS"] = "16"
os.environ["MKL_NUM_THREADS"] = "16"
os.environ["VECLIB_MAXIMUM_THREADS"] = "16"
os.environ["NUMEXPR_NUM_THREADS"] = "16"

# export OMP_NUM_THREADS=16
# export OPENBLAS_NUM_THREADS=16
# export MKL_NUM_THREADS=16
# export VECLIB_MAXIMUM_THREADS=16
# export NUMEXPR_NUM_THREADS=16
```

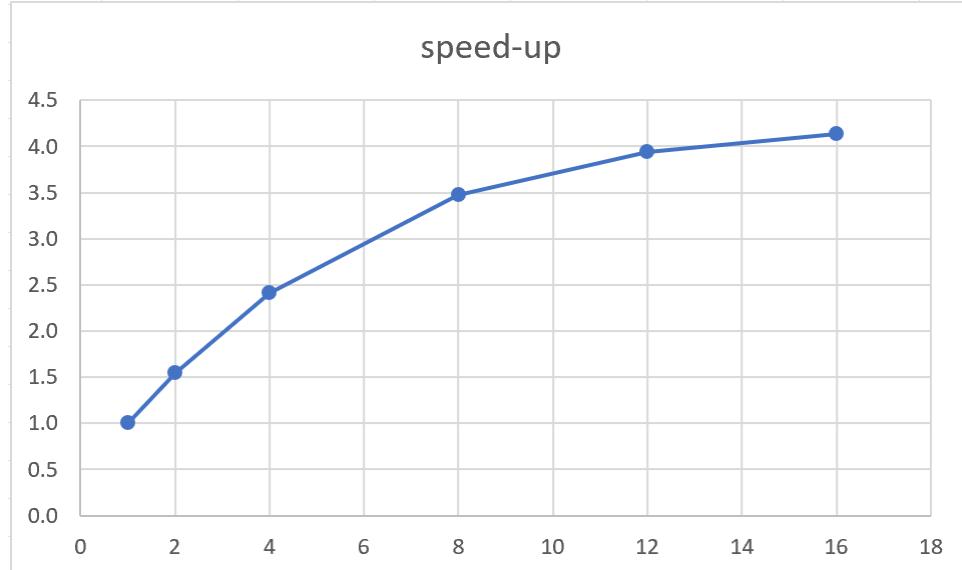
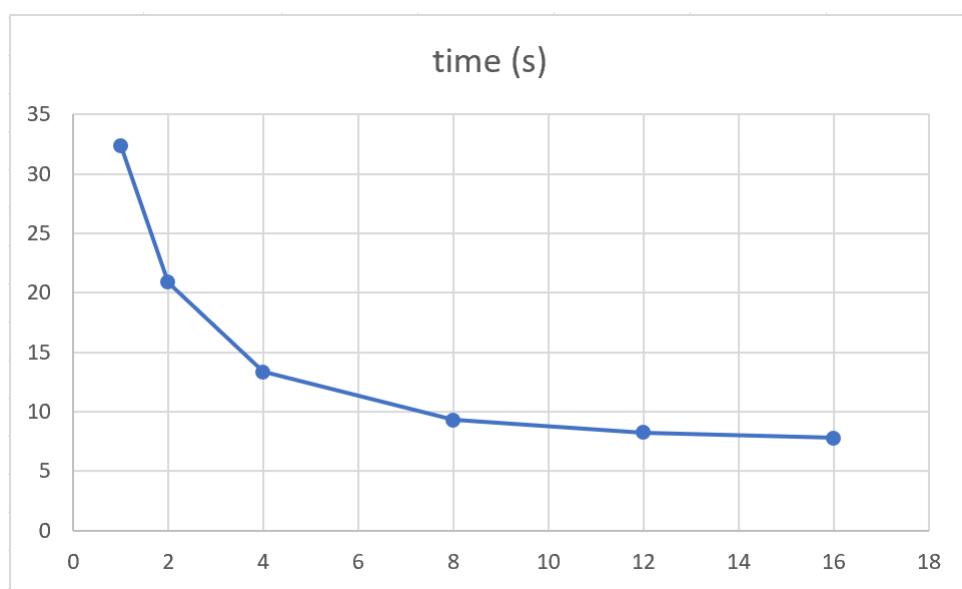
```
PS C:\Users\veera\Documents\2 Teaching\625 Data Sci Arch\programs> python .\dotproduct.py
0:00:07.793126 h:mm:ss
```



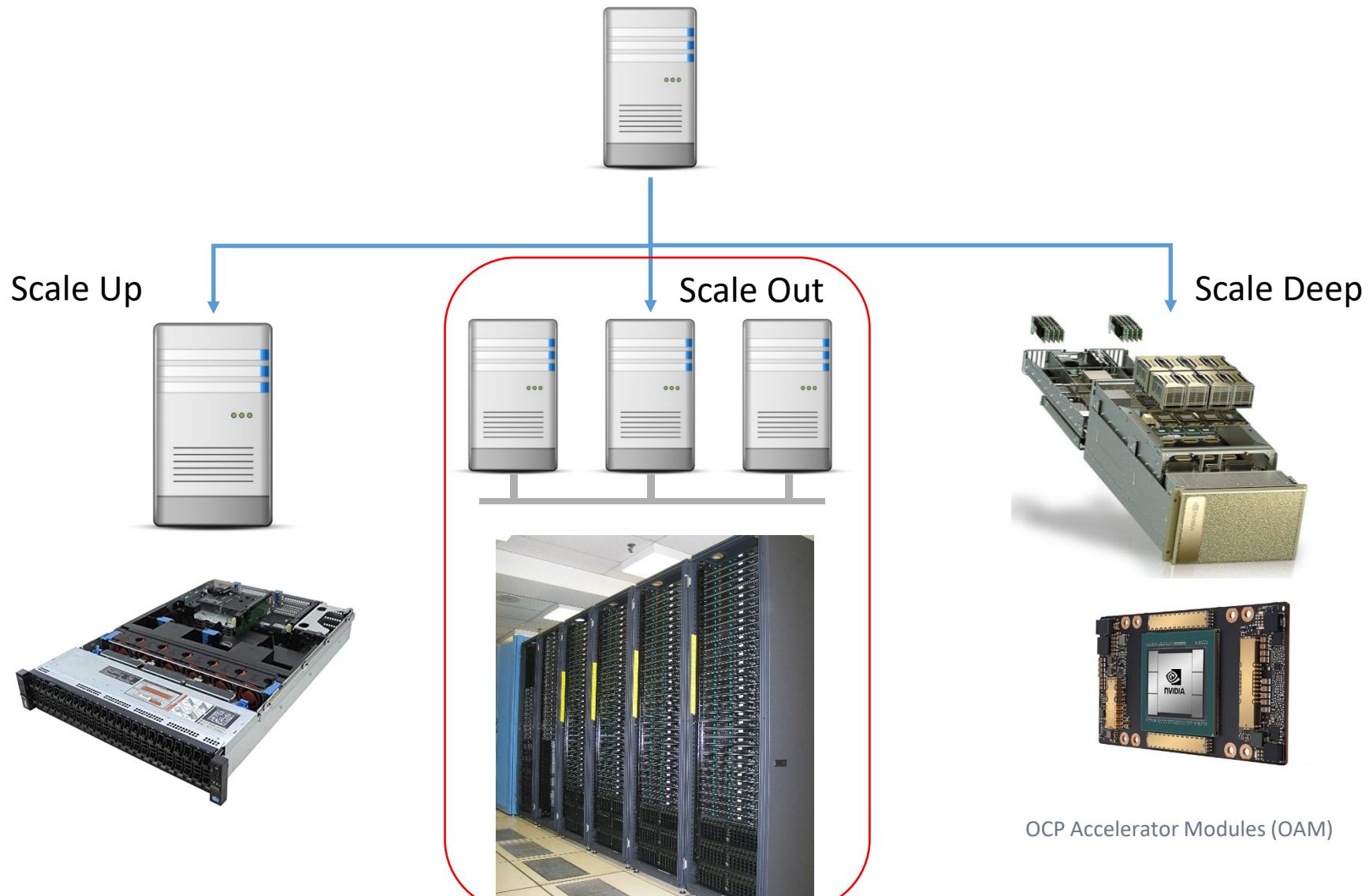
Performance Evaluation of Matrix Multiplication

threads	time (s)	speed-up
1	32.3	1.0
2	20.9	1.5
4	13.4	2.4
8	9.3	3.5
12	8.2	3.9
16	7.8	4.1

Max speed-up is about 4X



Scaling Out provides more scalability



Scale-out: Distributed Computing

Distributed computing is about making distributed components work together as one.

- **Distributed Processors** (CPU សម្រាប់ Pool CPU ដ៏អំពីរគុណភាព)

- A collection of stand-alone computers connected over a local network and working together as a single computer system

- **Distributed Shared Memory** (Spark)

- Providing a shared memory programming model over distributed memory computers.

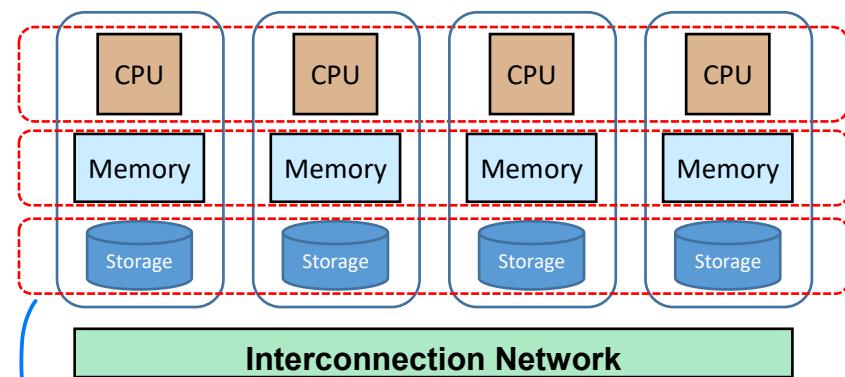
- **Distributed File System** (Hadoop)

- A shared file system consisting of storage devices that are physically distributed across machines.

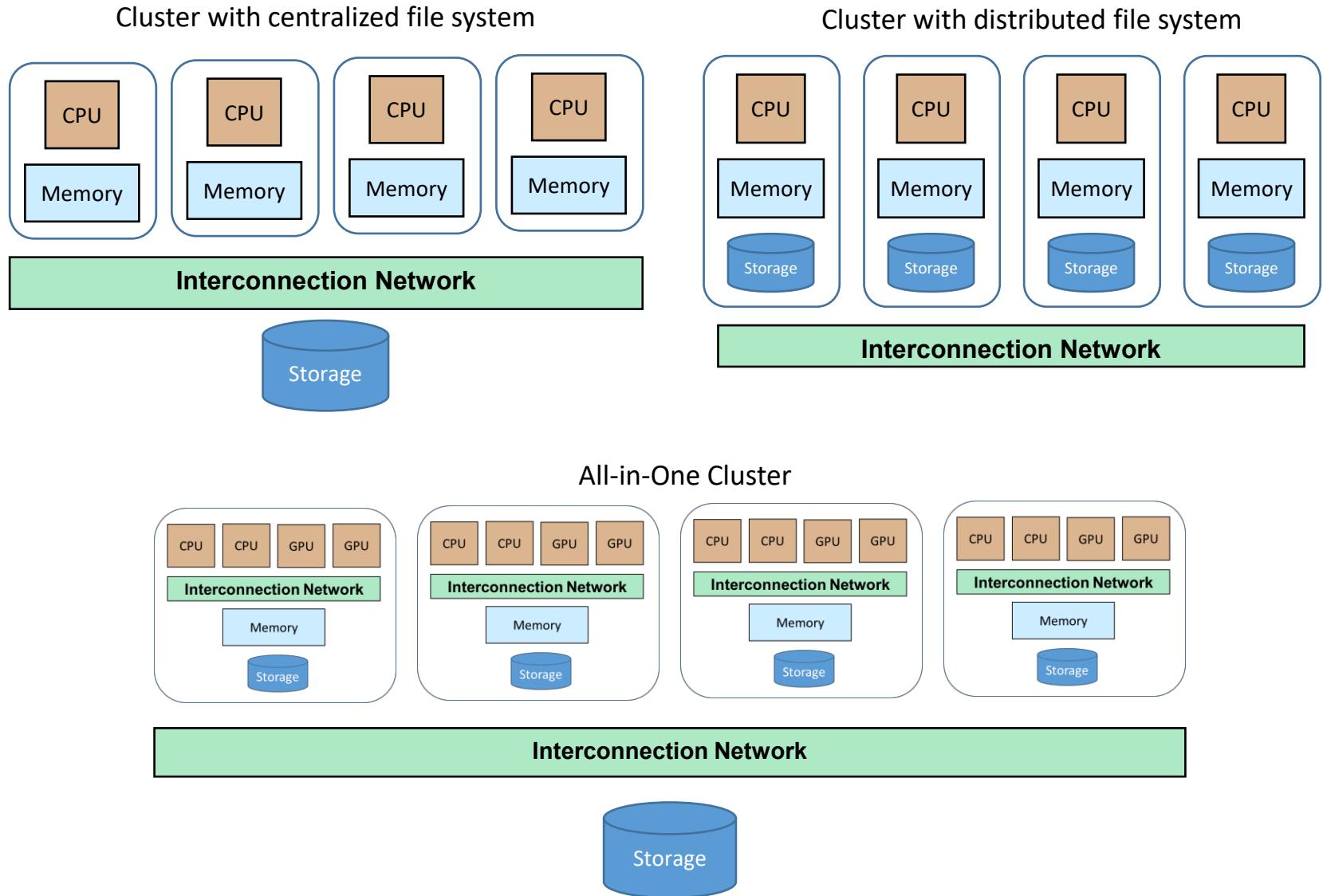
ការងារ storage ក្នុងរបៀបប្រើប្រាស់ប្រចាំថ្ងៃ

អាជីវកំណើនទិន្នន័យ

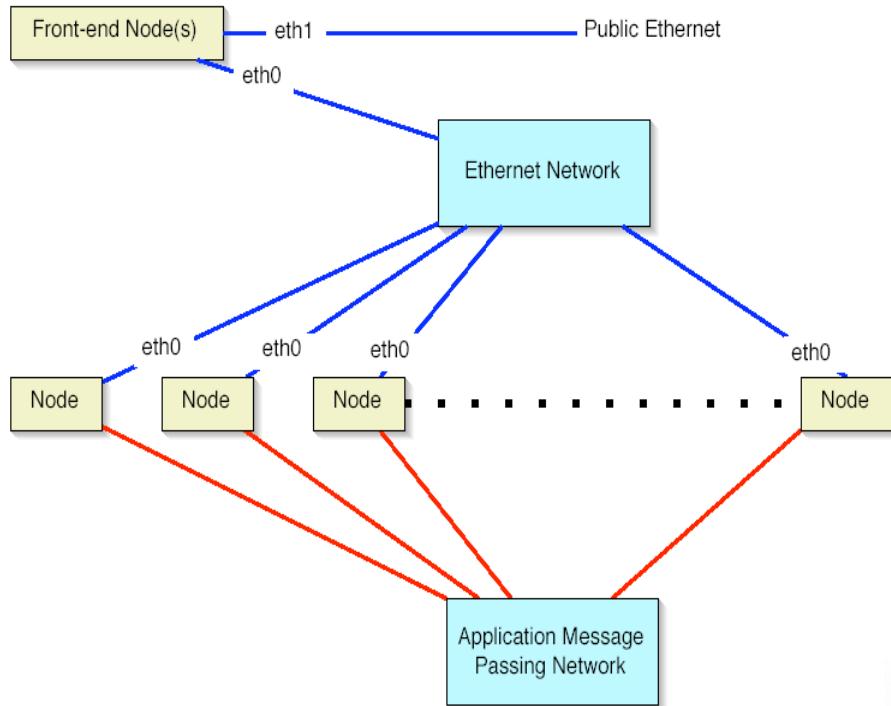
នៅ Hadoop ការងារ Parallel ចូល



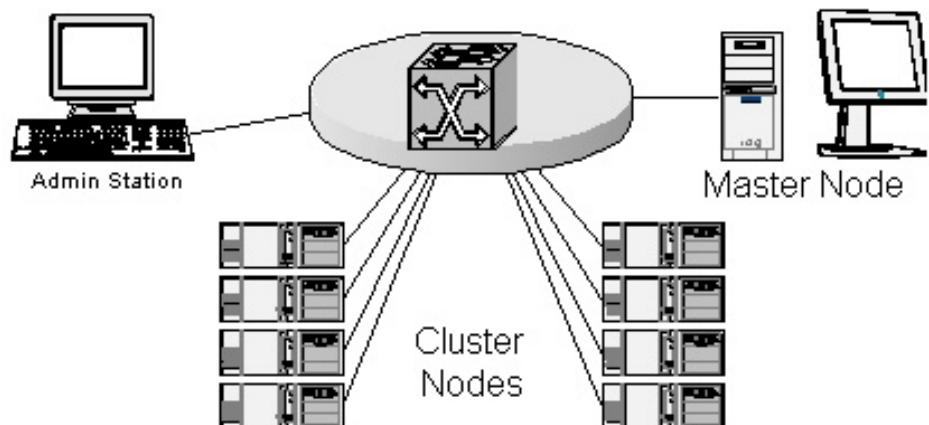
Clusters



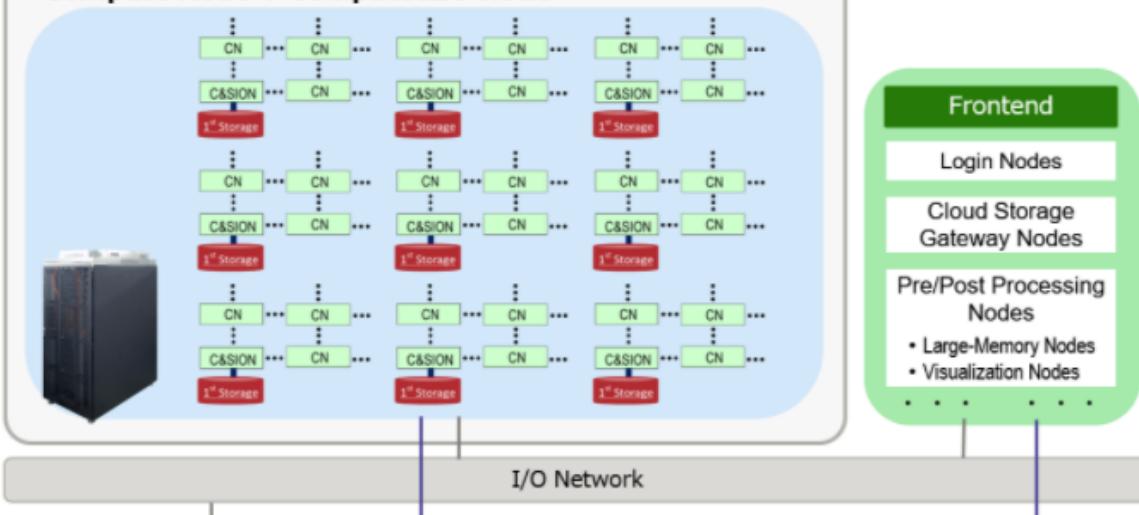
High-Performance Cluster



- CPU-intensive tasks
- Parallel computing
- Processes are distributed to low-cpu-utilization nodes
- Processes communicate via messages
- Move data to computing
- Programming Model: Message Passing



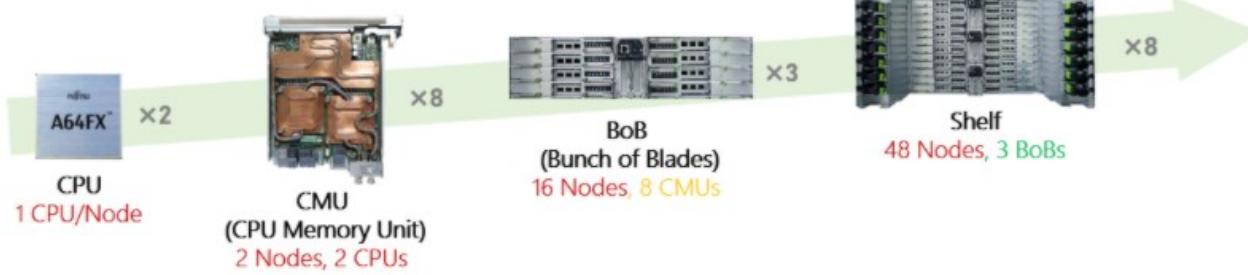
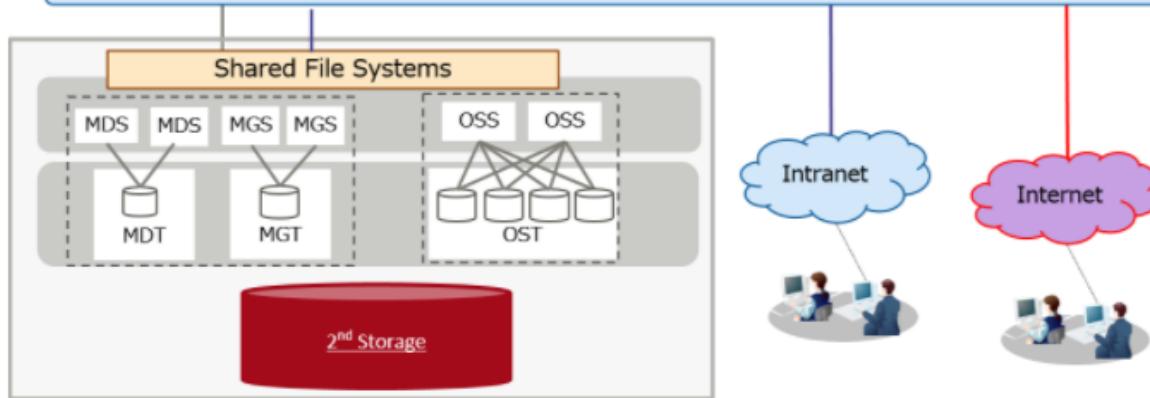
Compute Node + Compute&IO Node



Fugaku



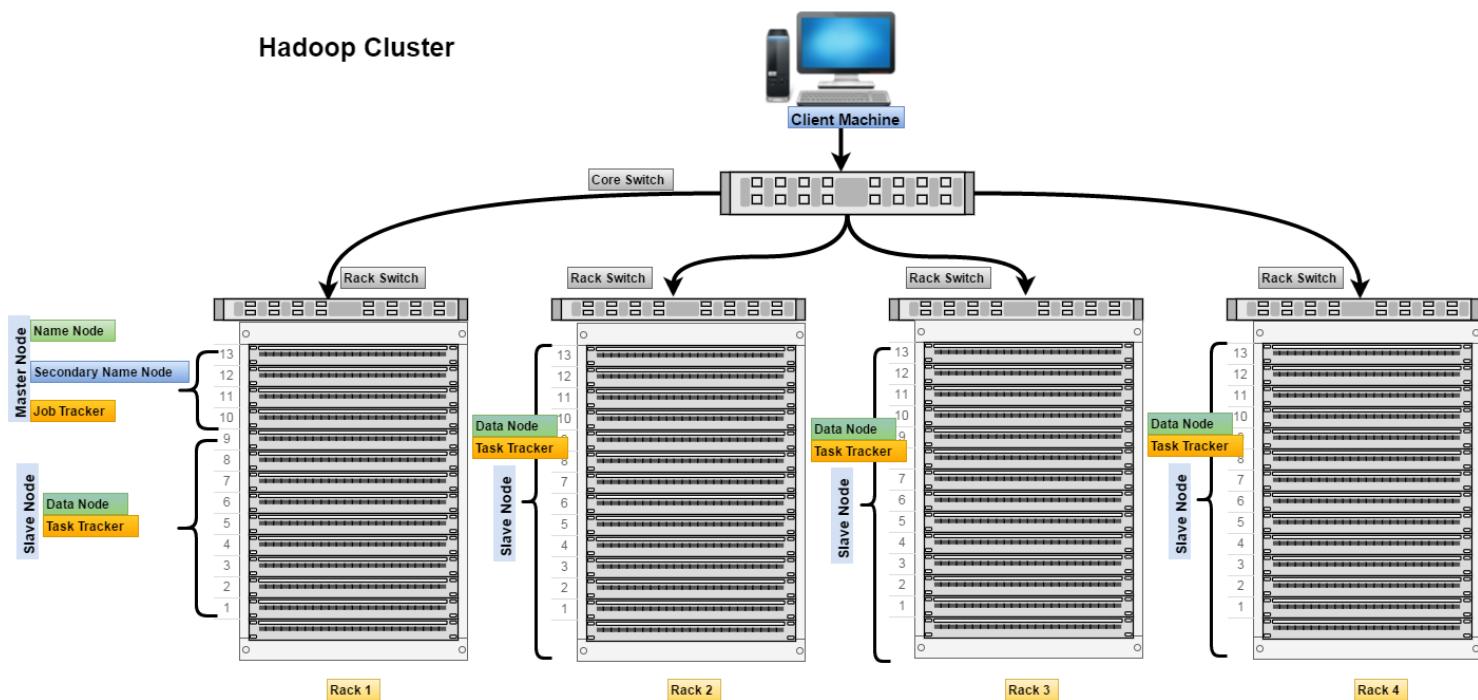
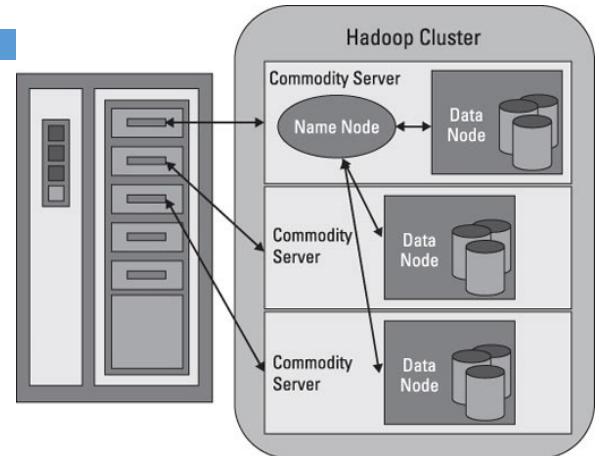
Network



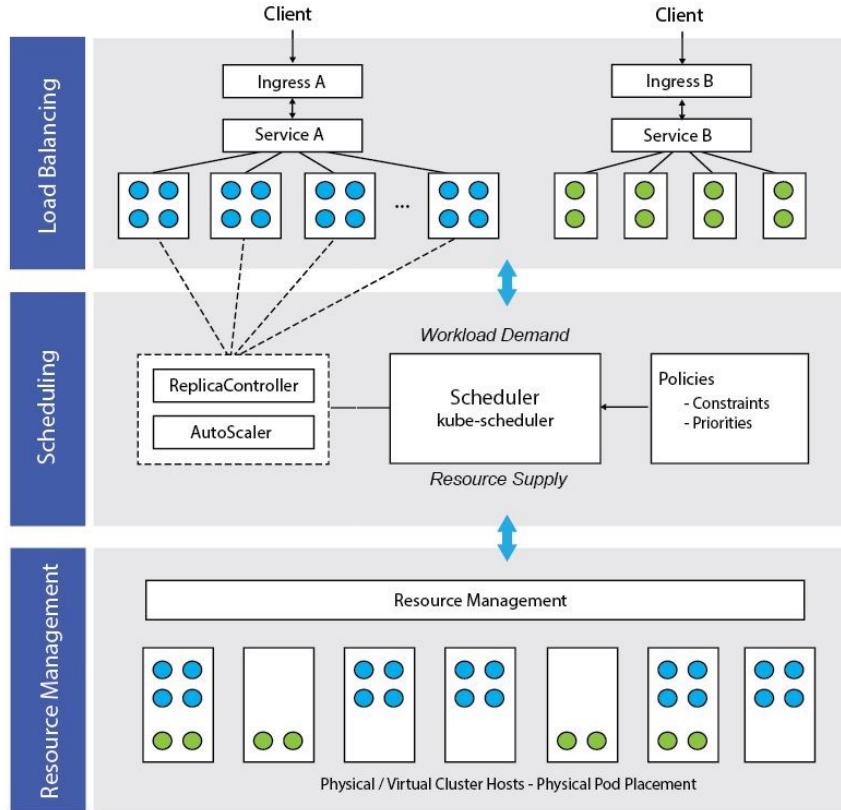
Rack
384 Nodes, 8 Shelves

Big Data Cluster

- Data-intensive computation
- High-availability distributed file system
 - Each node is a part of a distributed file system
- Jobs are distributed to nodes that contain required data
- Move computing to data
- Programming models: MapReduce, Resilient Distributed Data
- E.g. Hadoop, Spark



Container Cluster



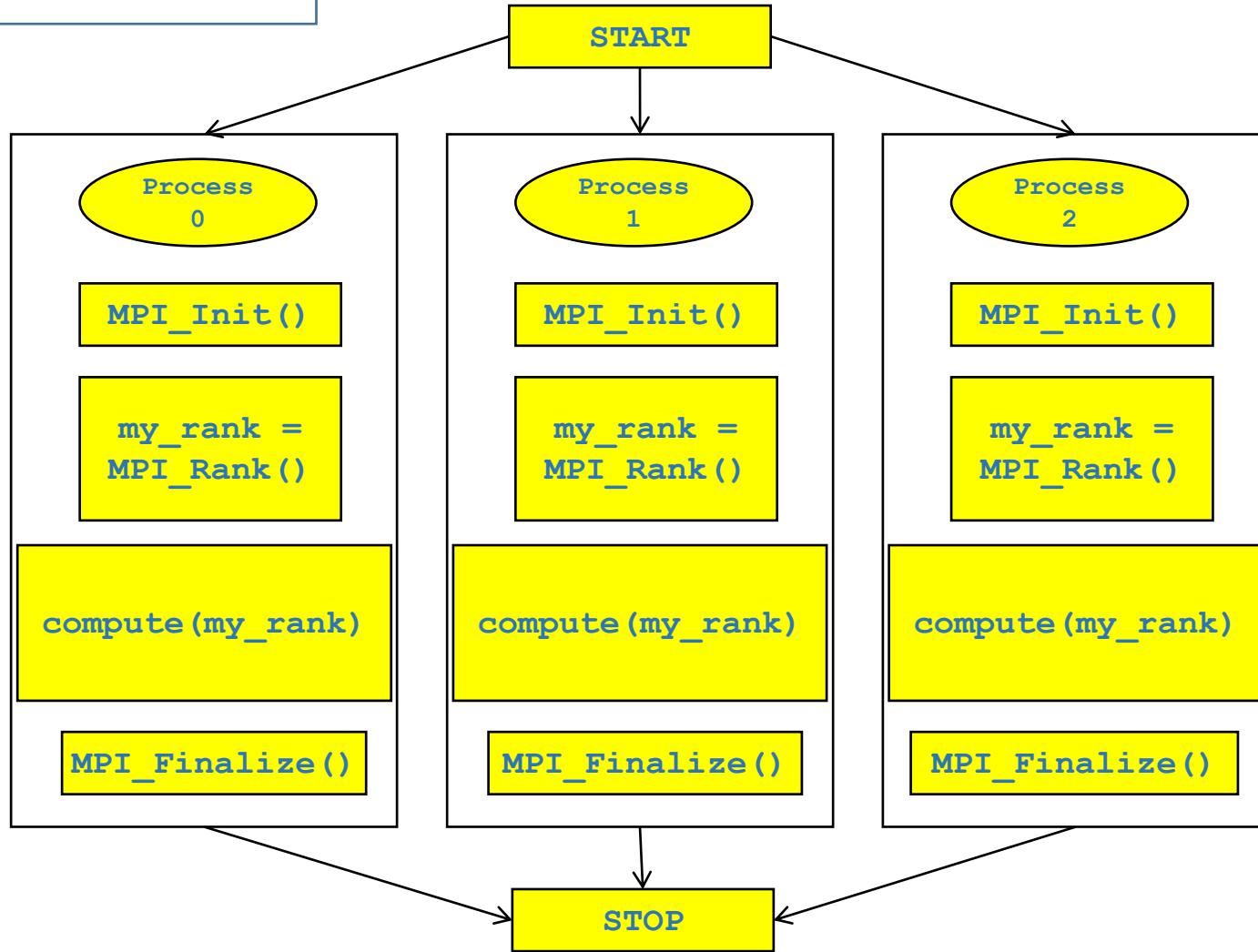
- Service-oriented applications
- Container-based
 - Process-level virtualization
- Focus on
 - Availability
 - Scaling
 - Load balancing
- Multiple virtual clusters on a single physical cluster
- E.g. Kubernetes cluster

Message-Passing Programming Model

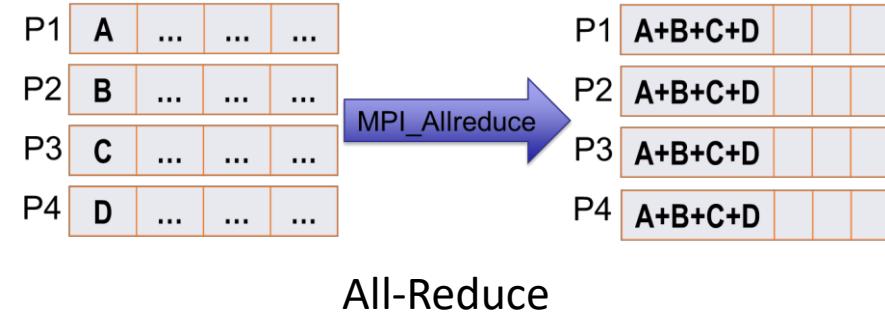
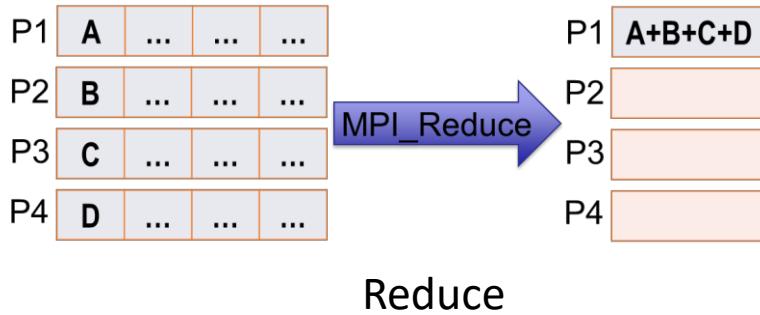
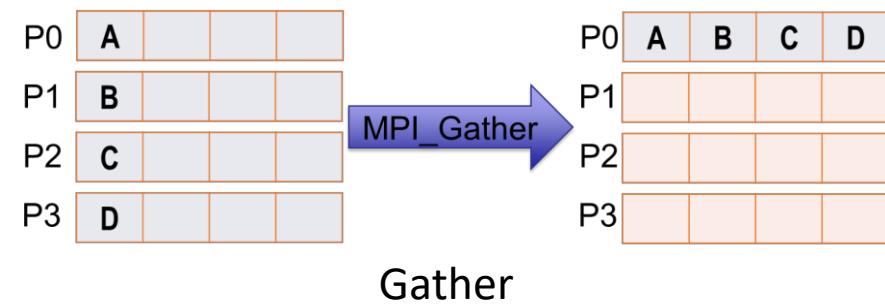
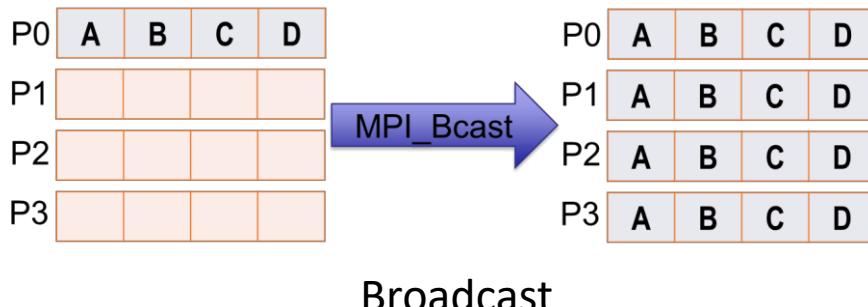
- A group of processes on one or more computers
- All processes execute the same program.
- Processes communicate and synchronize with each other by passing messages containing data.
- One-to-one and collective communication.
- MPI (Message Passing Interface) is a de facto standard.

MPI Processes

```
mpirun myprog -np 3
```



MPI Collective Communication



Example: Dot Product in MPI (in C language)

```
#include <stdio.h>
#include <mpi.h>

int main(int argc,char *argv[])
{
    double sum, sum_local, a[256], b[256];
    int n=256, i, numprocs, id, first, last;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &id);
    first = id * n/numprocs;
    last = (id + 1) * n/numprocs;
    for (i = 0; i < n; i++) {
        a [i] = i * 0.5;
        b [i] = i * 2.0;
    }
    sum_local = 0;
    for (i = first; i < last; i++) {
        sum_local = sum_local + a[i]*b[i];
    }
    MPI_Reduce(&sum_local, &sum, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
    if (id == 0) printf ("sum = %f\n", sum);
}
```

Example: Parallel Matrix Multiplication with MPI

```
from mpi4py import MPI
import numpy as np

comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()

# Dimensions of the matrices A, B
n = 200

# Generate matrices A, B on root process
if rank == 0:
    A = np.random.rand(n, n)
    B = np.random.rand(n, n)
else:
    A = None
    B = None

# Broadcast matrix B to all processes
B = comm.bcast(B, root=0)

# Scatter rows of matrix A to all processes
rows = np.empty((n // size, n), dtype='float64') # Prepare a buffer for received rows
comm.Scatter(A, rows, root=0)

# Perform local computation of matrix multiplication
local_result = np.dot(rows, B)

# Gather the local results back to the root process
result = None
if rank == 0:
    result = np.empty((n, n), dtype='float64') # Prepare a buffer to receive the full result matrix
    comm.Gather(local_result, result, root=0)

# Print the result in the root process
if rank == 0:
    print("Result of matrix multiplication:")
    print(result)
```

matmul_mpi.py

Dask

- Dask is a Python library that allows parallel and distributed computing.
- It enables scaling from a single machine to a cluster of machines, making it suitable for processing large datasets and complex computations.

Matrix Multiplication in Dask

Multithread version

```
import dask.array as da
from dask.distributed import Client, LocalCluster

def parallel_matrix_multiplication(n, client):
    # Create random matrices
    A = da.random.random((n, n), chunks=(n // 10, n // 10))
    B = da.random.random((n, n), chunks=(n // 10, n // 10))

    # Perform matrix multiplication
    result = da.matmul(A, B)

    result_computed = result.compute()

    return result_computed

if __name__ == "__main__":
    # Setup a LocalCluster to use threads instead of processes
    with LocalCluster(processes=False, n_workers=4, threads_per_worker=2) as cluster:
        with Client(cluster) as client:
            n = 2000 # Size of the matrix
            result = parallel_matrix_multiplication(n, client)
            print(f"Result of the matrix multiplication: \n{result}")
```

Matrix Multiplication in Dask

Distributed computing version

```
import dask.array as da
from dask.distributed import Client, performance_report

def parallel_matrix_multiplication(n, client):
    # Create random matrices
    A = da.random.random((n, n), chunks=(n // 10, n // 10))
    B = da.random.random((n, n), chunks=(n // 10, n // 10))

    # Perform matrix multiplication
    result = da.matmul(A, B)

    # Compute the result with a performance report
    with performance_report(filename="dask_report.html"):
        result_computed = result.compute()

    return result_computed

if __name__ == "__main__":
    # Connect to your Dask cluster
    client = Client("tcp://scheduler:8786") # Replace with your scheduler address

    n = 2000 # Size of the matrix
    result = parallel_matrix_multiplication(n, client)
    print(f"Result of the matrix multiplication: \n{result}")
```

Dask Cluster on Docker

docker-compose.yml

```
version: '3'

services:
  scheduler:
    image: daskdev/dask:latest
    command: dask-scheduler
    ports:
      - "8786:8786" # Dask scheduler dashboard
      - "8787:8787" # Dask diagnostics dashboard
    volumes:
      - ./data:/data

  worker:
    image: daskdev/dask:latest
    command: dask-worker scheduler:8786
    depends_on:
      - scheduler
  deploy:
    resources:
      limits:
        cpus: "0.25"
        memory: 512M

  jupyter:
    image: daskdev/dask-notebook:latest
    ports:
      - "8888:8888"
    command: start-notebook.sh --NotebookApp.token='mytoken'
    environment:
      - DASK_SCHEDULER_ADDRESS=tcp://scheduler:8786
      - JUPYTER_TOKEN=mytoken
    volumes:
      - ./notebooks:/home/jovyan
    depends_on:
      - scheduler
```

When running on a single machine, setting limits of resources for each worker will better simulate cluster scaling (adding workers → more resources). Without this setting, total resources are divided among workers. When running on a real cluster with multiple machines, this can be removed.

This configuration sets up:

- A Dask scheduler service, exposing ports 8786 (Dask dashboard) and 8787 (Diagnostics dashboard).
- A Dask worker service, which connects to the scheduler.
- An optional Jupyter notebook server with Dask integration, exposing port 8888 for web access. The notebook server automatically connects to the Dask scheduler.

Running Dask Cluster

```
docker-compose up --scale worker=4
```

Remove Dask Cluster

```
docker-compose down
```

```
docker-compose up --scale worker=1
```

localhost:8787/workers

Status Workers Tasks System Profile Graph Groups Info More...

CPU Use (%)

Memory Use (%)

name	address	nthreads	cpu	memory	limit	memory	managed	unmanag	unmanag	spilled	# fds	net read	net write	disk read	disk write
Total (1)		16	2 %	154.6 MiB	6.6 GiB	2.3 %	0.0	154.6 MiB	0.0	0.0	22	286 B	1 KiB	0	24 KiB
tcp://172.17.0.1:8787	tcp://172.17.0.1:8787	16	2 %	154.6 MiB	6.6 GiB	2.3 %	0.0	154.6 MiB	0.0	0.0	22	286 B	1 KiB	0	24 KiB

```
docker-compose up --scale worker=4
```

localhost:8787/workers

Status Workers Tasks System Profile Graph Groups Info More...

CPU Use (%)

Memory Use (%)

name	address	nthreads	cpu	memory	limit	memory	managed	unmanag	unmanag	spilled	# fds	net read	net write	disk read	disk write
Total (4)		64	3 %	619.3 MiB	26.5 GiB	2.3 %	0.0	235.4 MiB	383.9 MiB	0.0	88	1 KiB	6 KiB	0	192 KiB
tcp://172.17.0.1:8787	tcp://172.17.0.1:8787	16	2 %	155.1 MiB	6.6 GiB	2.3 %	0.0	58.8 MiB	96.3 MiB	0.0	22	286 B	1 KiB	0	48 KiB
tcp://172.17.0.1:8788	tcp://172.17.0.1:8788	16	6 %	153.4 MiB	6.6 GiB	2.3 %	0.0	58.7 MiB	94.7 MiB	0.0	22	286 B	1 KiB	0	48 KiB
tcp://172.17.0.1:8789	tcp://172.17.0.1:8789	16	2 %	155.7 MiB	6.6 GiB	2.3 %	0.0	59.1 MiB	96.6 MiB	0.0	22	287 B	1 KiB	0	48 KiB
tcp://172.17.0.1:8790	tcp://172.17.0.1:8790	16	2 %	155.1 MiB	6.6 GiB	2.3 %	0.0	58.7 MiB	96.3 MiB	0.0	22	287 B	1 KiB	0	48 KiB

127.0.0.1:8888/lab/tree/matmul.ipynb

File Edit View Run Kernel Tabs Settings Help

matmul.ipynb Python 3 (ipykernel)

```
[2]: import dask.array as da
from dask.distributed import Client, performance_report

def parallel_matrix_multiplication(n, client):
    # Create random matrices
    A = da.random.random((n, n), chunks=(n // 10, n // 10))
    B = da.random.random((n, n), chunks=(n // 10, n // 10))

    # Perform matrix multiplication
    result = da.matmul(A, B)

    # Compute the result with a performance report
    with performance_report(filename="dask_report.html"):
        result_computed = result.compute()

    return result_computed

if __name__ == "__main__":
    # Connect to your Dask cluster
    client = Client("tcp://scheduler:8786") # Replace with your scheduler address

    n = 2000 # Size of the matrix
    result = parallel_matrix_multiplication(n, client)
    print(f"Result of the matrix multiplication: \n{result}")

/opt/conda/lib/python3.10/site-packages/distributed/client.py:1392: VersionMismatchWarning: Mismatched versions found
+-----+-----+-----+
| Package | Client | Scheduler | Workers |
+-----+-----+-----+
| python | 3.10.13.final.0 | 3.10.12.final.0 | 3.10.12.final.0 |
| toolz | 0.12.1 | 0.12.0 | 0.12.0 |
+-----+-----+-----+
warnings.warn(version_module.VersionMismatchWarning(msg[0][“warning”]))
/opt/conda/lib/python3.10/site-packages/dask/array/routines.py:444: PerformanceWarning: Increasing number of chunks by factor of 10
    out = blockwise(
Result of the matrix multiplication:
[[497.91665321 506.06414028 496.76232952 ... 499.67452945 495.14205927
 505.43544277]
 [495.197243 493.83357347 489.99604106 ... 495.20964541 482.28825918
 498.33713429]]
```

dask_report.html

Summary Task Stream System Scheduler Logs Worker Profile (compute) Worker Profile (administrative) Scheduler Profile (administrative) Bandwidth (Workers) Bandwidth (Types)

Dask Performance Report

Select different tabs on the top for additional information

Duration: 5.97 s

Tasks Information

- number of tasks: 1600
- compute time: 65.32 s
- disk-read time: 266.02 ms
- disk-write time: 322.09 ms
- transfer time: 128.91 s

Scheduler Information

- Address: tcp://172.19.0.2:8786
- Workers: 8
- Threads: 128
- Memory: 4.00 GiB
- Dask Version: 2024.2.1
- Dask.Distributed Version: 2024.2.1

Calling Code

```
def parallel_matrix_multiplication(n, client):
    # Create random matrices
    A = da.random.random((n, n), chunks=(n // 10, n // 10))
    B = da.random.random((n, n), chunks=(n // 10, n // 10))

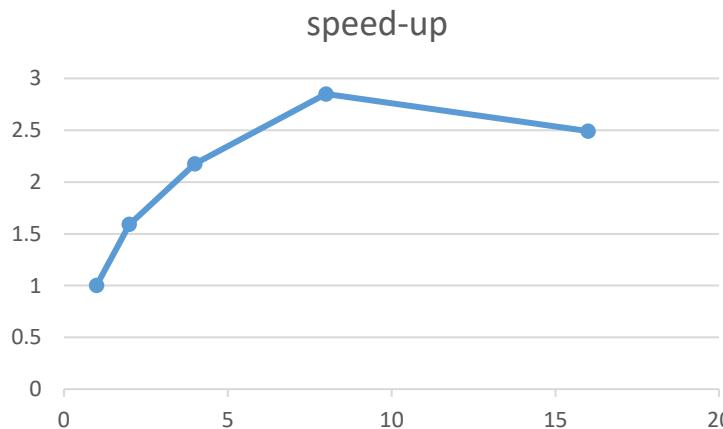
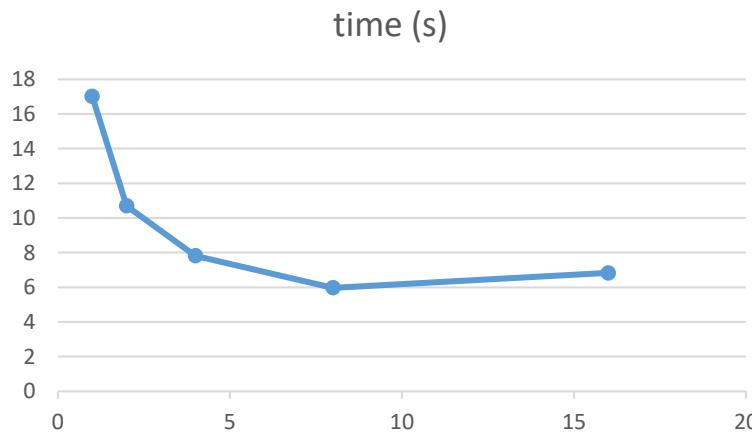
    # Perform matrix multiplication
    result = da.matmul(A, B)

    # Compute the result with a performance report
    with performance_report(filename="dask_report.html"):
        result_computed = result.compute()

    return result_computed
```

Speed-up

workers	time (s)	speed-up
1	17.01	1
2	10.7	1.58972
4	7.82	1.368286
8	5.97	1.309883
16	6.83	0.874085



Computing intensive : คำนวณมาก ๆ

Parallel Computing for Big Data

- Data-Parallel Model
 - send, receive, broadcast
ต้องรู้ชื่อและ type ของ จัดการในการถ่ายทอดข้อมูล
- Instead of explicit communication between nodes (as in Message-Passing Model), the focus is on partitioning large datasets and applying operations on these partitions in parallel. แทนที่จะมีการส่ง msg ระหว่าง dataset เป็นสองที่ ไม่ต้องแลกเปลี่ยน operation แต่dataset ทั้งหมดจะ
- The programming model works on top of distributed file systems and distributed data processing frameworks.
 - Ex: Apache Hadoop (with HDFS), Apache Spark (with RDDs)
- Key concepts:
 - Large datasets are partitioned and distributed across cluster nodes แปลง dataset และ กากจ่ายไปยัง cluster
 - Same operation or function is applied to each partition concurrently ที่ operation แบบขนานๆ
 - Framework handles distribution, communication, and coordination
 - ↓ ไม่ต้องสนใจการสื่อสาร รูปแบบการประมวลผล (implicit communication) กรณีทางเดียว
 - FW ต้องจัดการเรื่องการสื่อสาร

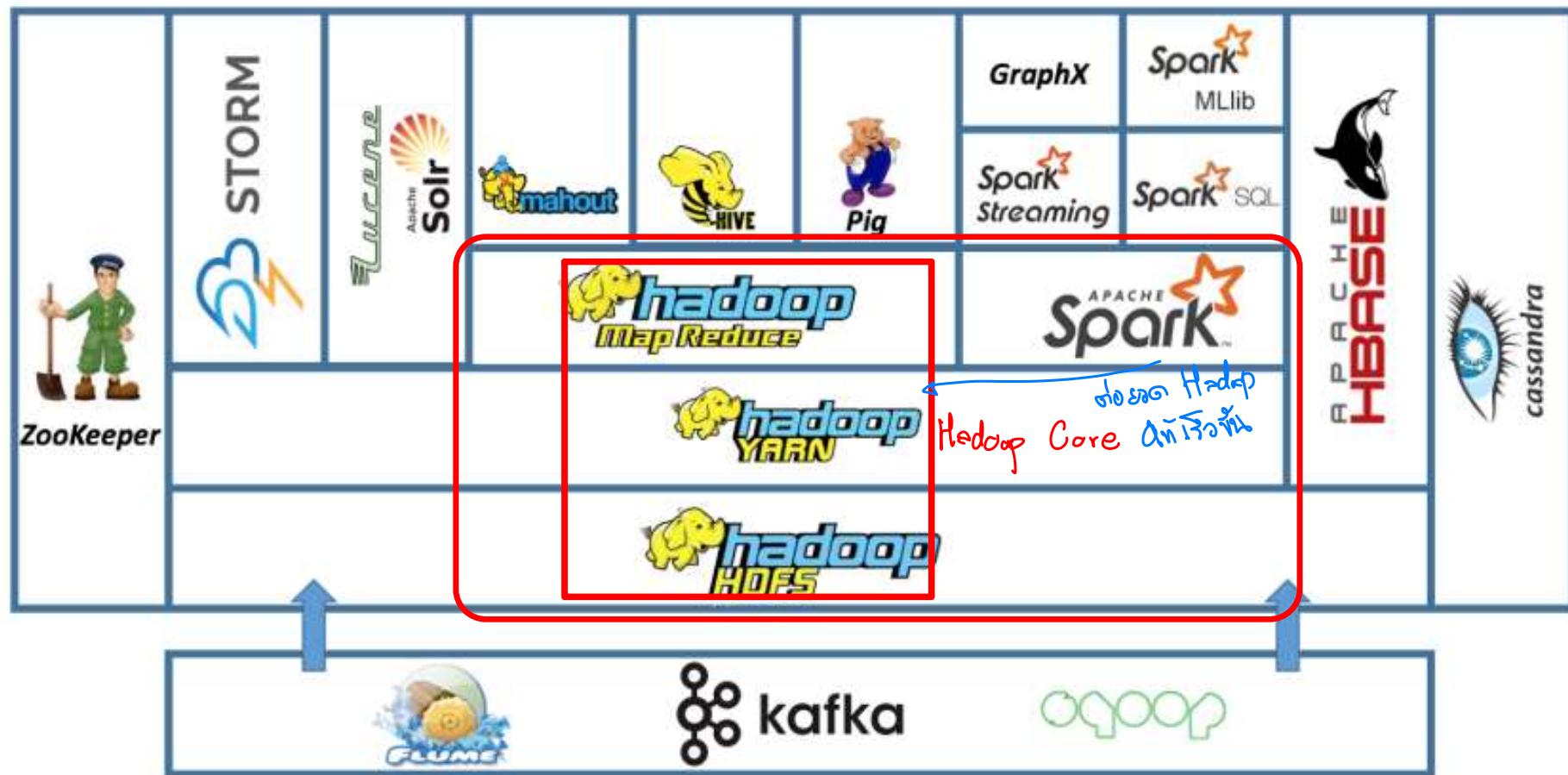
Apache Hadoop → Big Data Computing หลักการทำงาน

- An open-source framework for distributed storage and processing that led to Big Data era FN ក្នុងការងារនេះ distribute storage សមស្របតាមរយៈការពារ
- Key Features:
 - HADOOP កំណត់ឡើងខ្លួន ហើយអាមេរិក Parallel
 - Scalability: Easily scales out by adding more nodes to the cluster. គឺជាគិច្ចការងារថ្មី ដែល CPU, node
 - Fault Tolerance: Automatically handles failures by replicating data blocks. ផ្ទាល់អាមេរិក Hadoop ឲ្យការងារនេះ នូវការងារបានការពារឡើង
- Core Components:
 - Data storage → HDFS (Hadoop Distributed File System): Scalable, highly available storage system that splits large data files into blocks distributed across nodes. ទំនួលកំណត់ឡើង HDFS file ធ្លាប់បាននៅ block មួយ ឬក្នុងការងារបានការពារឡើង node
 - Resource Manager → YARN (Yet Another Resource Negotiator): Manages and schedules resources across clusters. កំណត់ឡើងការងារ គែងការងារបានការពារឡើង node Yarn
 - Programming → MapReduce: Programming model for processing large data sets with a parallel, distributed algorithm. Based on map and reduce operations.

ឯកសារ map, reduce function លើក្នុងការងារ

គឺជាការងារកំណត់ឡើង , លើក្នុងការងារ

Apache Hadoop Ecosystem



Hadoop Distributed File System (HDFS)

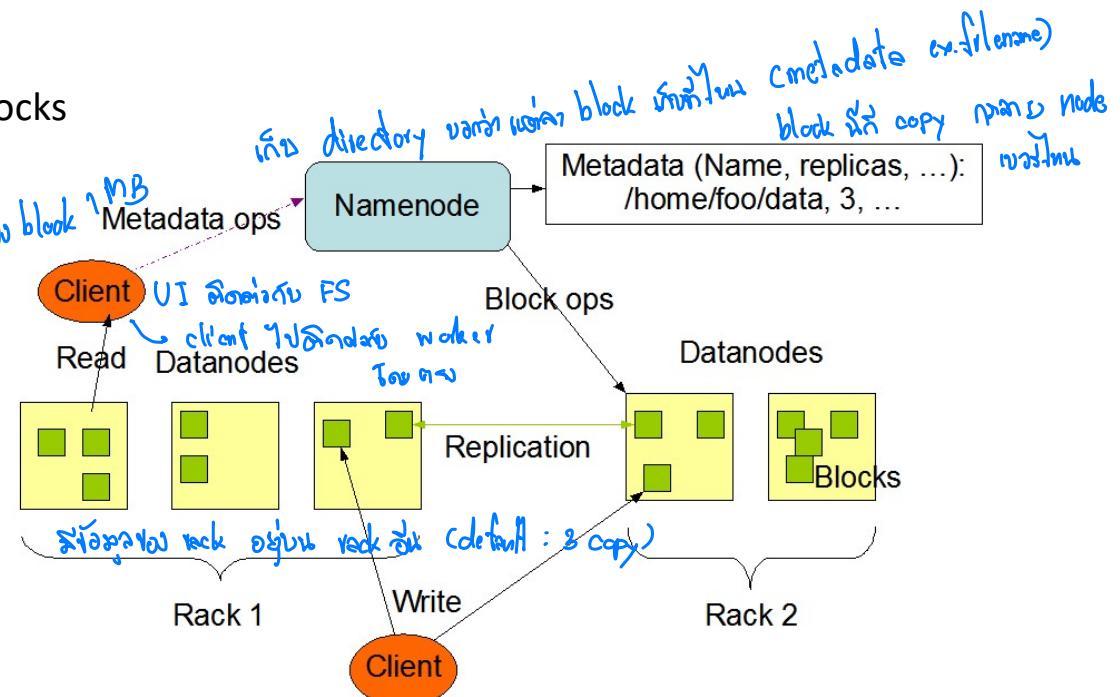
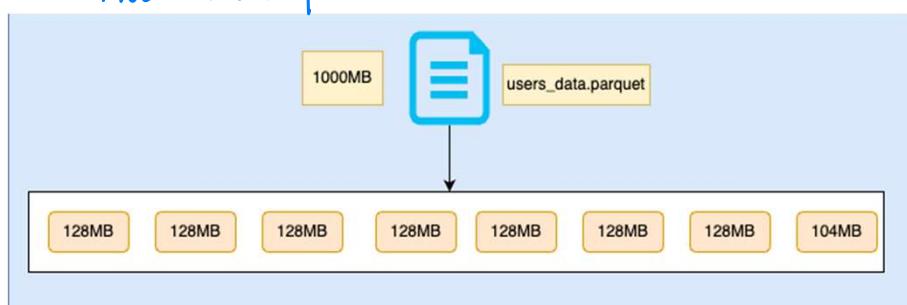
ສະ储age ឧបម node ទាំង នៅក្នុង Hadoop នឹងរឿងធម៌នដែលវារណា

- Distributed file system for storing large datasets across a cluster
- Support a traditional hierarchical file organization (directories and files)
- Run on commodity hardware អាជីវ អិន បូត្រូវ ពារិនិមាតុ ឬ Proprietary HW តើចេចរាបូចិន
- Focus on fault-tolerance and high throughput
- Optimized for
 - large files ទូទស័ព្ទ file ទូទស័ព្ទ តារាងនិងនិភ័យនិងនិរនោះ
 - write-once-read-many access pattern log data , archive ការគ្រប់គ្រងឯកសារ ឬសិទ្ធិការ
 - batch processing ពេន្យាកំណើនក្នុង real time

HDFS Architecture

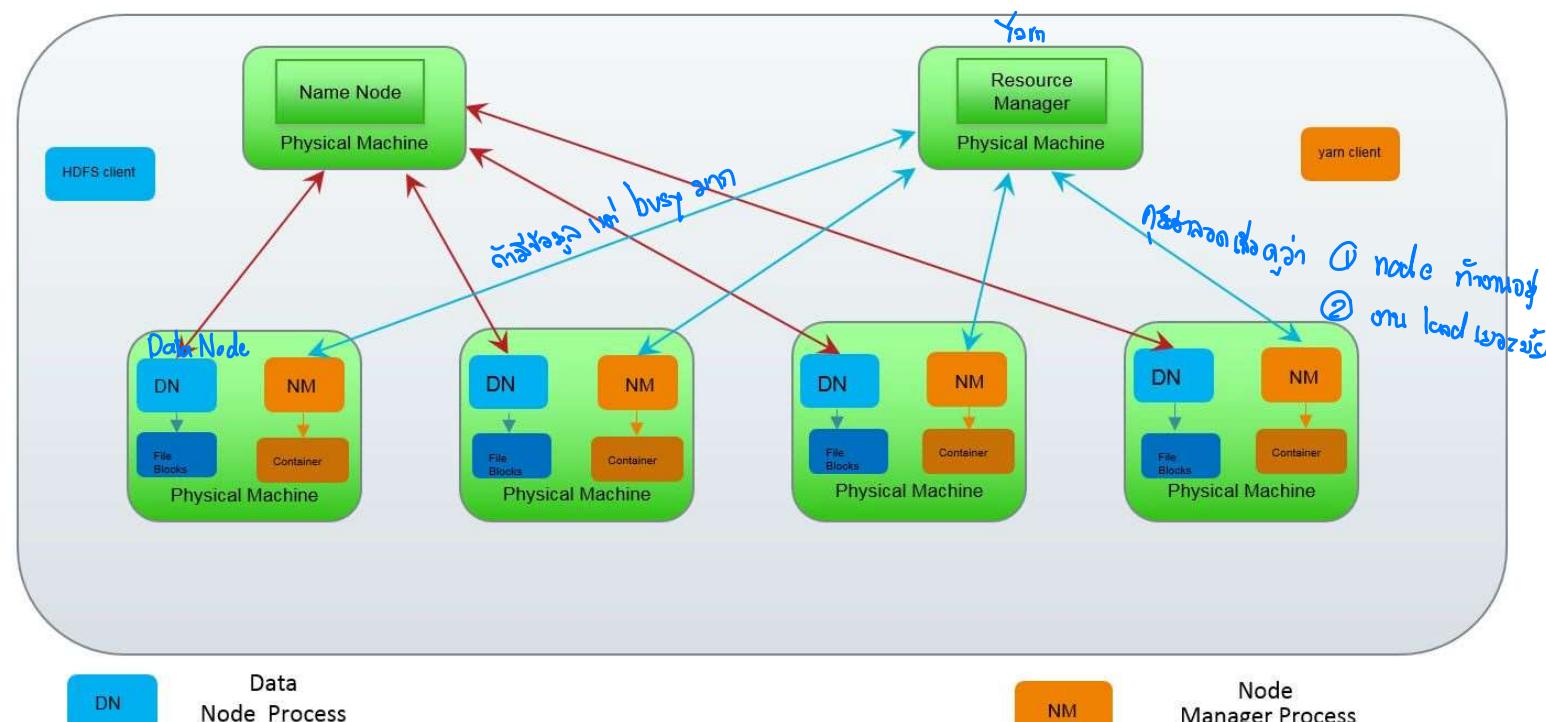
HDD ຈາງແມ່ນທີ່ມີເກີນອຸປະກອດ
ຕາວີ່ທີ່ມີກົດເຫັນກະຈາຍກຳໄປ

- A file is divided into blocks (default size = 128MB). ດັ່ງ file ເຊັ່ນເຂົ້າໃນ block 128 MB ສໍາມັນຢູ່ລົງ
- Each block is replicated and stored on multiple nodes. *ແທນໃຈ block ຕ່ອງສ້າງຫຼາຍ copy ເຊິ່ງການມານຸຍາ node ເພື່ອການກຳ fault tolerance ຂອງ copy ໃນ rack ເພື່ອກຳນົດການ communication*
- Files are write-once (except for appends and truncates)
- Components
 - file ຖື່ນທີ່ມີກົດເຫັນ ພາບສັນກະດົບທີ່ມີກົດເຫັນກຳຕົວທັນ*
 - ເຫັນວ່າມີກົດເຫັນຖຸກ copy ຈຶ່ງເຮັດ concept write-once*
- Components
 - HDFS Client
 - Name Node (master)
 - Manages file system metadata and coordinates accesses
 - Keeps track of block locations
 - Data Nodes (workers)
 - Store and serve the actual file data blocks



YARN (Yet Another Resource Negotiator) ඔහුගැනීමාපන Parallel

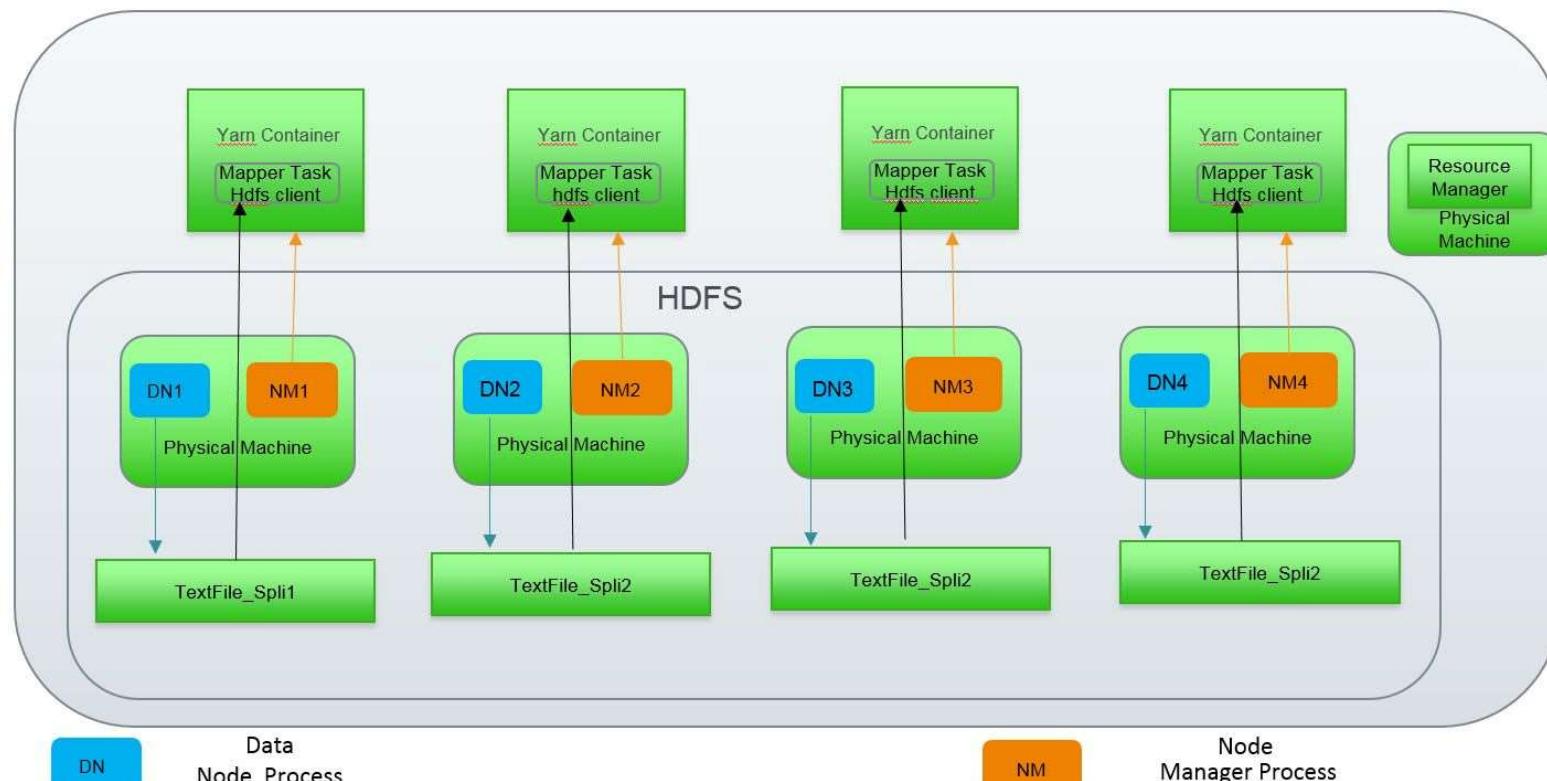
- Cluster resource management layer in Hadoop
- YARN schedules tasks across cluster nodes, considering resource availability, data locality for input splits, and efficient utilization of cluster resources.
- ResourceManager (RM):
 - Central process that manages and allocates cluster resources. Schedules jobs and allocates containers
- NodeManager (NM):
 - Runs on each node in the cluster and manages the resources on that node. Monitors the life-cycle of containers running on the node.



YARN Scheduling

YARN schedules MapReduce tasks (in this case, Mapper tasks) across containers running on different nodes.

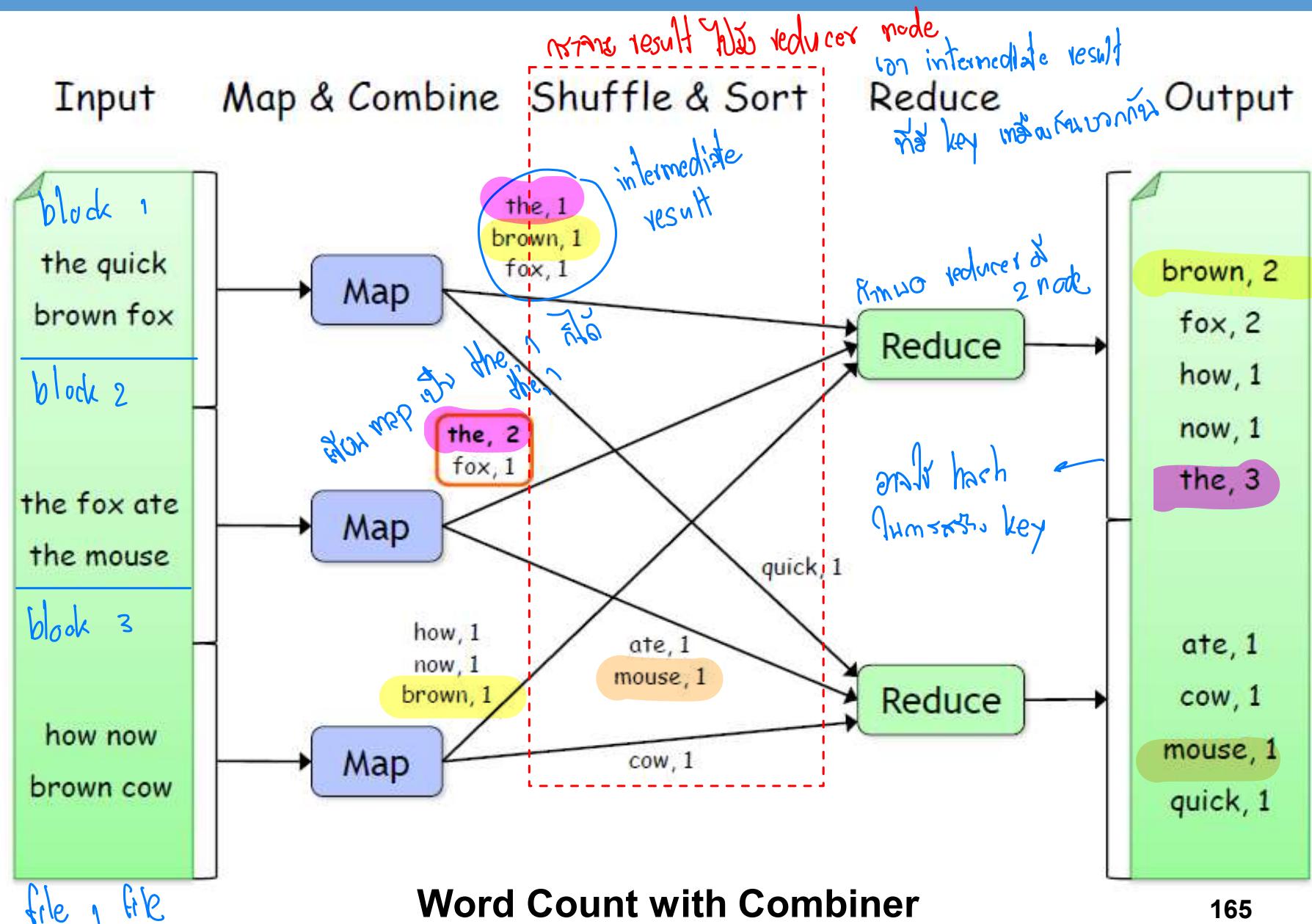
- ResourceManager allocates YARN containers to NodeManagers across the cluster nodes.
- Each container runs a process (Mapper Task) that reads input data from HDFS splits distributed across DataNodes.
- YARN tries to achieve data locality by scheduling Mapper Tasks on nodes where the input data resides, to minimize data movement across the network.
- If data locality is not possible, YARN schedules tasks on available nodes, and the Mapper Tasks read data from remote DataNodes. *ถ้าไม่สามารถ assign บนหอดีต้องส่งไป assign บนหอดูซึ่งมี*



MapReduce Programming Model

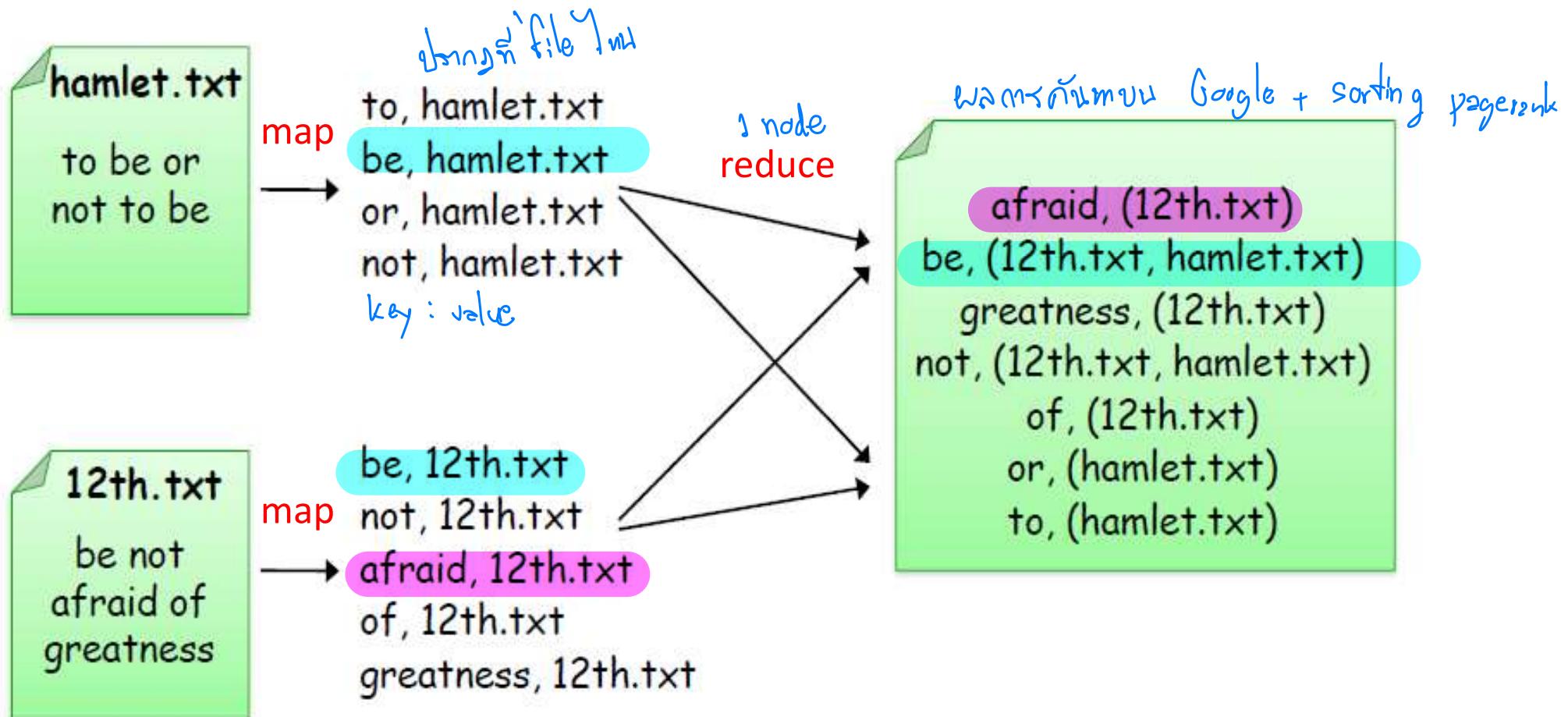
- A framework for processing large datasets in a parallel and distributed manner across a cluster. *ក្នុងអេក្រង់នេះមាន Map Reduce រាយការ។*
- It consists of two main phases: the Map phase and the Reduce phase.
- **Map Phase:**
 - The input data is split into multiple chunks (key-value pairs) and distributed across the cluster nodes. *និងចូលរាយការឡើង នានា key : value*
 - The *Map function* is applied to each key-value pair, producing intermediate key-value pairs as output. *ឯកសារតិច intermediate គឺជាការសំណួរ*
 - The intermediate key-value pairs are sorted and partitioned based on their keys.
- **Shuffle and Sort:** *ដំឡើង map / reduce ដែលធ្វើឡើងឡើង*
 - The intermediate key-value pairs are redistributed across the cluster nodes based on their keys, so that all values with the same key are sent to the same node.
 - The values with the same key are sorted and prepared for the Reduce phase.
- **Reduce Phase:** *ទាក់ទងយ៉ាងនៃ map*
 - The *Reduce function* is applied to each unique key and its corresponding values (the sorted values from the Map phase). *ឯកសារពីរការស្ថិតិយោគ*
 - The Reduce function performs operations (such as aggregation, filtering, or transformation) on the values and produces the final output key-value pairs.
 - The final output is written to the distributed file system (e.g., HDFS in Hadoop).
ឯកសារតិចនៅក្នុង ឯកសារពីរការស្ថិតិយោគ HDFS

Ex: word count with MapReduce



Keyword ជូនក្នុងឯកសារ, រាយការ

Ex. Inverted Index with MapReduce



Google សិទ្ធិ key : url

be . www . facebook . com

Map-Reduce Programming Model

- **Input data type:** file of K/V records
- **Map function:** $(K_{in}, V_{in}) \rightarrow \text{list}(K_{inter}, V_{inter})$
- **Reduce function:** $(K_{inter}, \text{list}(V_{inter})) \rightarrow \text{list}(K_{out}, V_{out})$
- Example:

```
def mapper(line):
```

```
    foreach word in line.split();
```

```
        output(word, 1);
```

$\{(the, 1)$
 $(the, 1)$
 $(the, 1)$
 $(be, 1)\}$

$// (K_{inter}, V_{inter})$

```
def reducer(key, values):
```

```
    output(key, sum(values));
```

$\{(the, 3)$
 $(be, 1)\}$

$// \text{list}(K_{out}, V_{out})$

Word Count with MapReduce in Python

```
from mrjob.job import MRJob
from mrjob.step import MRStep

class WordCount(MRJob):
    def mapper(self, _, line):
        """
        The mapper function splits the input line into words and emits (word, 1) pairs.
        """
        words = line.split()
        for word in words:
            yield word, 1

    def combiner(self, word, counts): fn mssng(The,1)(The,1)往ugha (The,2)
        """
        The combiner function performs a local aggregation of counts for each word on the mapper node.
        """
        yield word, sum(counts)

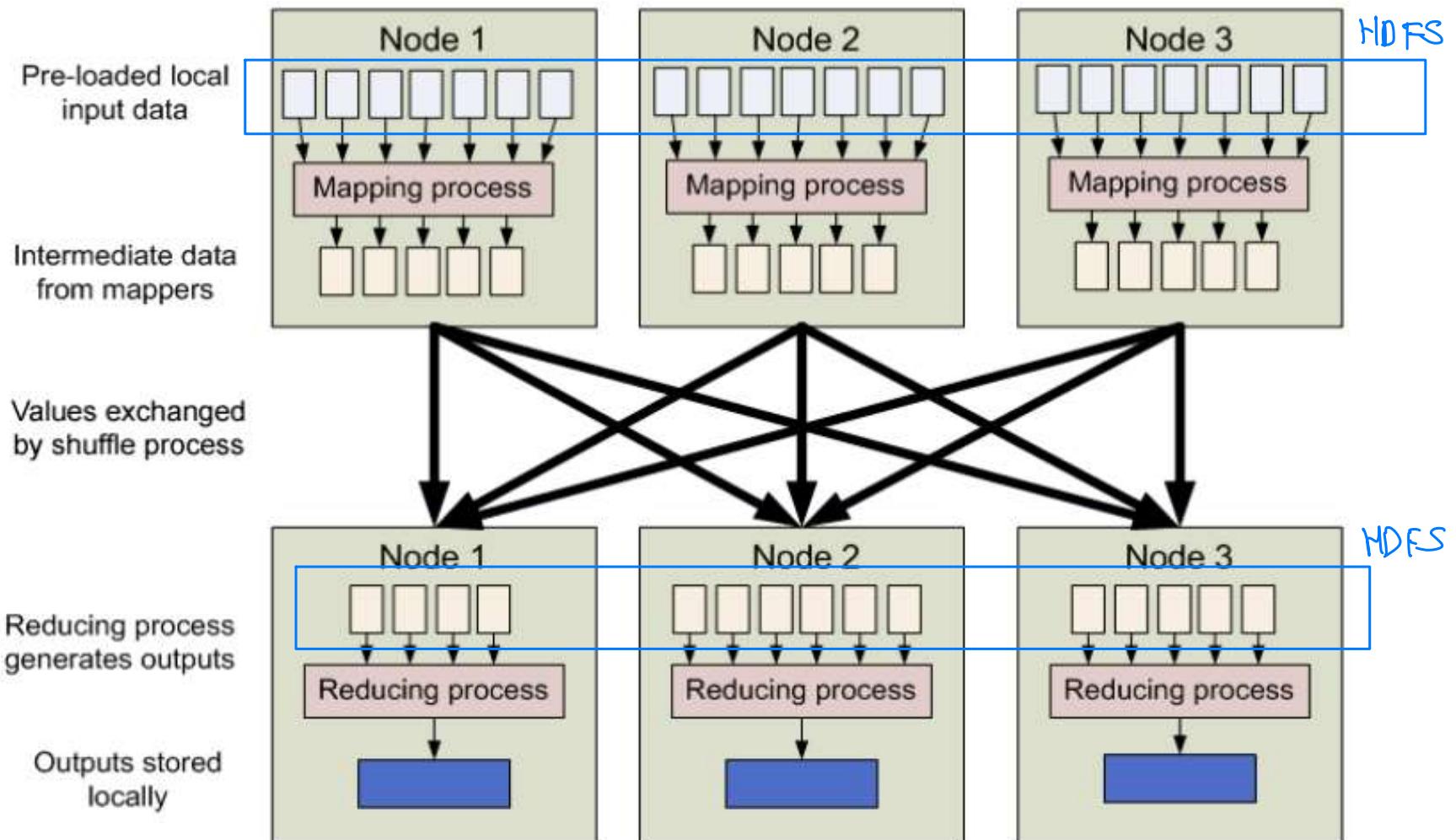
    def reducer(self, word, counts):
        """
        The reducer function aggregates the counts from all mapper nodes and emits the final (word, count) pairs.
        """
        yield word, sum(counts)

    def steps(self):
        """
        The steps method defines the MapReduce job workflow.
        """
        return [
            MRStep(mapper=self.mapper,
                   combiner=self.combiner,
                   reducer=self.reducer)
        ]

if __name__ == '__main__':
    WordCount.run()
```

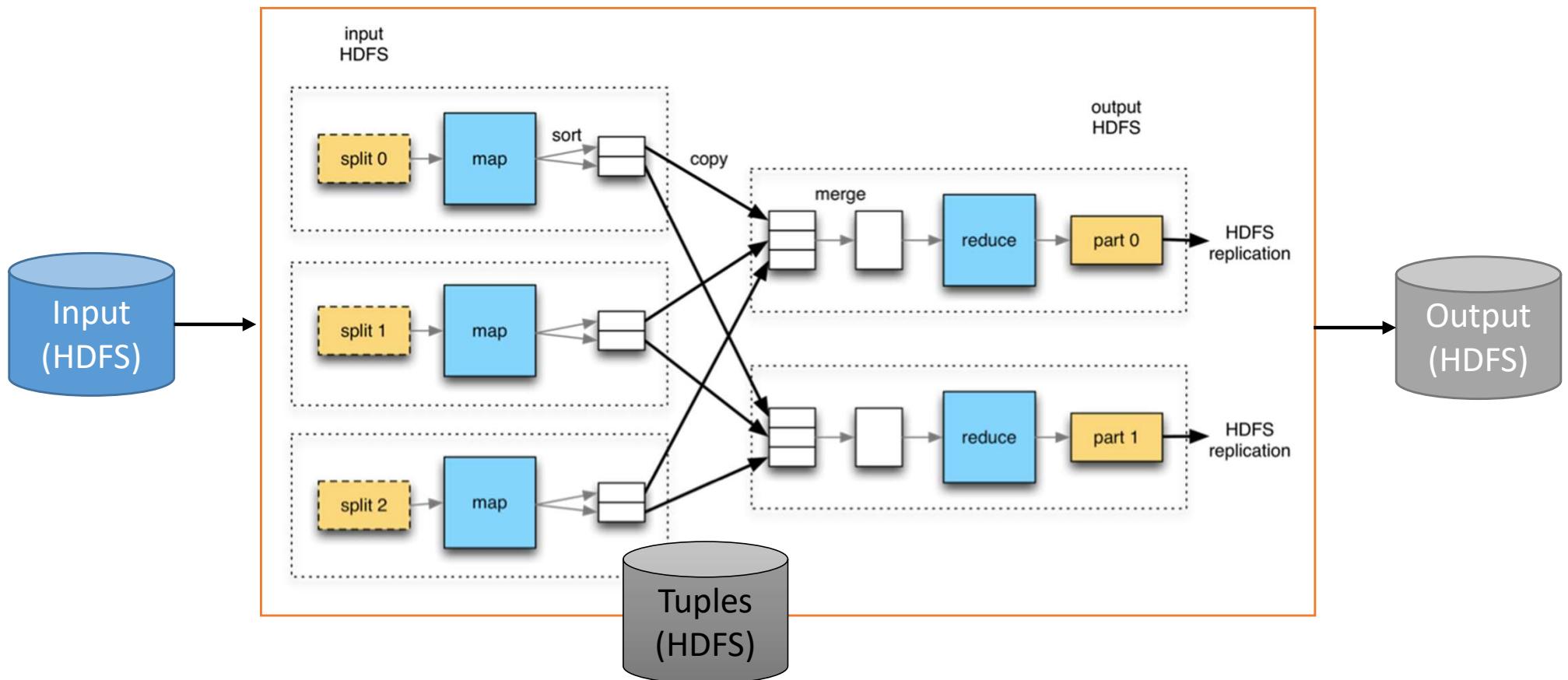
```
python word_count.py -r hadoop hdfs://path/to/your/input_file.txt --output-dir hdfs://path/to/output_directory/
```

Map-Reduce Processing



Hadoop Disk-oriented Processing

input, output, result ማህተም ዘመን



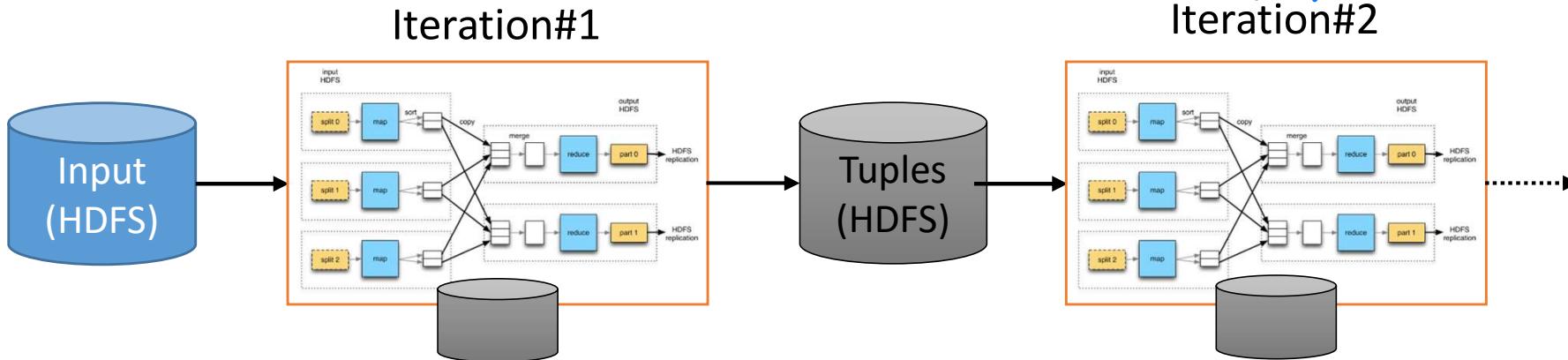
Hadoop Problems

- HDFS is **slow** on write (due to replication) *ກົດໝາຍ copy ທົ່ວລະ ວິທີ*
- Not suitable for iterative processing *ຖຸກຮອບຈັດເປັນພາຍລັງລົງ ມີ HDFS*
- Many problems aren't easily described as map-reduce

ໄຟເຫດກົດ map reduce ຖຸກສານ ເພີ້ມ dot product

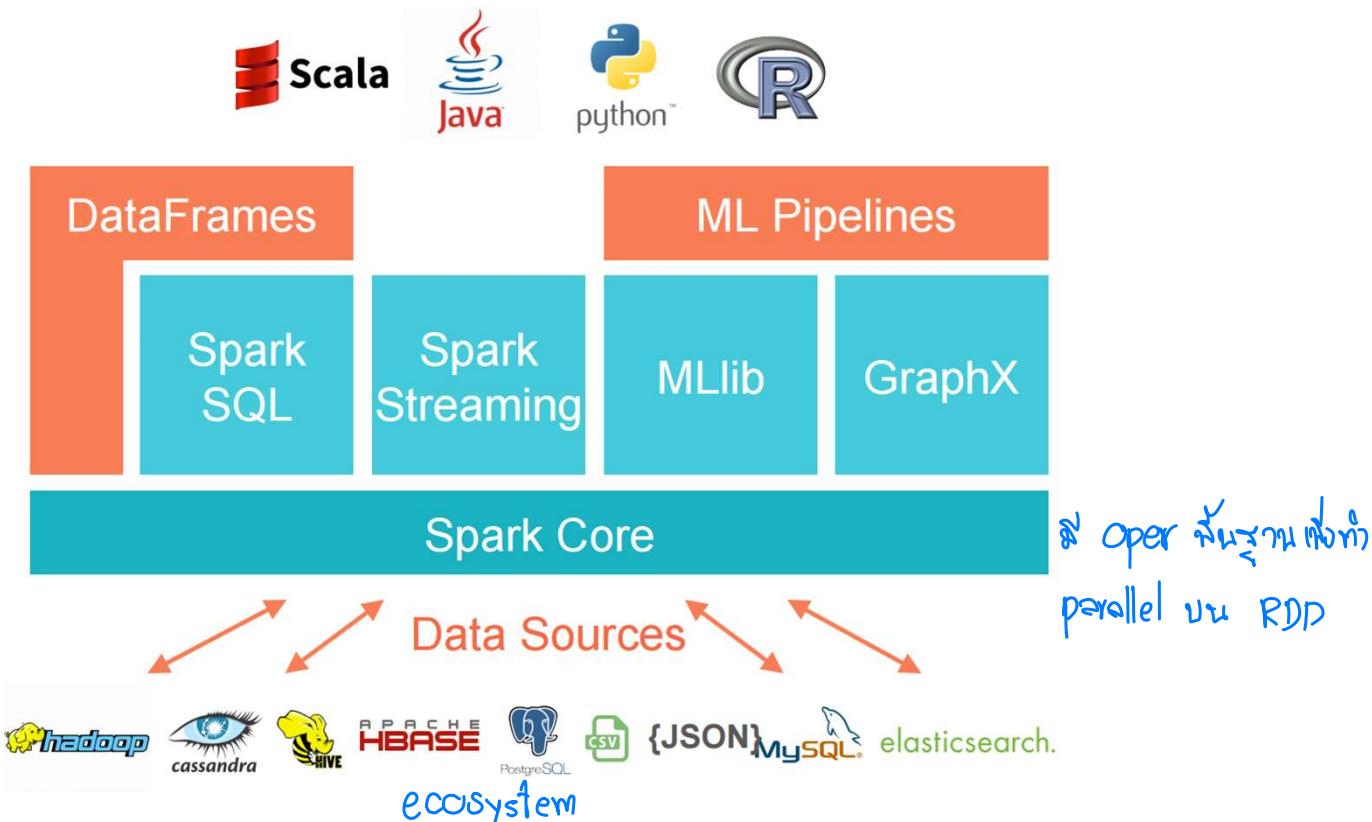
*ກົດຄູນ matrix key ດັ່ງ row number
value ດັ່ງລາຍລະອຽດ*

Iteration#2



Apache Spark

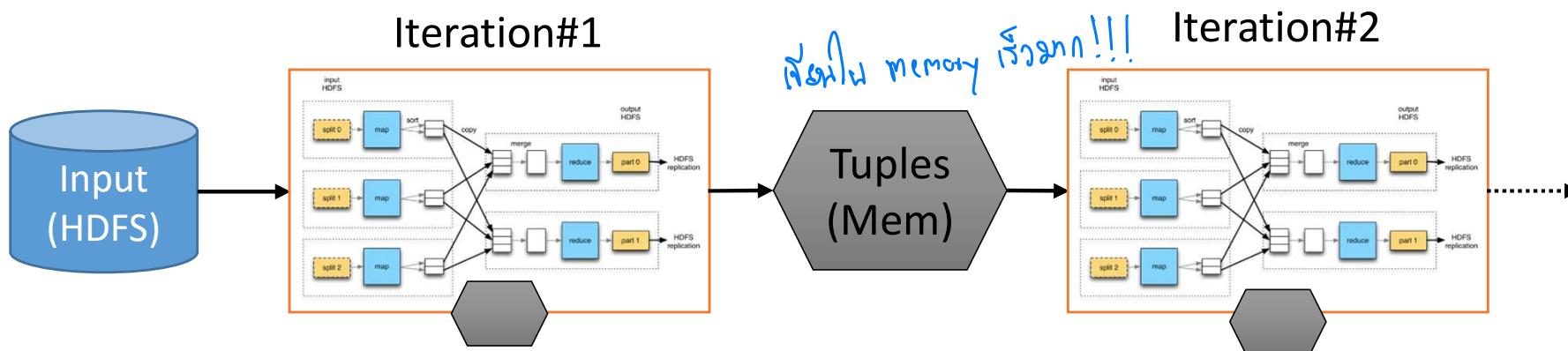
- In-memory data processing *ການໂປ່ງ memory ຖືມຂອງມາດລະບານ HDD*
- Resilient Distributed Data (RDD) *in memory ອັນຈະເຕັມຂອງມາດລະບານໄດ້ຮັບຮັດ*
- Large set of operations (transformations & actions) *ກຳນົດສິ່ງ operation ມາດພົນກວດຈາກ mapReduce*



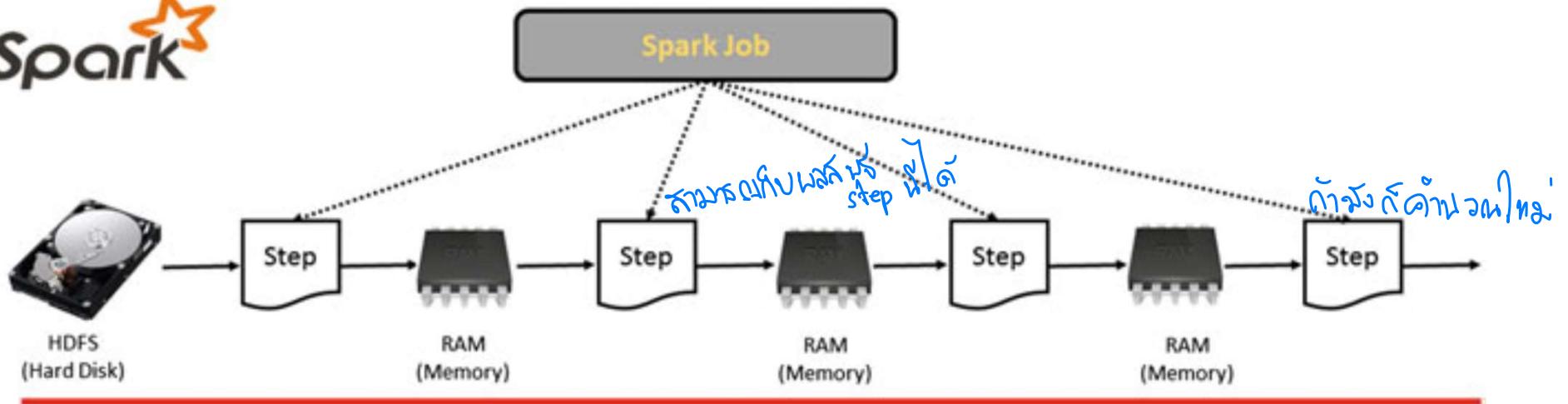
Source: M. Zaharia, "New Directions for Spark in 2015", Spark Summit East 2015, 18 March 2015.

Spark Memory-oriented Processing

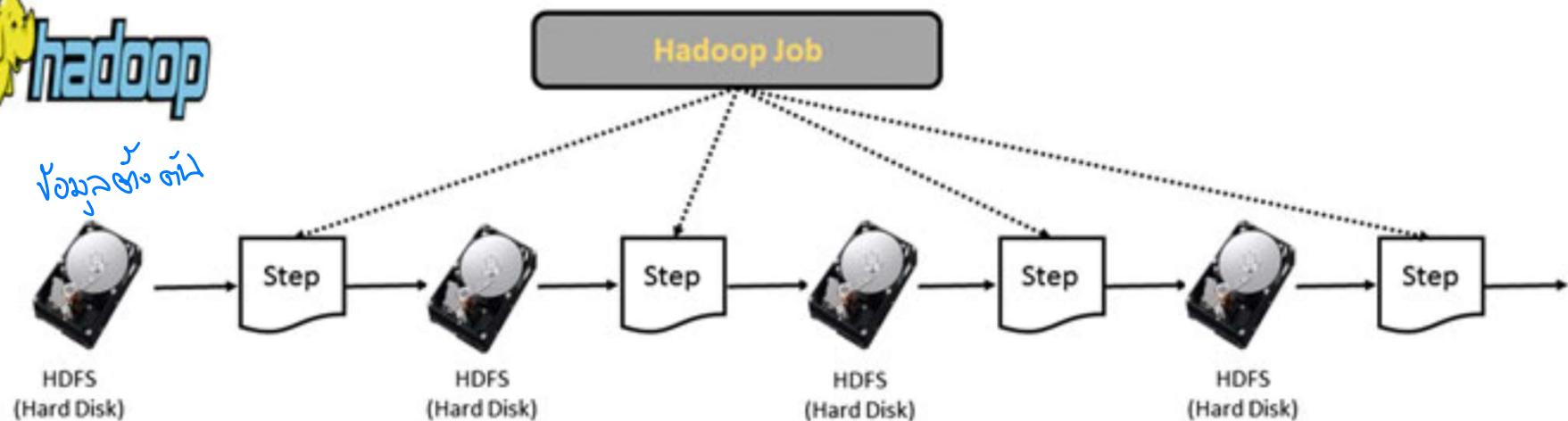
- Intermediate data are kept in memory. *ສ່າງຮອກຳນົງການ iterative ຢໍາ*
- Only read from disk for input and write to disk for output *input, output*
ໃຊ້ຈາກ HDFS ຢໍາ



Spark vs Hadoop



ข้อมูลต่อ



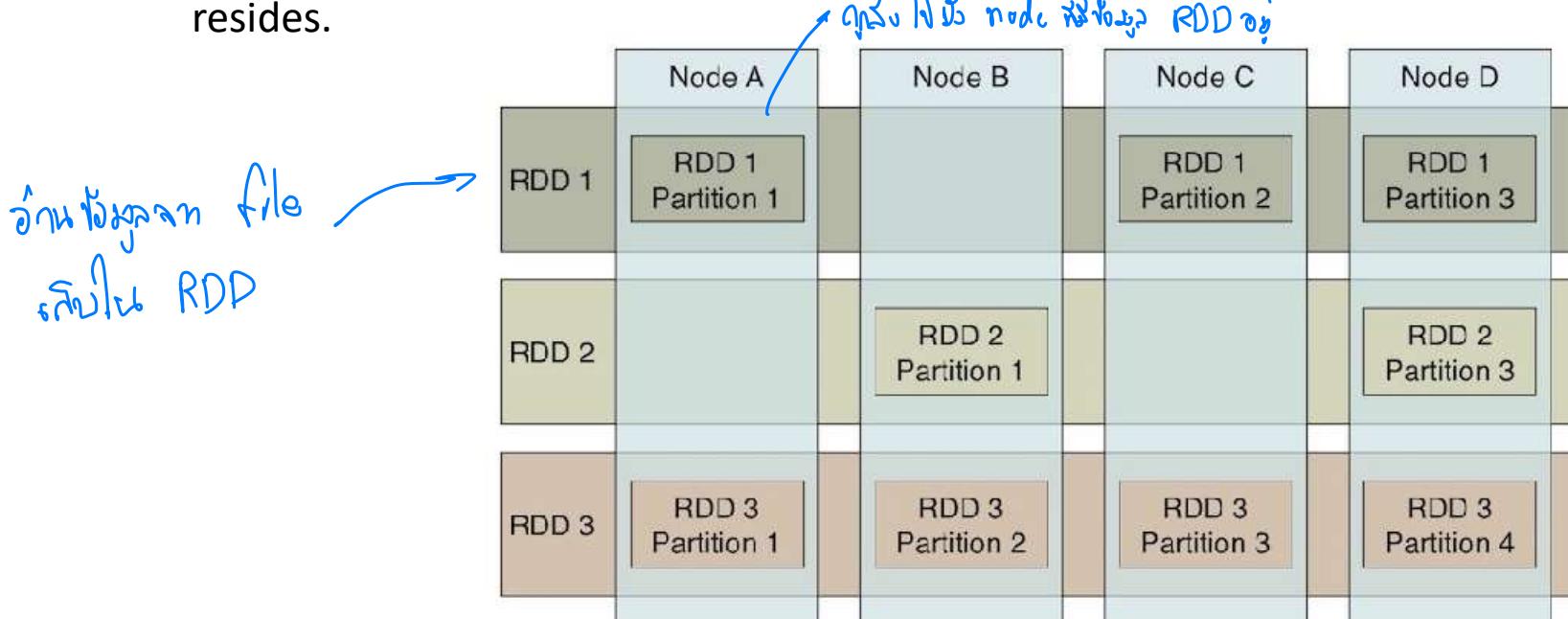
Why Spark is Faster? --> Data Abstraction

គេចងក្រងការបង្កើតផ្ទុល

- Hadoop's basic data unit is a block in HDFS *(គេង block នៃ HDFS)*
 - Slowwww on write → must completely replicated before continue
 - Fault tolerance by prevention (block replication) *(វិធានការរំលែកក្នុង (replication))*
- Spark's basic data unit is RDD
 - Array of “read-only” data resided in memory *(តារាងពិសេស នៅលើក្នុង memory នៃ RDD)*
 - Fault tolerance by recovery (recompute RDDs) *(វិធានរក្សាទុក្ខ (តារាងពិសេស នៃ RDD))*
- Why must RDD be read-only (immutable)?
 - In a cluster, multiple nodes may operate on the data in parallel. If the data were mutable, nodes would need to synchronize to each other to ensure consistency, which can be a significant overhead.
*(memory & cache នៅលើ multiple memory នឹងមាន inconsistency
នៅក្នុង HDFS ដែលមានសម្រាប់ការប្រើប្រាស់)*

RDD: Resilient Distributed Data

- **Immutable** Distributed Collection of Objects
- RDD can be created by parallelizing an existing collection (data structure) or loading a dataset from a storage.
- RDD partitioning is based on the specified number of partitions, cluster configuration, or storage partitioning.
- Partitions are distributed across different nodes in the Spark cluster.
- The Spark driver sends tasks to executors. Each task processes data in one partition.
- Spark tries to place partitions on nodes in a way that respects data locality preferences. If data is read from HDFS, Spark will attempt to perform computations on the nodes where the data resides.



Creating RDDs

sparkContext

list, DataFrame

- `sc.parallelize(data)` – create an RDD from data

RDD

```
lines = sc.parallelize(List("a", "an", "the"))
```

តម្លៃកុង RDD នៃ file

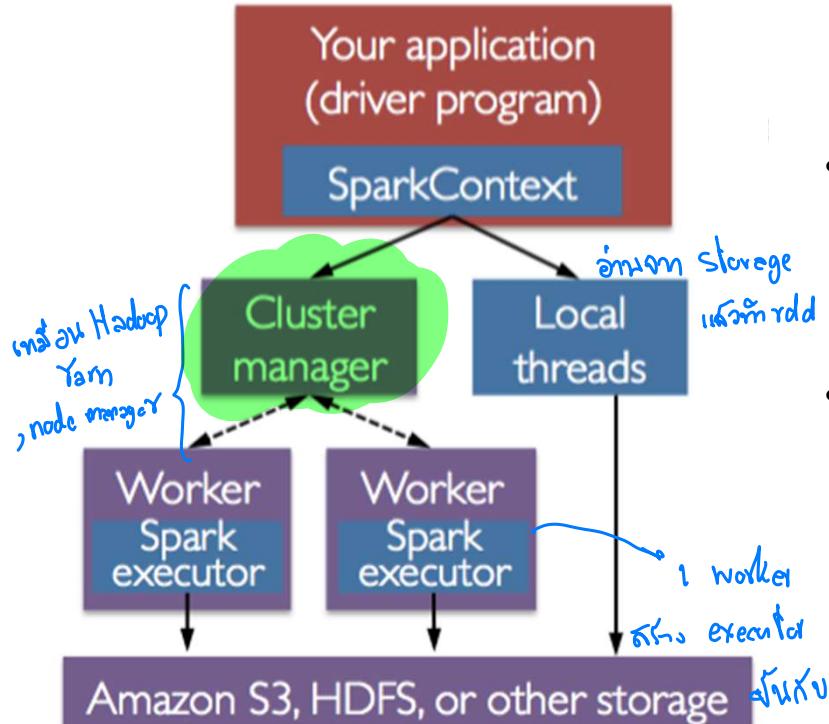
- `sc.textFile(filename)` – create an RDD from
 - HDFS, text files, Hypertable, Amazon S3, Apache Hbase, SequenceFiles, any other Hadoop InputFormat, and directory
 - Can be wildcard: /data/file*

ការពាក្យណា data ដែលត្រួតពីរឿងទិន្នន័យ file

```
rdd = sc.textFile("/home/user01/data/file*)")
```

```
rdd.count() បែងចាយរាយការណ៍
```

Spark components



- Driver program → Python, R
 - The user program that creates `SparkContext` and communicates with it to submit jobs and receive the results.
- `SparkContext` process
 - Entry point to Spark functionality.
 - Coordinates execution, communicates with cluster manager for resource allocation and task scheduling.
- Worker program
 - Executes tasks assigned by the driver program.
 - Run on cluster nodes or local threads (cluster mode vs local mode)
 - If you run a Spark program from a Jupyter notebook on your laptop, it will run in local mode. This means Spark driver, executor, and all tasks will run on the same machine as Jupyter notebook.
 - If you want to run the program in cluster mode, you need to configure Spark to connect to a remote Spark cluster and submit your program to the cluster using a cluster manager like YARN or Spark Master.
- Cluster Manager allocates and manage resources across the cluster
 - E.g. YARN, Mesos, Kubernetes, Spark Master (in standalone mode)
- Executor
 - A process that runs on a worker node, responsible for executing a task.
 - A worker program can host multiple executors.
 - Each executor can have multiple cores and run multiple tasks in parallel.

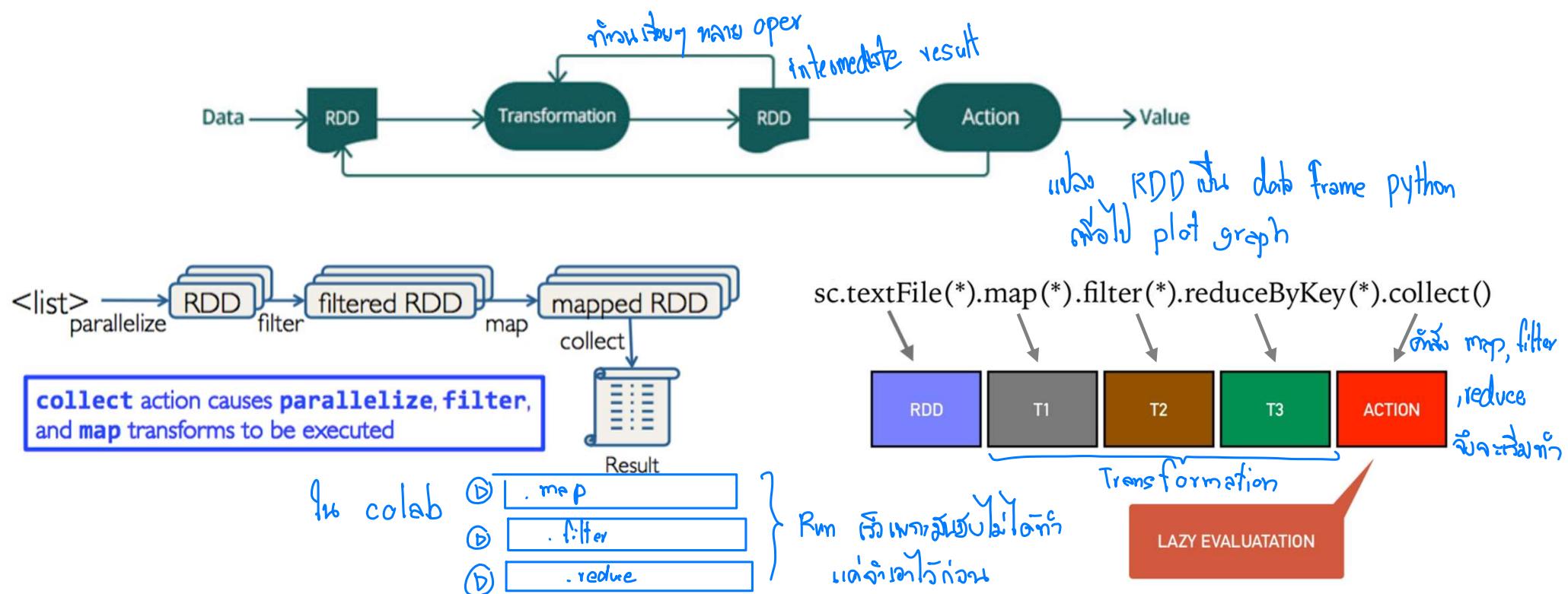
Why Spark is Faster? --> Execution Planning

- Hadoop is very straight-forward, ^{ก็ทำงานไปเรื่อยๆ ไม่ต้องล็อป} Map-then-Reduce, executed immediately
- Spark is **lazy-execution**, but constructs DAG (Directed Acyclic Graph) of RDDs (Sequence of tasks to be executed on RDDs)
 - Schedule DAG when needed

^{กราฟที่ execute บน RDD}
lazy - execution คือ apply function ยังไม่ทําหน้าที่ ทำจริงๆ เมื่อต้องการลําเลงร

Typical Spark Program

- Create RDD from external dataset (e.g. file from HDFS) or existing collection of objects (List, Set, DataFrame)
- Transformation (Create a new RDD from existing RDDs)
- Action (Compute result from an RDD) ការងារ RDD កំណត់
- (Maybe) persist (cache) RDDs to disk or memory RDD ផ្គាល់នូវ disk និងផ្គាល់នូវ data បានឡើង
- *Transformations are lazy* (not computed immediately) ផលិតផលនេះ មិនបានត្រួតពិនិត្យឡើង
- Transformed RDD is executed when an *action* runs on it



Sample Spark Transformations

- **map(func)**: Return a new distributed dataset formed by passing each element of the source through a function func.
- **filter(func)**: Return a new dataset formed by selecting those elements of the source on which func returns true
- **union(otherDataset)**: Return a new dataset that contains the union of the elements in the source dataset and the argument.
- **intersection(otherDataset)**: Return a new RDD that contains the intersection of elements in the source dataset and the argument.
- **distinct([numTasks]))**: Return a new dataset that contains the distinct elements of the source dataset
- **join(otherDataset, [numTasks])**: When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key. Outer joins are supported through leftOuterJoin, rightOuterJoin, and fullOuterJoin.

Sample Spark Actions

இலு execution transformation நிலை

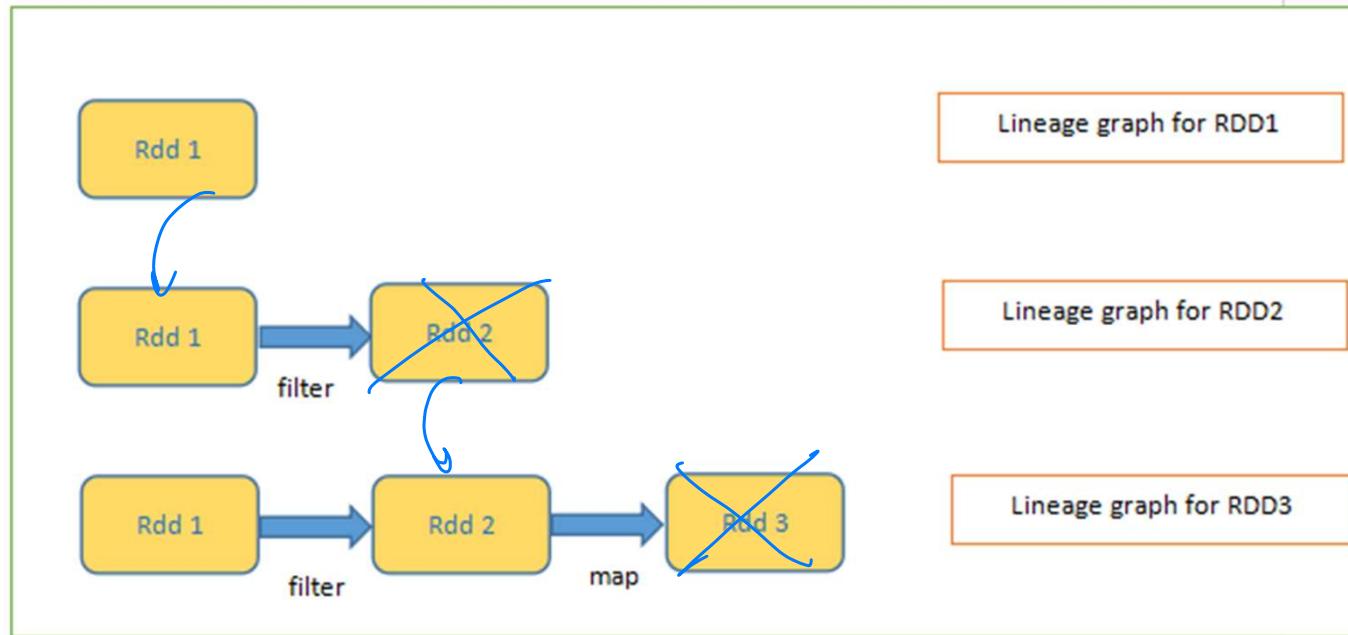
- **reduce(func)**: Aggregate the elements of the dataset using a function func (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
- **collect()**: Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data. தாபகமாகவுள்ளது
- **count()**: Return the number of elements in the dataset.

Actions cause calculations to be performed;
transformations just set things up (lazy evaluation)

RDD Fault Tolerance

- RDD can be kept in memory or persisted to disks
- If a node fails, only the partitions on that node need to be recomputed, as Spark keeps track of the lineage of transformations of each RDD.

ถ้าต้อง recompute → spark จะทำการ record ใน rdd's lineage graph
ที่มี 3 ประเภทคือ filter, filter + map



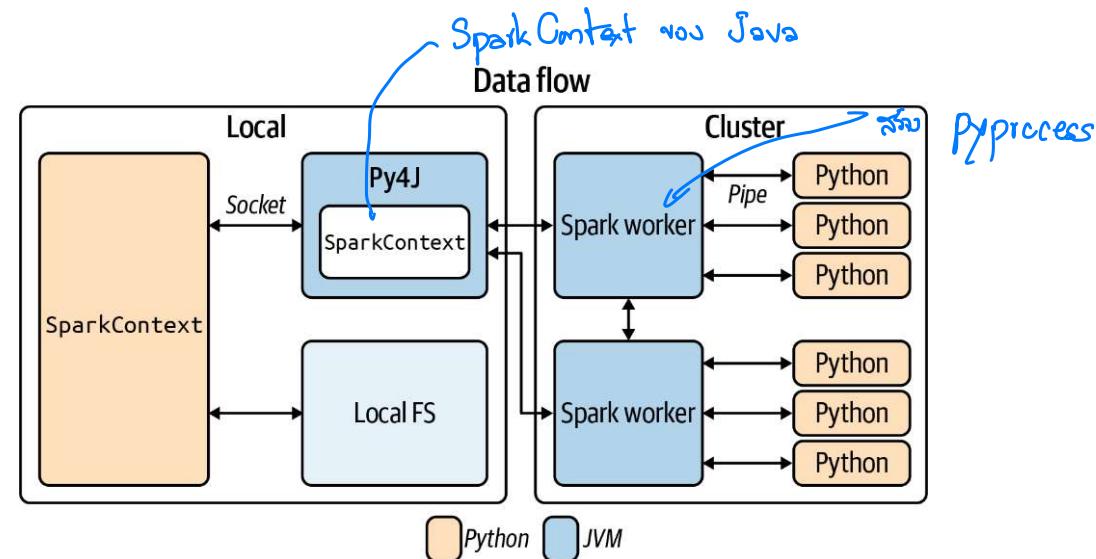
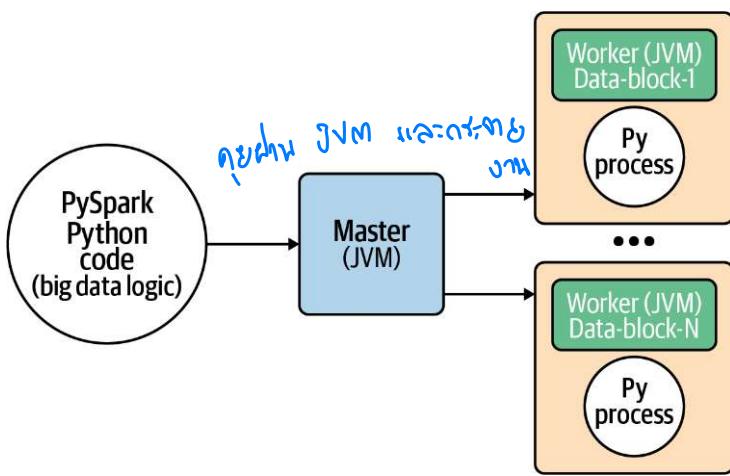
Common Communication Patterns in Spark Program

RDD operations can cause implicit communication. For example:

- Shuffling: Redistribution of data across partitions for operations like reduceByKey, groupByKey, join, sortBy, and repartition, causing significant network and disk I/O. RDD គ្មាន node ទាំង
- Broadcasting: Efficiently sending a large, read-only variable to all nodes, commonly used in join operations with disparate dataset sizes. ការបាយចែងតាមការដែលត្រូវបានផ្តល់នៅលម្អិត node ទាំងអស់ នៃការប្រើប្រាស់ Join ការវាយតម្លៃ dataset នៃក្នុងមួយក្រុង node
- All-to-All: Each node communicates with all others, often triggered by repartition.
- Map-Side: Operations like *map* and *filter* apply functions to RDD elements independently, with no cross-node communication. apply map function
- Reduce-Side: Aggregations like *reduce* and *fold* bring data together for processing, requiring data movement for combiners and reducers.
- Gather to Driver: The *collect* operation pulls data from all worker nodes back to the driver program. ទៅលើ collect តើក្នុងក្រុងការប្រើប្រាស់ driver program

PySpark Architecture

- PySpark is a Python API for Apache Spark, built on top of Spark's Java API.
- PySpark uses the Py4J library, which acts as a bridge between the Python driver program and the Spark Java core.
- RDD operations in user's Python program are sent to Spark workers, which create Python processes to execute the tasks.



<https://spark.apache.org/docs/latest/api/python/index.html>

<https://www.oreilly.com/library/view/data-algorithms-with/9781492082378/ch01.html>

PySpark Examples

USERS.CSV (MOVIELENS DATABASE)

```
id | age | gender | occupation | zip  
1|24 M technician|85711  
2|53 F other|94043  
3|23 M writer|32067  
4|24 M technician|43537  
5|33 F other|15213  
6|42 M executive|98101  
7|57 M administrator|91344  
8|36 M administrator|05201
```

M = 670
F = 273

Spark df \neq python df

With RDD and MapReduce

```
from pyspark import SparkContext  
  
sc = SparkContext.getOrCreate()  
  
print sc.textFile( "users.csv" ) \  
.map( lambda x: (x.split("|")[2], 1) ) \  
.reduceByKey(lambda x,y:x+y).collect()  
  
sc.stop()
```

M, 1
M, 1
F, 1
M, 1

[(u'M', 670), (u'F', 273)]

With Spark DataFrame and built-in functions

```
from pyspark import SparkContext  
from pyspark.sql import SparkSession  
sc = SparkContext.getOrCreate()  
spark = SparkSession(sc)  
  
spark.read.load( "users.csv", format="csv", sep="|" ) \  
.toDF( "id","age","gender","occupation","zip" ) \  
.groupby( "gender" ).count().show()  
  
sc.stop()
```

Word Count with PySpark

```
from pyspark import SparkConf, SparkContext

# create a SparkConf object and set the application name
conf = SparkConf().setAppName("WordCount")

# create a SparkContext object
sc = SparkContext.getOrCreate(conf)

# load the input text file as an RDD
input_file = sc.textFile("input.txt")

# split the text into individual words transform
words = input_file.flatMap(lambda line: line.split())

# count the occurrence of each word
word_count = words.map(lambda word: (word,
1)).reduceByKey(lambda a, b: a + b)

# print the word count result
for word, count in word_count.collect():
    print("{}: {}".format(word, count))

# stop the SparkContext
sc.stop()
```

- Load the input text file as an RDD using the `textFile()`
- Split the text into individual words using `flatMap()`
- Create a new RDD of word-count pairs using `map()`
- Group the word-count pairs by word and sum the counts using `reduceByKey()`
- Collect word count result using `collect()` and print

Dot Product with PySpark

```
from pyspark import SparkConf, SparkContext
import findspark
import random
import datetime

# findspark sets environment variables for pyspark path
findspark.init()

# Create a SparkConf and SparkContext
conf = SparkConf().setAppName("DotProduct")
sc = SparkContext.getOrCreate(conf)

n = 10
# Define two vectors as RDDs
rdd1 = sc.parallelize([random.random() for i in range(n)])
rdd2 = sc.parallelize([random.random() for i in range(n)])

start = datetime.datetime.now()

# Compute dot product using zip, map, and reduce operations
dot_product = rdd1.zip(rdd2).map(lambda x:
x[0]*x[1]).reduce(lambda x, y: x+y)

end = datetime.datetime.now()
print(end - start, "h:mm:ss")

# Print the dot product
print("The dot product is:", dot_product)

# Stop the SparkContext
sc.stop()
```

- Two vectors are defined as RDDs using the parallelize().
- Dot product is computed by:
 - zip() pairs up corresponding elements of the two vectors
 - Then, map() multiplies the paired elements together
 - Finally, reduce() sums the resulting products.
- Try increasing n to 1 million and more
- Notice the execution time and overhead

0:00:09.607832 h:mm:ss
The dot product is: 2.7423063387967446

Dot Product with PySpark

```
from pyspark import SparkConf, SparkContext
import findspark
import random
import datetime

findspark.init()

# Create a SparkConf and SparkContext
conf = SparkConf().setAppName("DotProduct")\
    .set("spark.driver.memory", "16g") \
    .set("spark.executor.memory", "16g")
sc = SparkContext.getOrCreate(conf)

n = 100000000
# Define two vectors as RDDs
rdd1 = sc.parallelize([random.random() for i in range(n)])
rdd2 = sc.parallelize([random.random() for i in range(n)])

start = datetime.datetime.now()

# Compute the dot product using zip, map, and reduce
operations
dot_product = rdd1.zip(rdd2).map(lambda x:
x[0]*x[1]).reduce(lambda x, y: x+y)

end = datetime.datetime.now()
print(end - start, "h:mm:ss")

# Print the dot product
print("The dot product is:", dot_product)

# Stop the SparkContext
sc.stop()
```

0:00:26.120495 h:mm:ss
The dot product is: 24999888.363398004

- If data is too large, an error like “*OutOfMemoryError: Java heap space*” may occur.
- Default size of driver memory and executor memory is 1 GB, try increasing them.
- pyspark does not support restarting the Spark context.
So, to change the settings, you will need to restart the Jupyter kernel.

Matrix Multiplication with PySpark

```
from pyspark import SparkConf, SparkContext
import findspark
import numpy as np
import datetime

findspark.init()

# Create a SparkConf and SparkContext
conf = SparkConf().setAppName("MatrixMultiplication")
sc = SparkContext.getOrCreate(conf)

n = 100
# Define the matrices as numpy arrays
matrix1 = np.random.rand(n, n)
matrix2 = np.random.rand(n, n)

# Broadcast the second matrix to all nodes
broadcast_matrix = sc.broadcast(matrix2)

start = datetime.datetime.now()

# Convert the first matrix to an RDD of rows
rows_rdd = sc.parallelize(matrix1)

# Perform the matrix multiplication using map and reduce operations
result = rows_rdd.map(lambda row: np.dot(row,
broadcast_matrix.value)).reduce(lambda x, y: x + y)

end = datetime.datetime.now()
print(end - start, "h:mm:ss")

# Print the result
print("The result of matrix multiplication is:\n", result)

# Stop the SparkContext
sc.stop()
```

0:00:09.086623 h:mm:ss
The result of matrix multiplication is:
[2512.5996501 2674.34912392 2493.76738227 2669.46479685 2706.17381393

- First, define two 2-dimensional matrices matrix1 and matrix2 as numpy arrays.
- Broadcast matrix2 to all nodes using sc.broadcast().
- Convert matrix1 to an RDD of rows using sc.parallelize().
- Perform matrix multiplication using map and reduce operations.
- In the map operation, use numpy.dot() to compute the dot product between each row of the first matrix and the entire second matrix (which has been broadcasted).
- In the reduce operation, add up all the results to get the final result of matrix multiplication.

Spark UI

<http://localhost:4040/jobs/>

The screenshot shows the Apache Spark 3.3.2 UI interface for the PySpark_Tutorial application. The top navigation bar includes tabs for Jobs, Stages, Storage, Environment, Executors, SQL / DataFrame, and the current tab, PySpark_Tutorial application UI. Below the navigation bar is a banner for the PySpark_Tutorial - Spark Jobs application.

Spark Jobs (?)

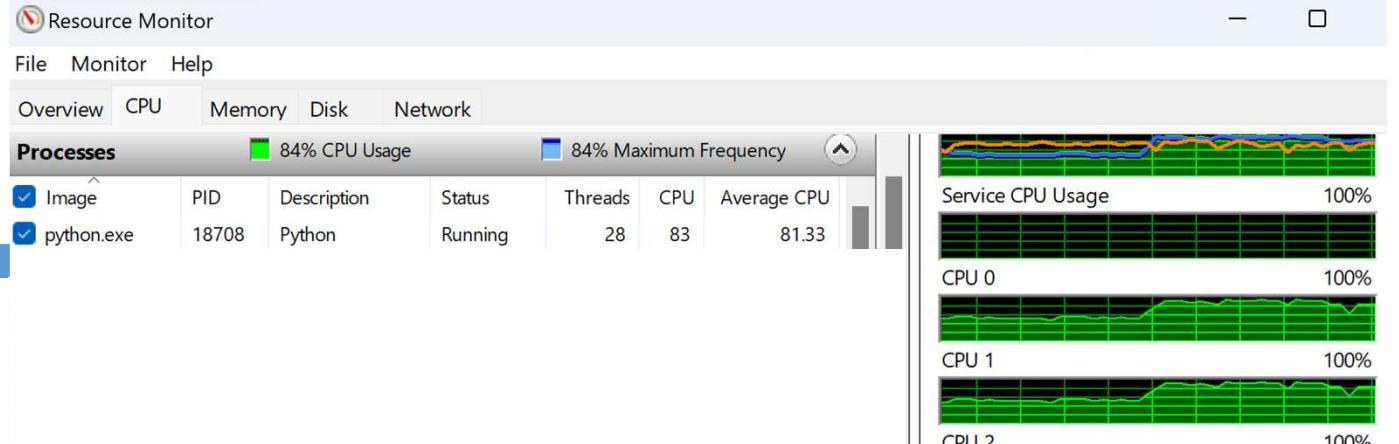
User: veera
Total Uptime: 2.0 min
Scheduling Mode: FIFO
Active Jobs: 1

Event Timeline: Executor driver added at 07:57 on Thu 16 March. A tooltip shows the command: sum at C:\Users\veera\AppData\Local\Temp\ipykernel_49108\2877661582.py:8 (kill).

Active Jobs (1)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0	sum at C:\Users\veera\AppData\Local\Temp\ipykernel_49108\2877661582.py:8 sum at C:\Users\veera\AppData\Local\Temp\ipykernel_49108\2877661582.py:8	2023/03/16 07:58:00	1.2 min	0/1	29/256 (17)

Matrix Multiplication with PySpark



```
!pip install -q pyspark
import numpy as np
from pyspark.mllib.linalg import DenseVector

size = 10000
A = DenseVector(np.random.rand(size, size))
B = DenseVector(np.random.rand(size, size))
%timeit C = DenseVector.dot(A, B)
```

12.6 s ± 361 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

- PySpark has many libraries including MLlib
- MLlib supports vector operations that is not implemented with RDD
- So, it does not create Spark jobs.