

2110625

Data Science Architecture

System architecture for data science solutions

Parallel computing

รศ.ดร. วีระ เหมืองสิน

Assoc. Prof. Veera Muangsin

Department of Computer Engineering

Chulalongkorn University

veera.m@chula.ac.th

Why study computer architecture?

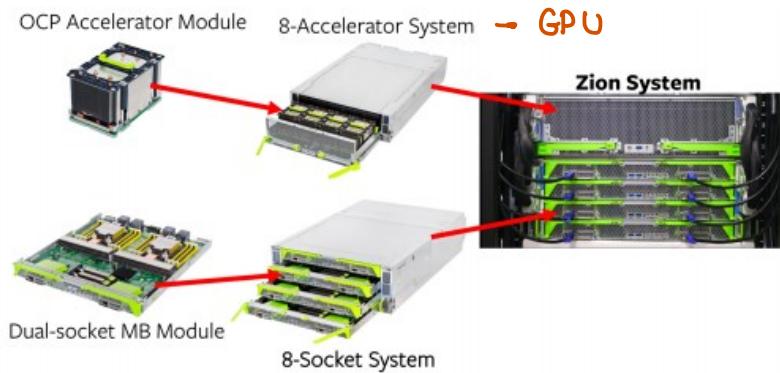
- To understand *how features of computer architecture impact applications* and, conversely, *how requirements of applications lead to the emergence of new hardware and architectural support* ឧប្បជ្ជនេយដារទាំង ១៦ ម៉ោង កំពើតិច អាន វិញ្ញាន សាស្ត្រ
ក្នុងគេង SA
- For system architect: *to design a system that can run applications efficiently* គោរពប្រចាំថ្ងៃ នូវ ការប្រគល់ការងារ
- For SW developer: *to develop software that efficiently utilizes the underlying system* dev ក្នុងក្រុមហ៊ែន ឱ្យប្រើប្រាស់
- For data scientist: *to make informed decisions about selecting hardware and programming tools for data processing and analysis tasks.* តើក អាន កំណែវារស្ម័គ្រាន់
- Examples
 - Software requirement/behavior: *Most memory accesses are to recently used data or nearby data.*
 - Architecture change: *memory hierarchy*
 - cache memory: *small and fast memory for the to-be-used-soon data* មិនគុច ក្នុងការប្រើប្រាស់ការងារទៅក្នុង cache
→ កិច្ចការក្រុមហ៊ែនក្នុងការ optimize cache
 - Software change: *Programmers write program to optimize cache usage.*

Requirements and architecture of a machine learning system

Deep Learning Training in Facebook Data Centers: Design of Scale-up and Scale-out Systems

Maxim Naumov*, John Kim†, Dheevatsa Muthukrishnan, Ming Tang, Yuxin Wang, Whitney Zhao, Serhat Yilmaz, Changkyu Kim, Hui Li, Ming Tang, Junwei Liu, Kumar Nair, Isabel Gao, Bor-Yiing Su, Jiayang Wang, and Ming Tang
Facebook, 1 Hacker Way, Menlo Park, CA 94034

Abstract—Large-scale training is important to ensure high performance and accuracy of machine-learning models. At Facebook we use many different models, including computer vision, video and language models. However, in this paper we focus on the deep learning recommendation models (DLRMs), which are responsible for more than 50% of the training demand in our data centers. Recommendation models present unique challenges in training because they exercise not only compute but also memory capacity as well as memory and network bandwidth. As model size and complexity increase, efficiently scaling training becomes a challenge. To address it we design Zion – Facebook’s next-generation large-memory training platform that consists of both CPUs and accelerators. Also, we discuss the design requirements of future scale-out training systems.



សម្រាប់ការ recommend

User FB

ជាន់ចំណែកផលិត

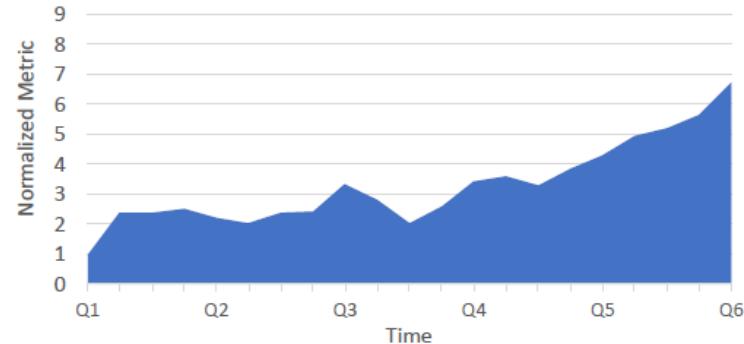
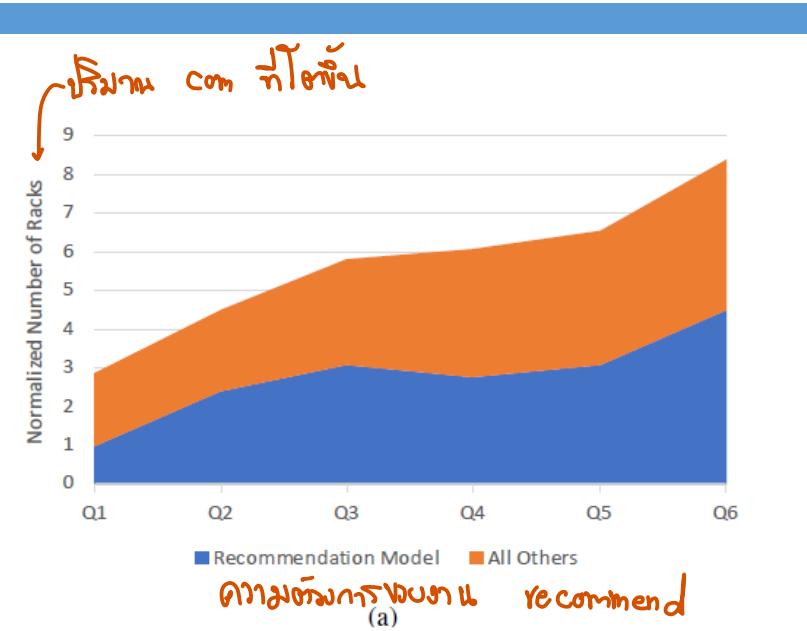


Fig. 1. (a) Server compute demand for training and (b) number of distributed training workflows across Facebook data centers.

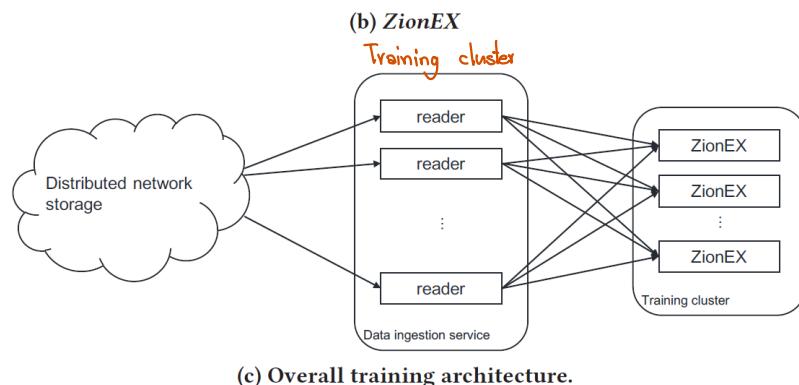
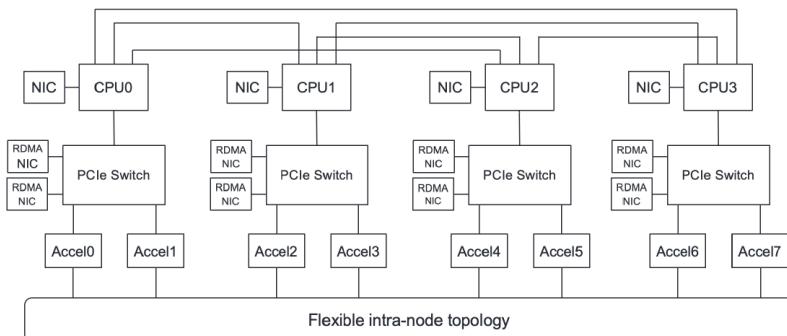
នាម workflow នៃ data centers

Requirements and architecture of a machine learning system

Software-Hardware Co-design for Fast and Scalable Training of Deep Learning Recommendation Models

Dheevatsa Mudigere^{†‡}, Yuchen Hao^{†‡}, Jianyu Huang^{†‡}, Zhihao Jia[§], Andrew Tulloch[‡], Srinivas Sridharan[‡], Xing Liu[‡], Mustafa Ozdal[‡], Jade Nie[‡], Jongsoo Park[‡], Liang Luo[‡], Jie (Amy) Yang[‡], Leon Gao[‡], Dmytro Ivchenko[‡], Aarti Basant[‡], Yuxi Hu[‡], Jiyuan Yang[‡], Ehsan K. Ardestani[‡], Xiaodong Wang[‡], Rakesh Komuravelli[‡], Ching-Hsiang Chu[‡], Serhat Yilmaz[‡], Huayu Li[‡], Jiuyuan Qian[‡], Zhuobo Feng[‡], Yinbin Ma[‡], Junjie Yang[‡], Ellie Wen[‡], Hong Li[‡], Lin Yang[‡], Chonglin Sun[‡], Whitney Zhao[‡], Dmitry Melts[‡], Krishna Dhulipala[‡], KR Kishore[‡], Tyler Graj[‡], Assaf Eisenman[‡], Kiran Kumar Matam[‡], Adi Gangidi[‡], Guoqiang Jerry Chen[‡], Manoj Krishnan[‡], Avinash Nayak[‡], Krishnakumar Nair[‡], Bharath Muthiah[‡], Mahmoud khorashadi[‡], Pallab Bhattacharya[‡], Petr Lapukhov[‡], Maxim Naumov[‡], Ajit Mathew[‡], Lin Qiao[‡], Mikhail Smelyanskiy[‡], Bill Jia[‡], Vijay Rao[‡]

[‡]Meta Platforms, [§]Carnegie Mellon University



ABSTRACT

Deep learning recommendation models have been used across many business-critical applications and are the single largest AI application in data-centers. In this paper, we present Neo, a software-hardware co-designed system for high-performance distributed training of large-scale DLRMs. Neo employs a novel 4D parallelism strategy that combines table-wise, row-wise, column-wise, and data parallelism for training massive embedding operators in DLRMs. In addition, Neo enables extremely high-performance and memory-efficient embedding computations using a variety of critical system optimizations, including hybrid kernel fusion, software-managed caching, and quality-preserving compression. Finally, Neo is paired with *ZionEX*, a new hardware platform co-designed with Neo's 4D parallelism for optimizing communications for large-scale DLRM training. Our evaluation on 128 GPUs using 16 *ZionEX* nodes shows that Neo outperforms existing systems by up to 40× for training 12-trillion-parameter DLRM models deployed in production.

கீழாண்மை நிதி போன்றுள்ள Neo

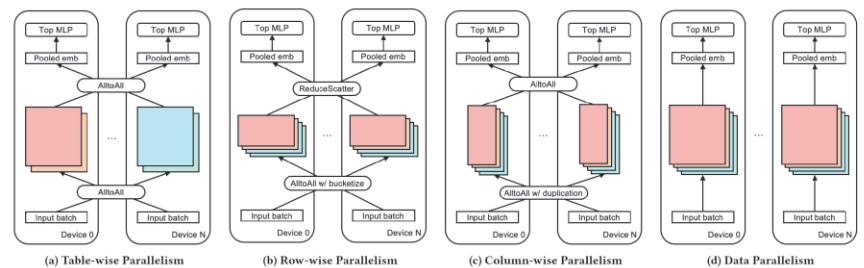
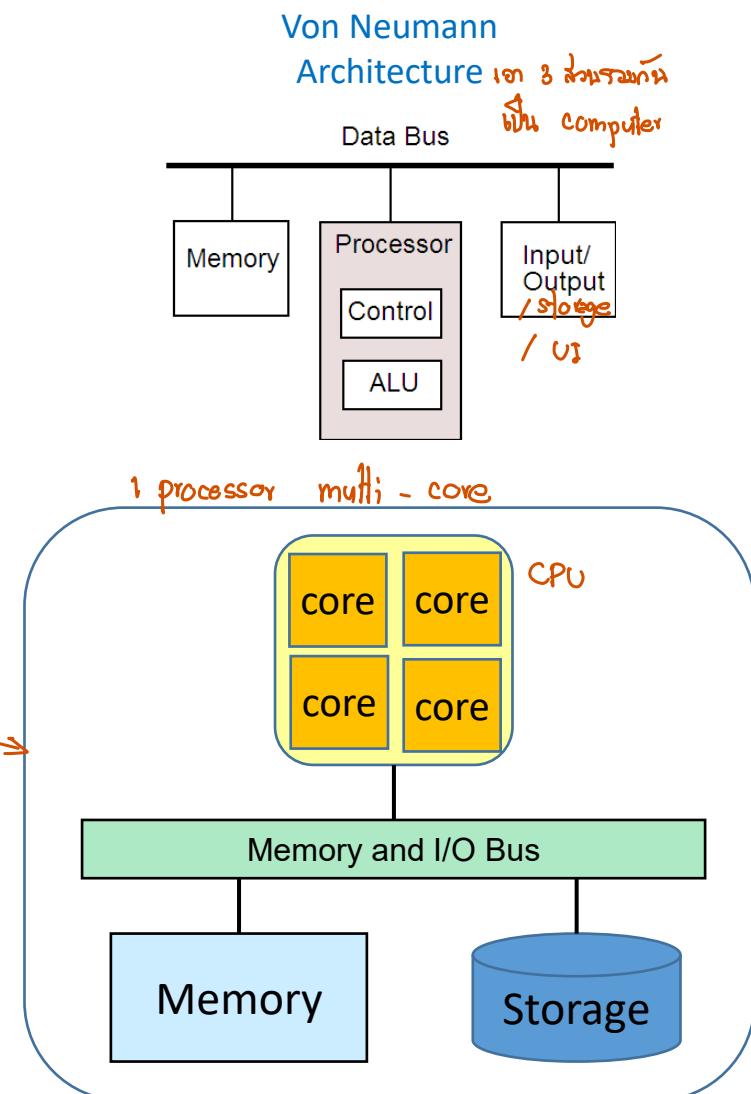
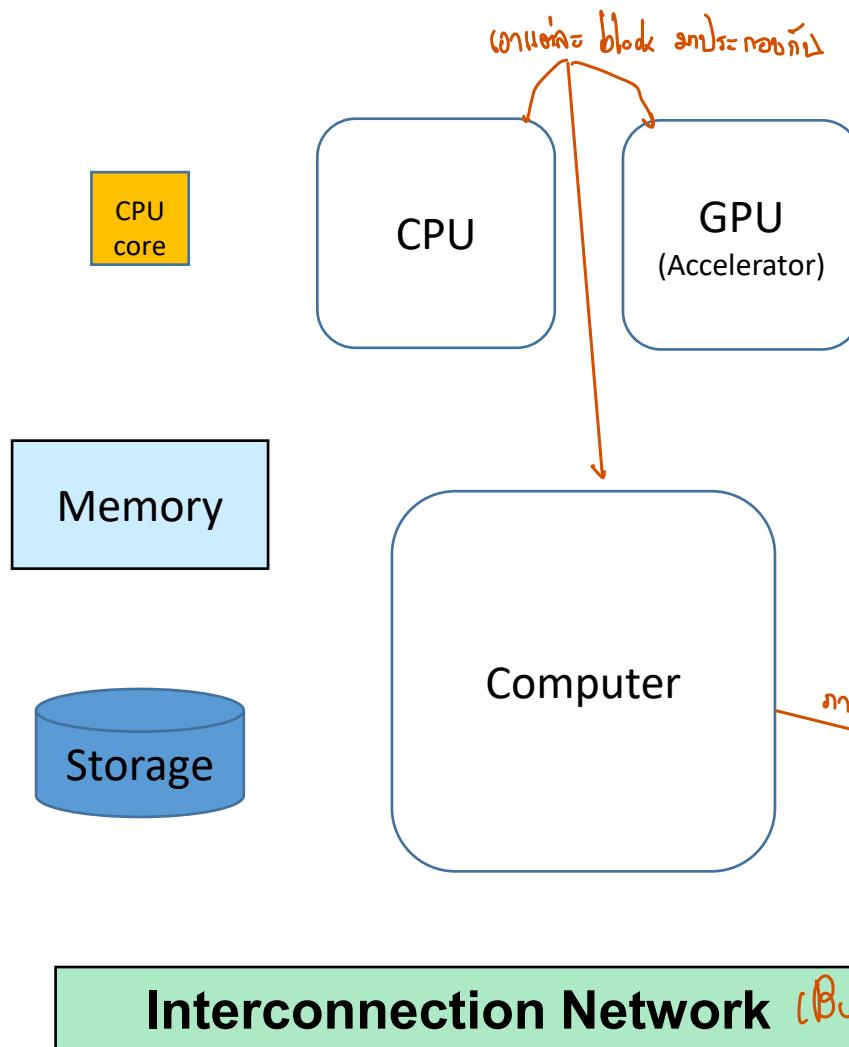


Figure 5: Embedding table sharding schemes with different implications on the communication cost, load balancing and memory requirement. Bottom MLP is omitted in this figure for simplicity of illustration.

Processing Architecture

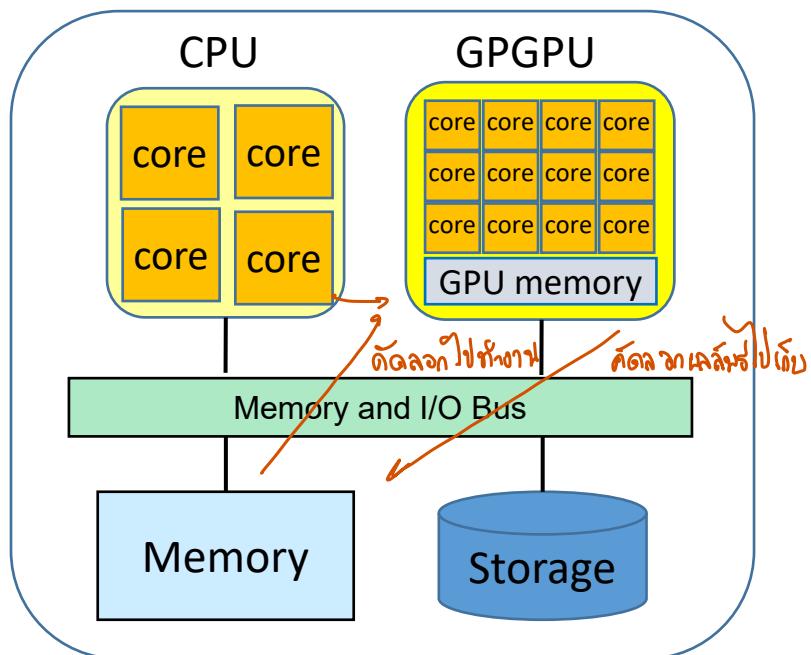
Building blocks of a computer system



GPGPU

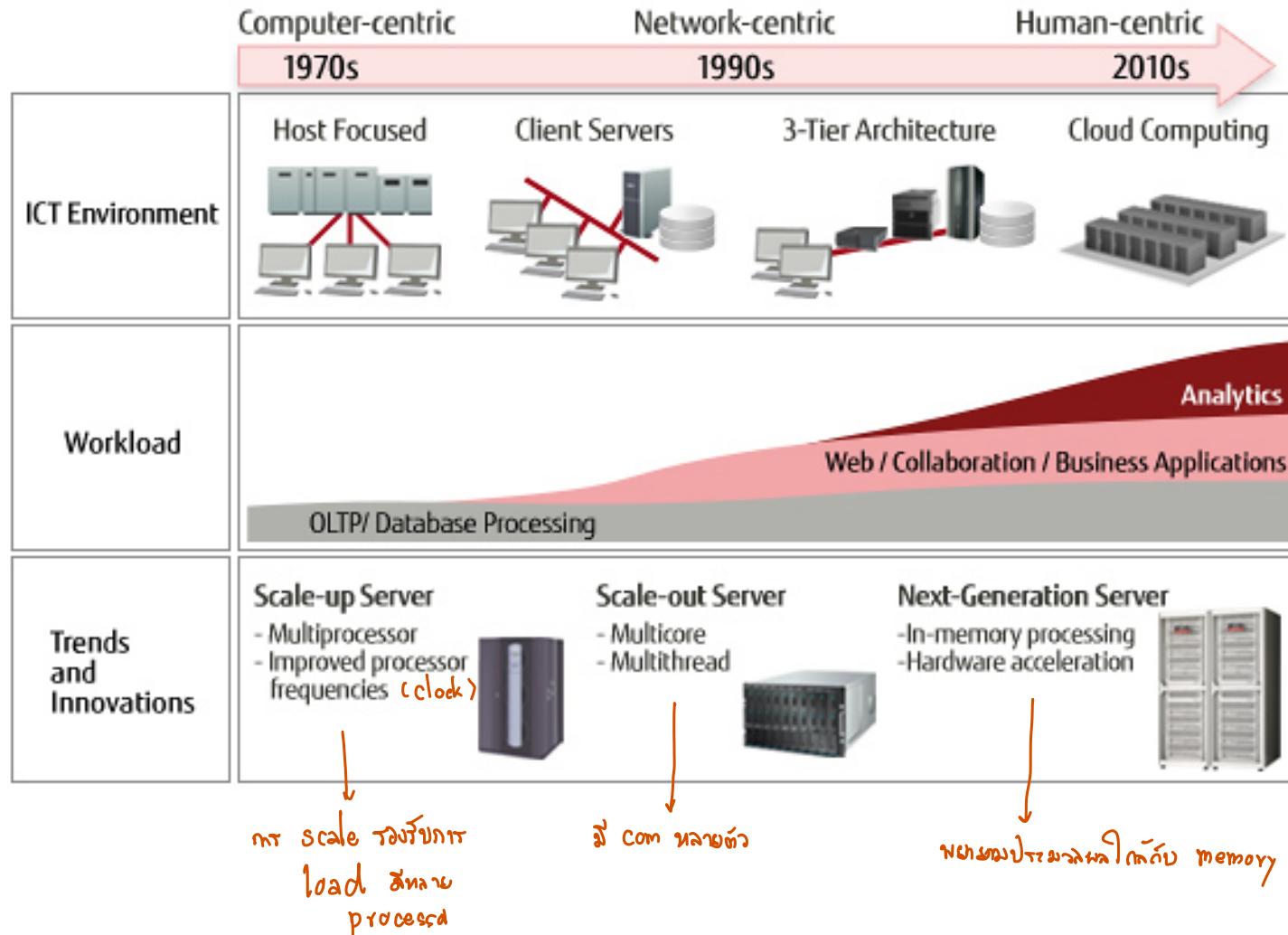
កំណត់ថា CPU បង្កិច

ដើម្បីរួម Video card ទូរសព្ទ



- GPU is an accelerator with vector processing architecture for graphics and image processing. GPU មែនគឺជាបច្ចេកទេសការងារក្នុង graphic, matrix
- GPU offers more parallelism via threads and vectorization.
- General-purpose GPU (GPGPU) provides programming model (e.g. CUDA, OpenCL) for other applications that exhibit data parallelism. ការស្រួលប្រព័ន្ធឌីជីថាមពេល GPGPU នៃបច្ចេកទេសការងារក្នុង CG
- CPU is needed to run OS and programs. CPU ដែលត្រូវមែនបច្ចេកទេសការងារក្នុង GPU
- CPU offloads some computation tasks to GPU. ការកែសម្រេច code និង data ទូទៅ GPU ប្រគល់
- GPU has its own memory. GPU មែនមាន memory ទូទៅ
- Data must be transferred between system memory and GPU memory
ការគ្រប់គ្រងទូទៅ overhead ឡើង

Evolution of Computer Systems



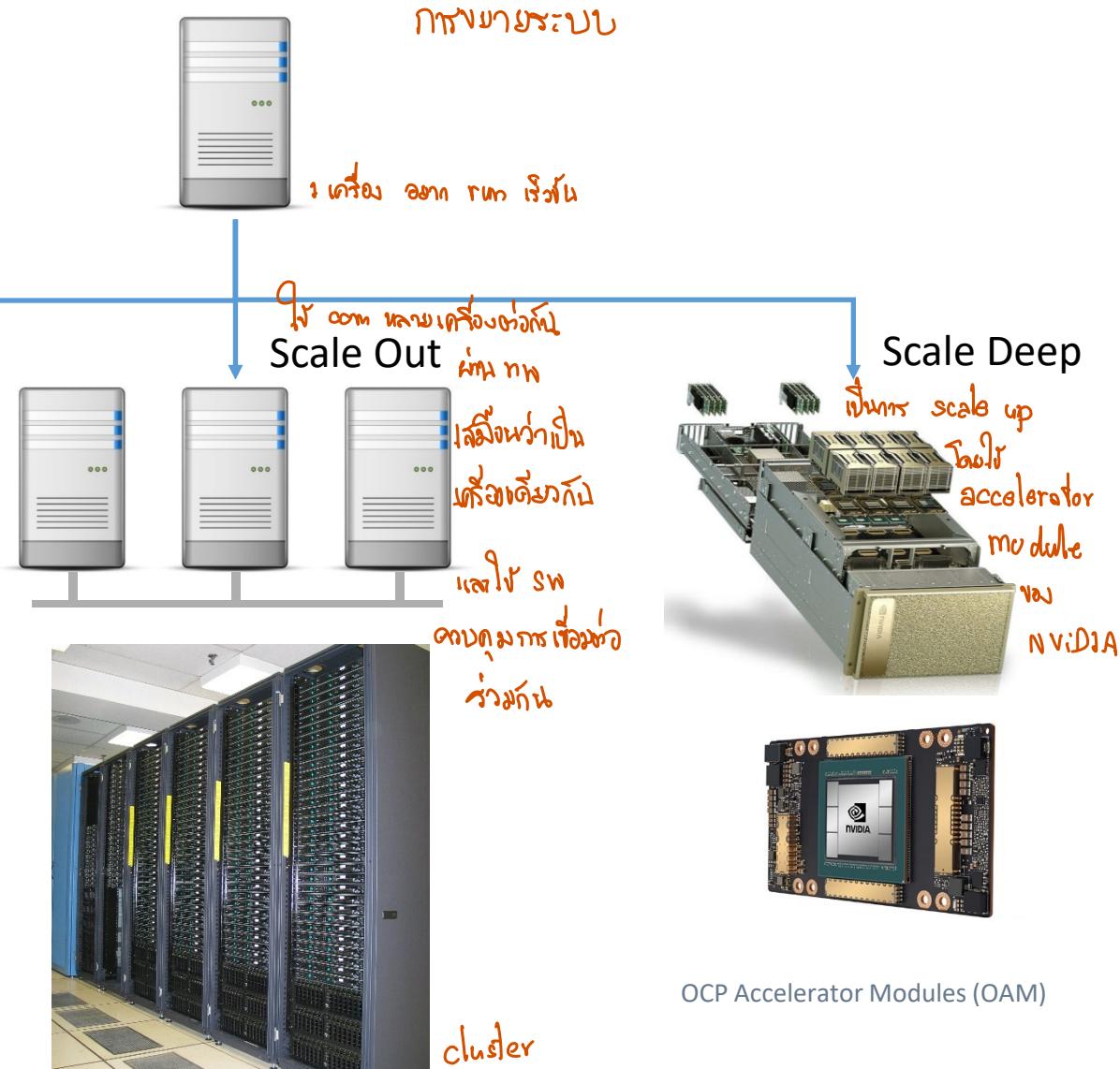
Scale Up, Scale Out, Scale Deep

เพิ่มประสิทธิภาพการทำงานของเครื่องเดียว
ที่มี คือ memory, core, processor

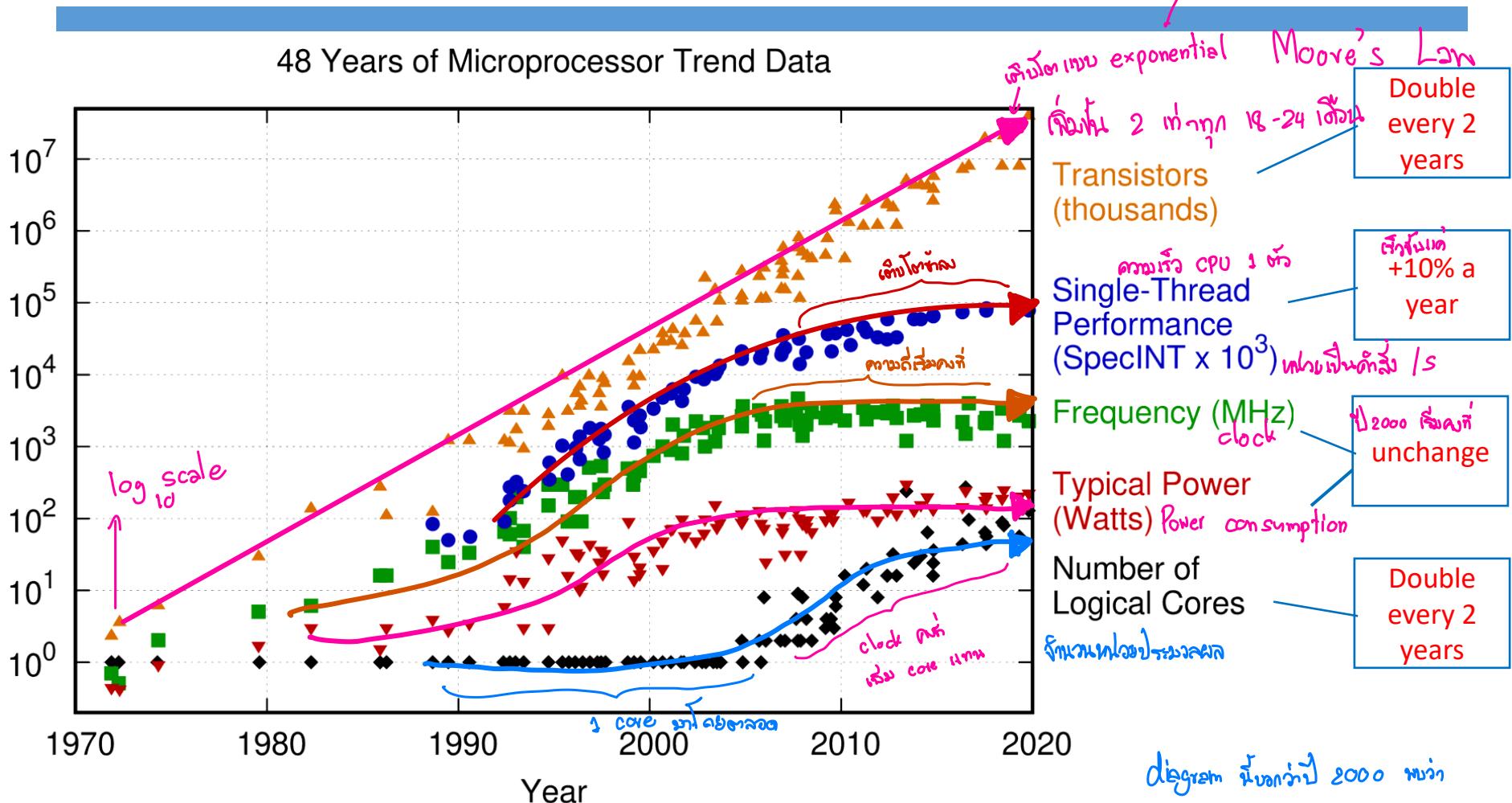
Scale Up, storage



ข้อดีคือ ได้ผลลัพธ์ทันที !!!



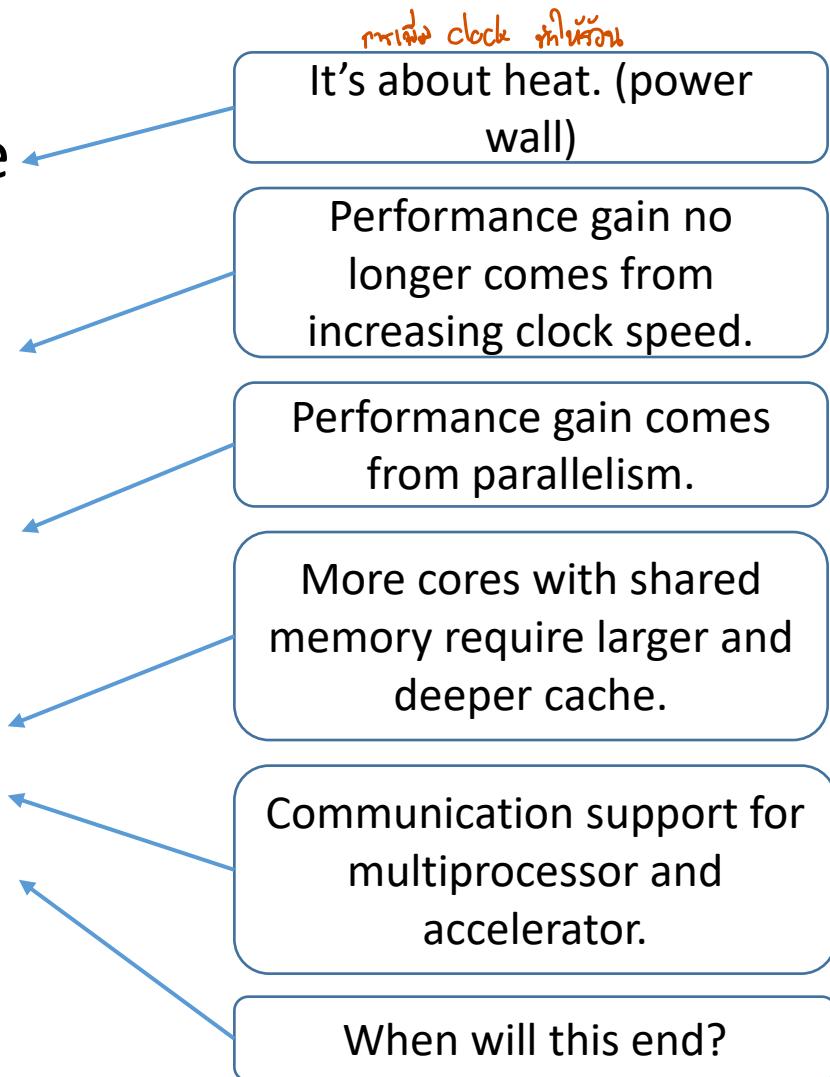
Microprocessor Trends



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2019 by K. Rupp

Microprocessor Trends

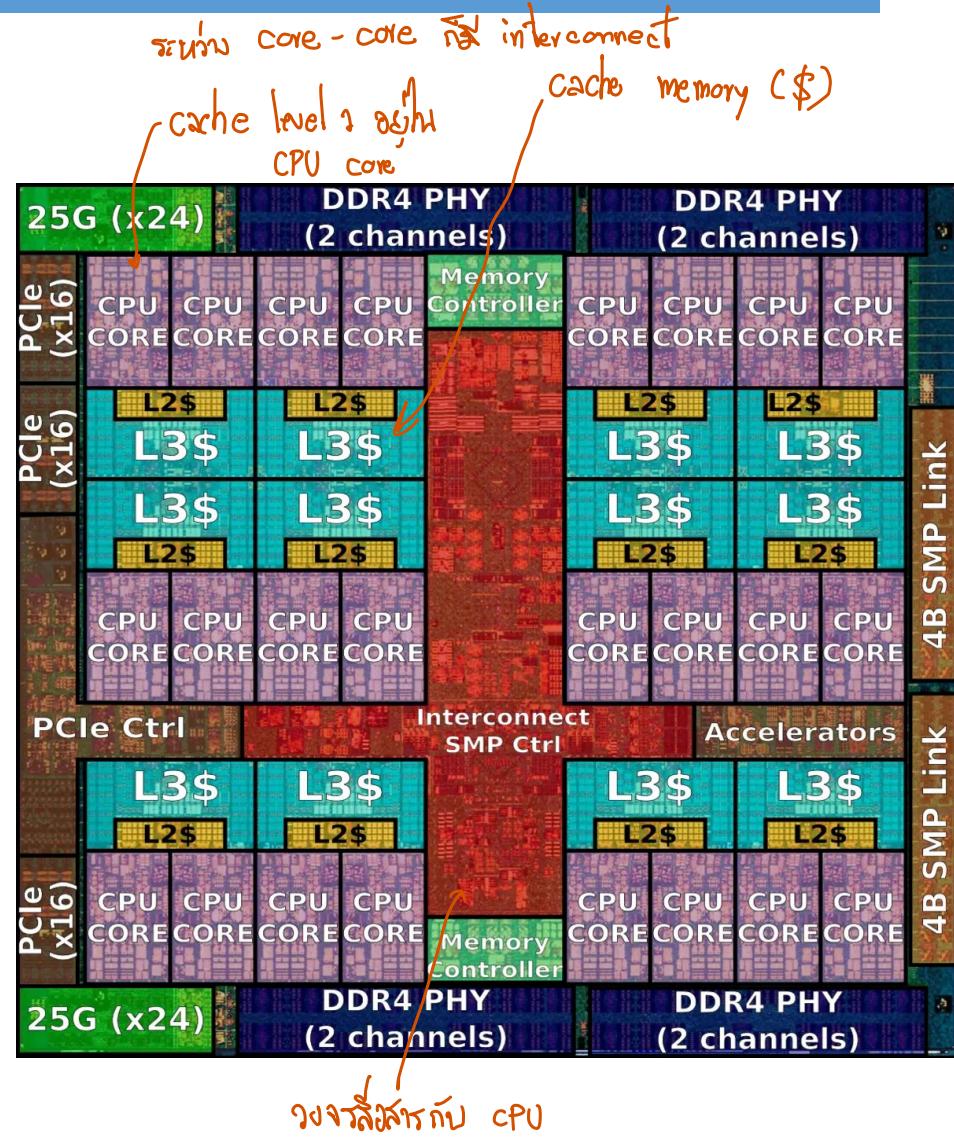
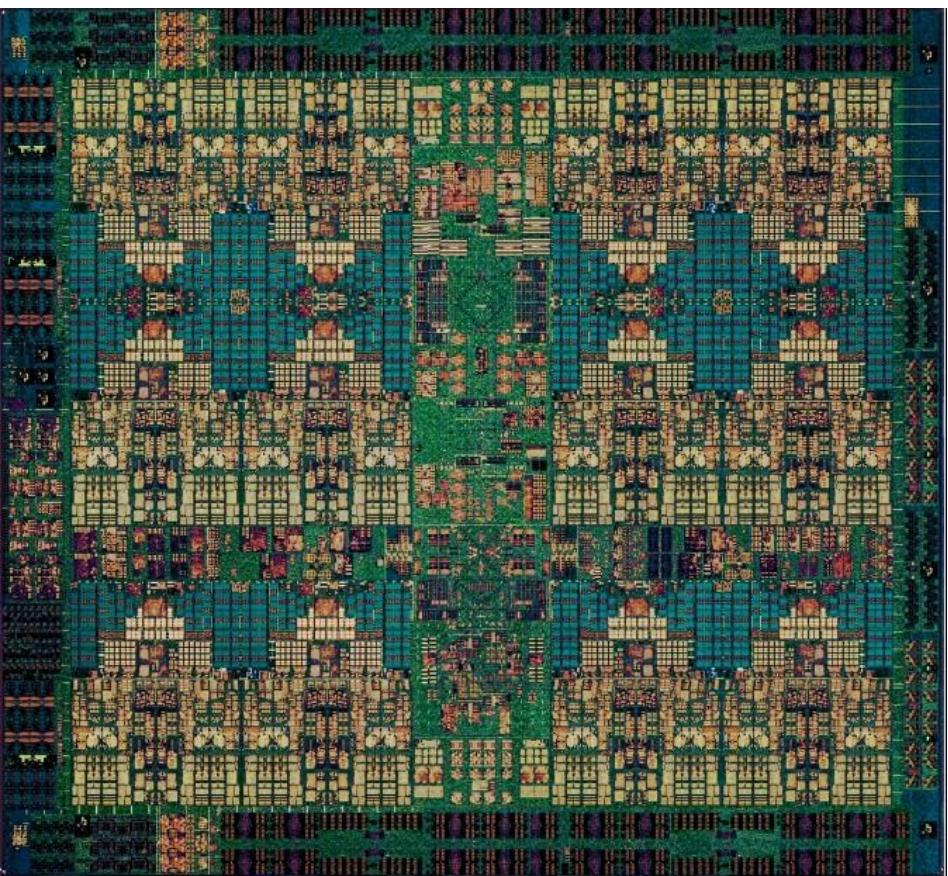
- Clock speed and power have stopped increasing.
- Single core speed still increase but slowly.
- Number of logical cores grows exponentially.
- Number of transistors still grows exponentially.



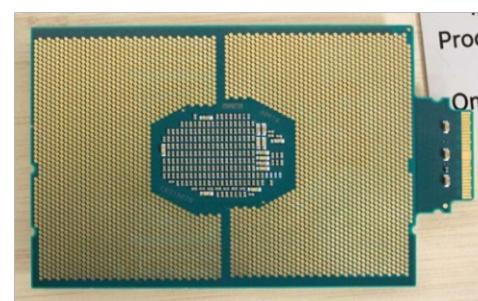
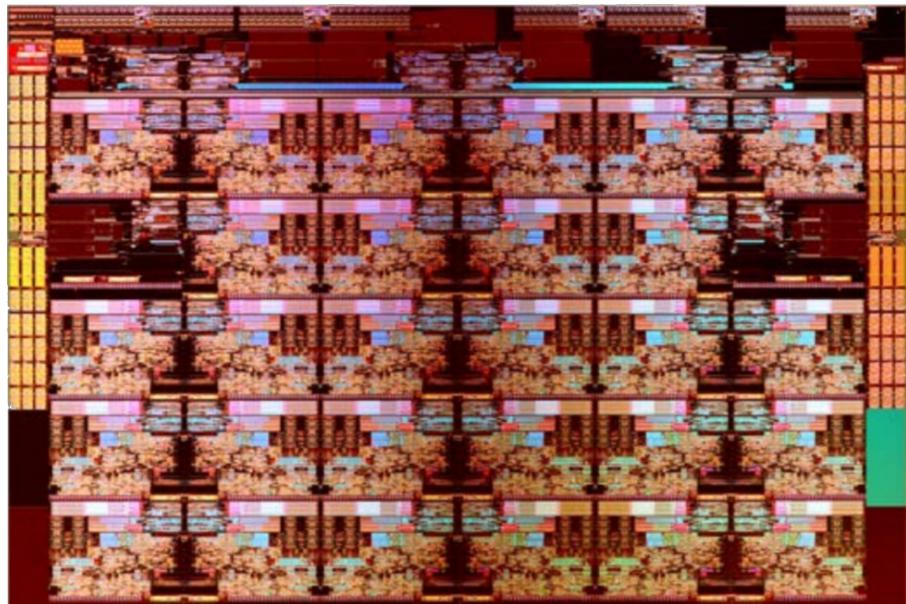
Other trends

- Accelerator processors with more parallelism
 - GPGPU (General Purpose Graphics Processing Unit)
 - Vector and array processing
- AI Accelerator (aka Deep Learning Processor, etc.)
 - TPU (Tensor Processing Unit) ทําชิ้นคําสั่งที่เน้นภารกิจ workload ตัว NN
 - Google
 - Optical Neural Network

IBM Power9 24-core Processor (2017)

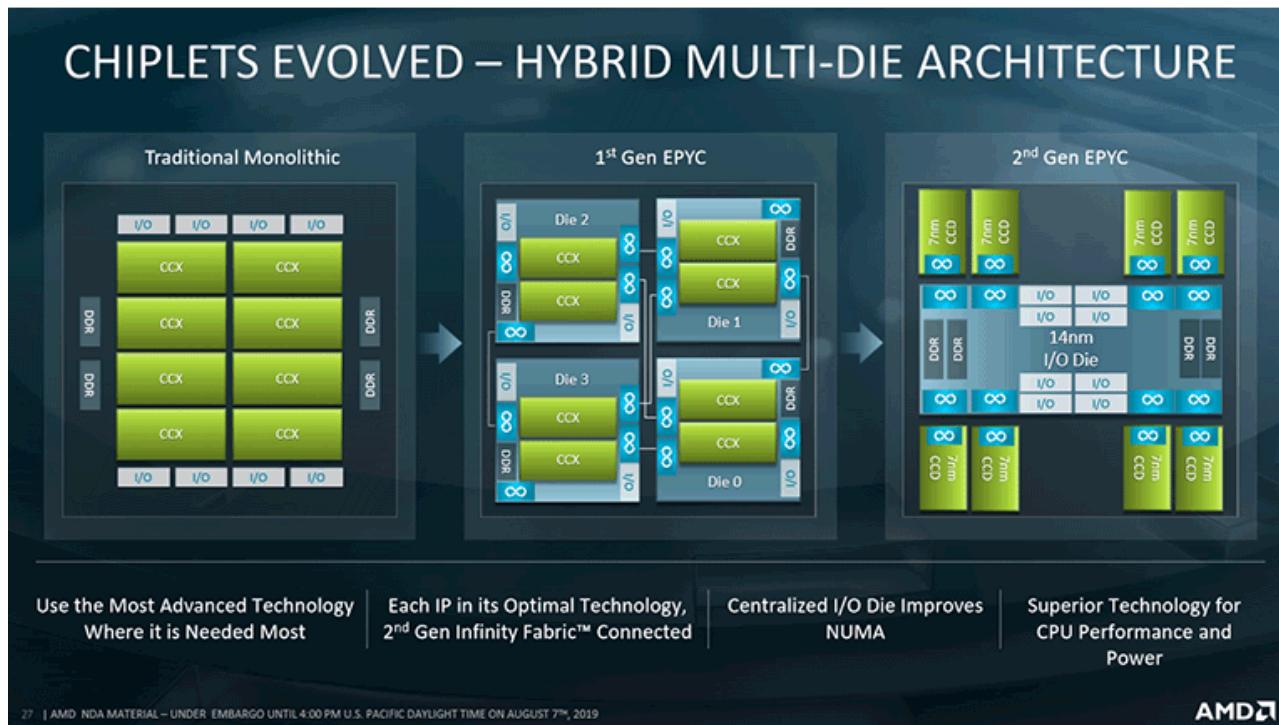


Intel Xeon Platinum 28-core Processor (2017)



Intel Xeon 8176F, Launch Date Q3'17, \$8874

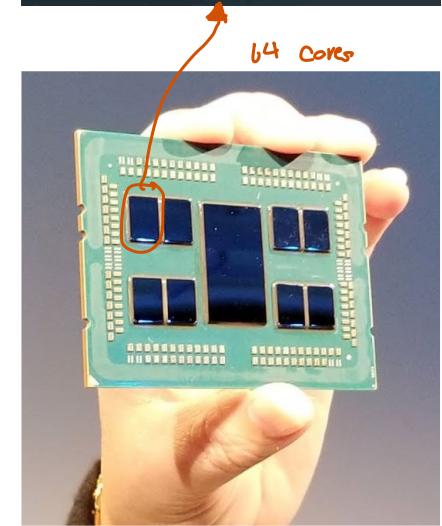
Multi-Die Architecture



27 | AMD NDA MATERIAL – UNDER EMBARGO UNTIL 4:00 PM U.S. PACIFIC DAYLIGHT TIME ON AUGUST 7th, 2019

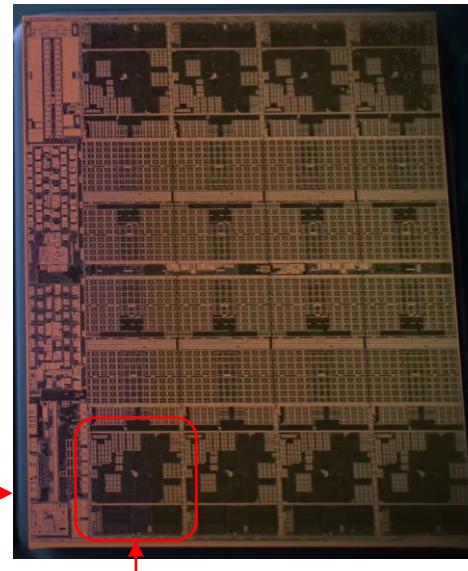
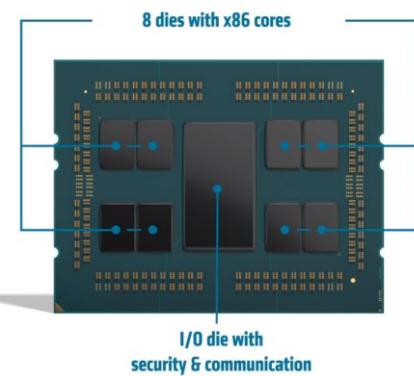
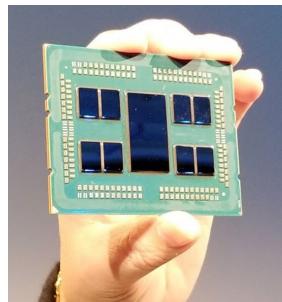
ក្រុមការ core ឬជាក្រុម chip ទាំងអស់នេះ ត្រូវបានចាយការ

Making a multi-die CPU is more economical than a monolithic-die CPU.

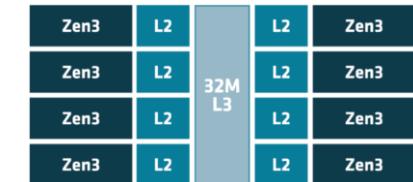
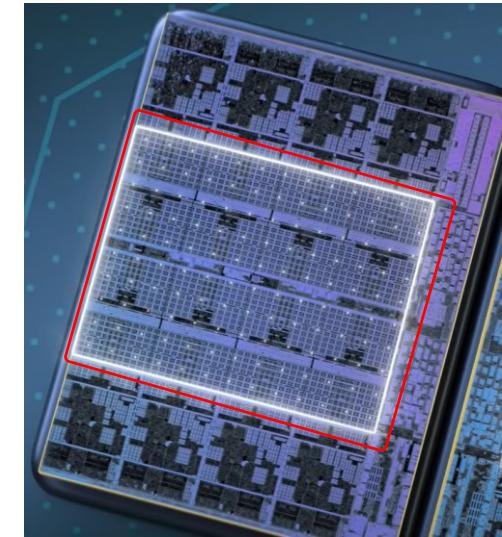


AMD EPYC 7763 Processor (2021)

- 64 cores / 128 threads
- Multi-die architecture
- 8 dies, 8 cores per die
- Each core has 32KB L1 cache, and 512 KB L2 cache.
- Each die has 32 MB L3 cache per die → total 256 MB L3 cache
- Transistor size 7nm
- 280 Watt 
- Launch 2021, \$8000



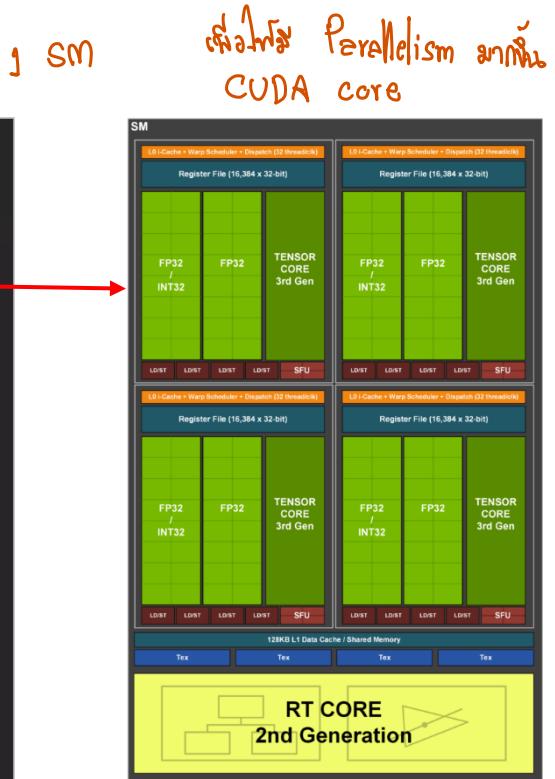
CPU core



Total 256 MB L3-cache

Nvidia GeForce RTX 3090 GPU (2020)

- 10,496 CUDA cores
- 82 SM (streaming multiprocessors)
- 24 GB RAM
- ~36 TFLOPS

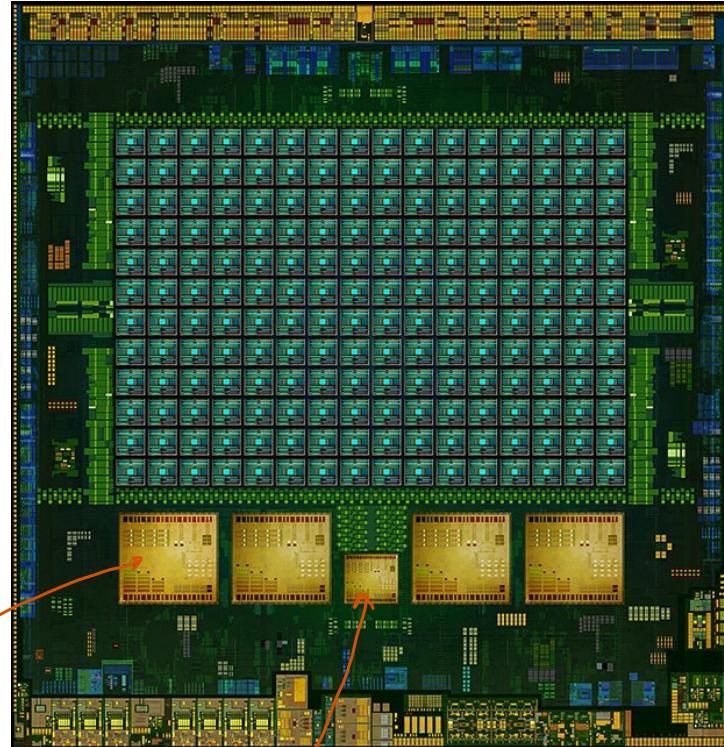


Integrated CPU GPU (2014) និង CPU + GPU នៃ chip តើម្បនា

- Nvidia Tegra K1
 - Mobile processor
 - System-on-a-chip: SoC
 - 4-Plus-1CPU cores
 - 192 GPU cores
 - 5W កន្លែង 5 watt

គ្រែសរែសក្ខុបិជ្ជ mobile device ដោយបានចំណែកជាមុនការ
តាតការងារឡាតក core ទាំងអស់ នឹងត្រួតពិនិត្យ

CPU



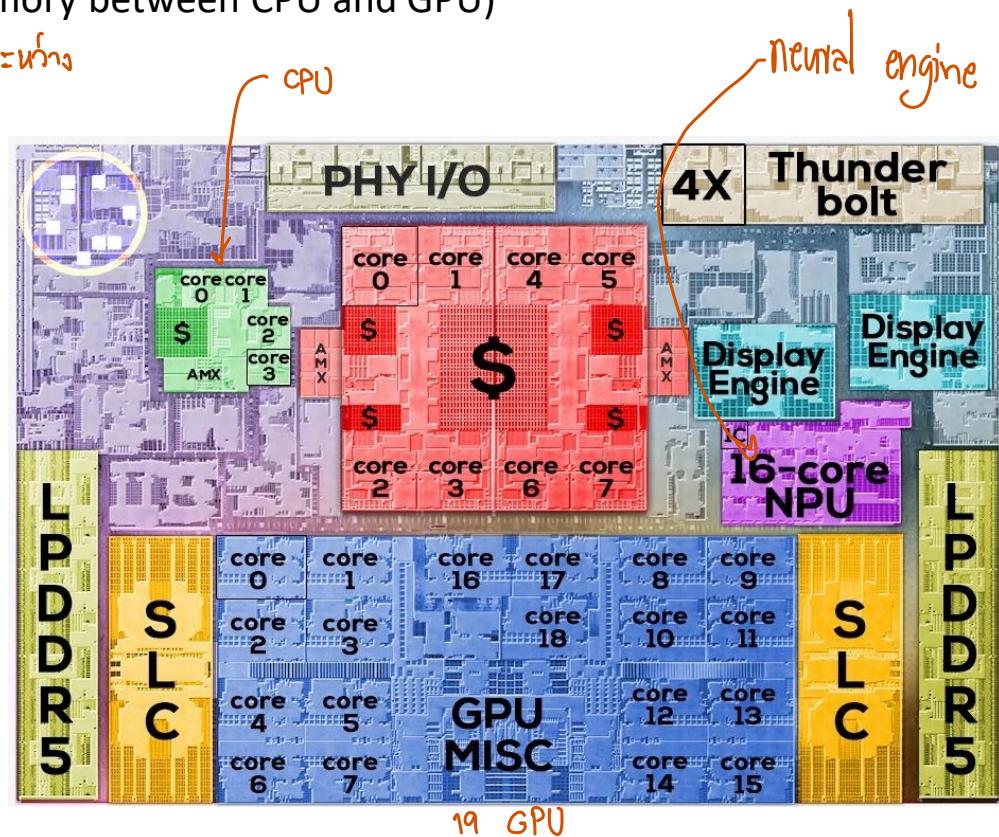
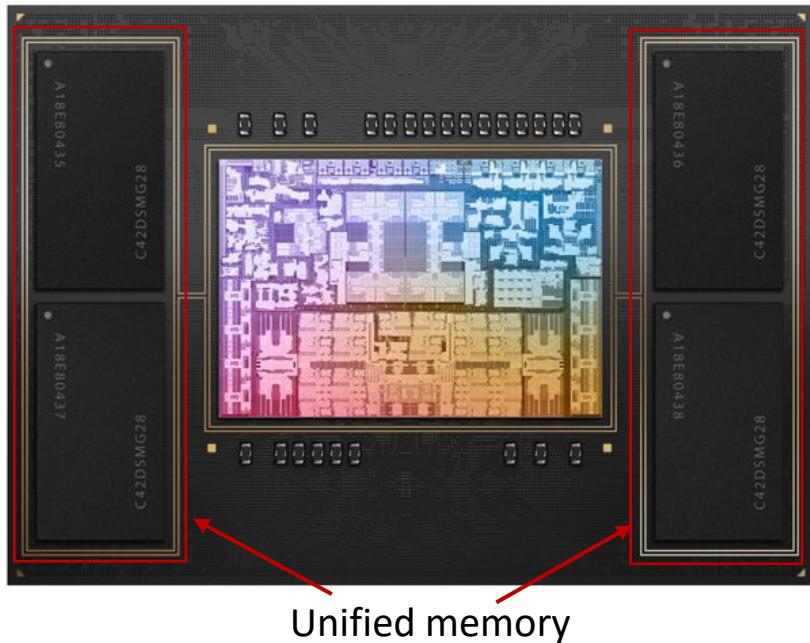
ការងារនៅក្នុងគ្រែ 1 CPU
ដើម្បីបានអាច
ឱ្យអាជីវកម្ម GPU, CPU ទាន់បានការងារ

System-on-Chip

- Apple M2 Pro
- 12 CPU cores (8 high-performance cores + 4 high-efficiency cores)
- 19 GPU cores + 16-core Neural Engine
- 32 GB Unified Memory (shared memory between CPU and GPU)

尽量减少 CPU 和 GPU 之间的数据传输开销

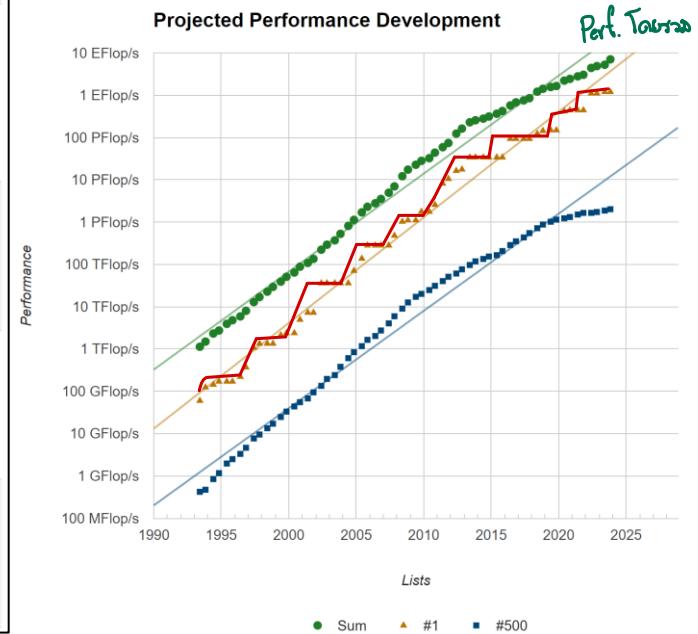
CPU 与 GPU



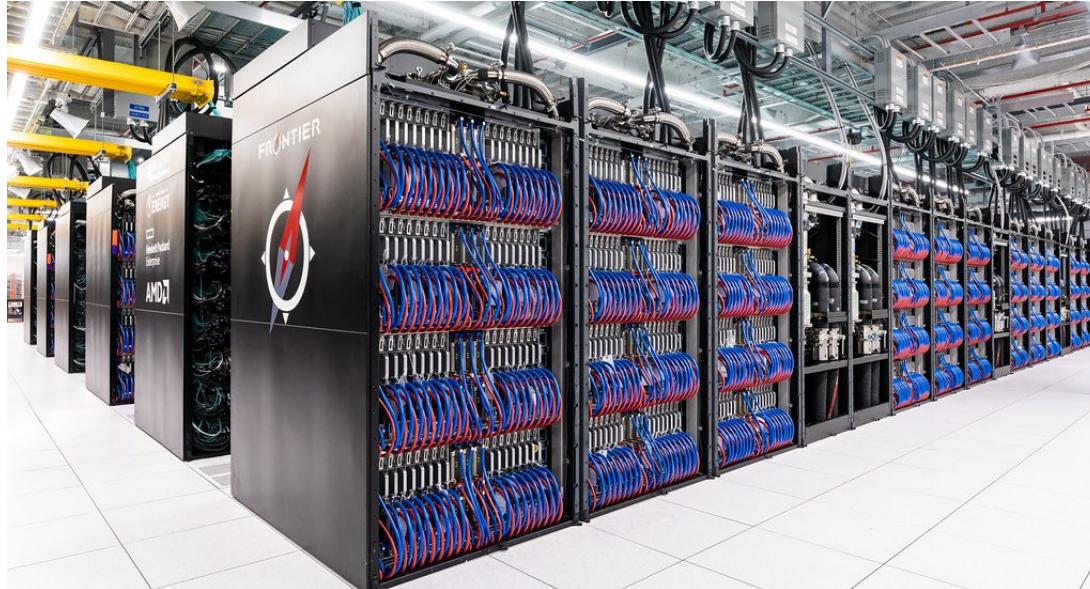
State of the Art: Top500 Nov 2023

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States	8,699,904	1,194.00	1,679.82	22,703
2	Aurora - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, Intel DOE/SC/Argonne National Laboratory United States	4,742,808	585.34	1,059.33	24,687
3	Eagle - Microsoft NDv5, Xeon Platinum 8480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR, Microsoft Microsoft Azure United States	1,123,200	561.20	846.84	
4	Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899
5	LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland	2,752,704	379.70	531.51	7,107
94	LANTA - HPE Cray EX235n, AMD EPYC 7713 64C 2GHz, NVIDIA A100 40GB, Slingshot-11, HPE NSTDA Supercomputer Center (ThaiSC) Thailand	87,296	8.15	13.77	310
500	TX-Green2 - PowerEdge C6420, Xeon Platinum 8260 24C 2.4GHz, 25G Ethernet, ACTION MIT Lincoln Laboratory Supercomputing Center United States	43,200	2.02	53.08	

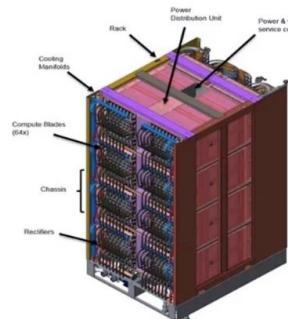
กินไฟเท่ากันบ้าน
20,000 กว่ากิกะวัตต์



Frontier Supercomputer



system

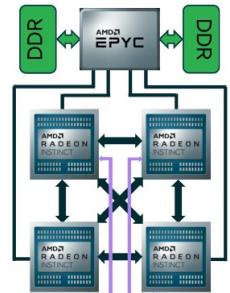


rack

Compute blade
• 2 AMD nodes



node



CPU + GPU

Fugaku Supercomputer



Topics ▾

Reports ▾

Blogs ▾

Multimedia ▾

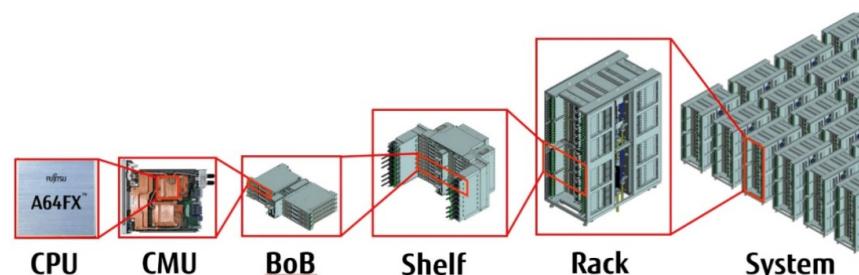
Tech Talk | Computing | Hardware

29 Jun 2020 | 16:00 GMT

Japan's Fugaku Supercomputer Completes First-Ever Sweep of High-Performance Benchmarks

Fugaku, based on the Arm architecture, races ahead in processing speed, running real-world applications, and in a new benchmark: evaluating calculations used in artificial intelligence

By John Boyd



<https://spectrum.ieee.org/tech-talk/computing/hardware/japans-fugaku-supercomputer-is-first-in-the-world-to-simultaneously-top-all-high-performance-benchmarks>

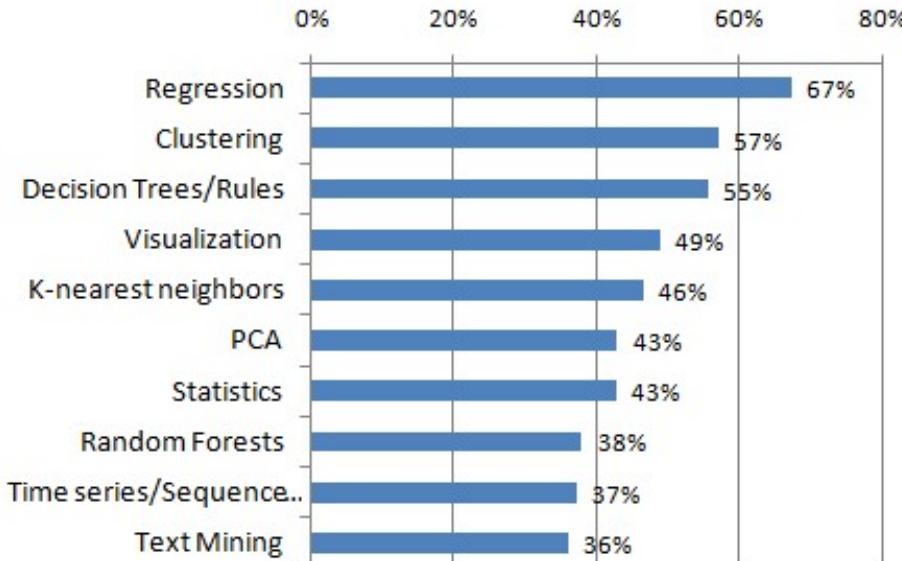
System Architecture Issues Software

Application Requirements

Requirements of Data Science Applications

និងខាងក្រោមគឺជាការប្រើប្រាស់នៃព័ត៌មានដើម្បី និងទទួលបាយព័ត៌មាន។

Top 10 Algorithms & Methods used by Data Scientists



<https://www.kdnuggets.com/2016/09/poll-algorithms-used-data-scientists.html>

Which methods/algorithms you used in the past 12 months for an actual Data Science-related application?

Common Data Analysis Methods

- Regression
 - least squares method *technique සංඛ්‍යා තුළුවා*
- Clustering
 - k-means clustering, hierarchical clustering
- Classification
 - decision tree, k-nearest neighbors, support vector machine, neural network
- Recommendation
 - collaborative filtering
- Text analysis
 - sentiment analysis, translation
- Network Analysis
 - centrality, community detection

Basic data structures and algorithms

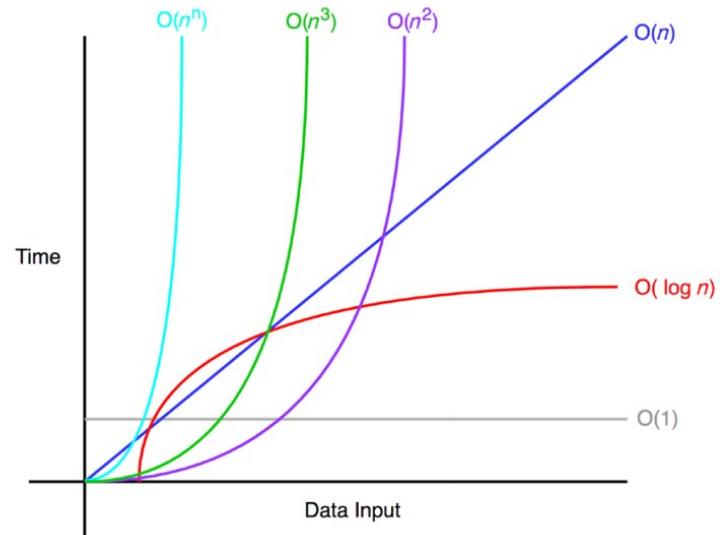
All analysis methods are constructed from basic data structures and algorithms: ຖຸກ method ອຸທະວາງໂອທ ການສະໜັບສະໜັດຂອງລົງທຶນເຄະຫຼິດ ການີ້ກັນ

- Array
 - sort, search, average
- Matrix
 - dot product, matrix multiplication
- Tree
 - search
- Graph
 - shortest path

Computational Complexity

ເຫດຄວາມຮັບອໍານວຍປະລຸງນາ, ຄວາມຫາກໃນການຄົ່ນຫຼຸດ

- $O(1)$
 - Read an element in an array
- $O(\log n)$ ວິທີມູນຄົວຈິນ ເລກໄໝໄວ້ເພີ່ມຂຶ້ນມາກ
 - Binary search
- $O(n)$
 - Average, vector-vector addition, vector-vector dot product
- $O(n \log n)$
 - QuickSort, merge sort
- $O(n^2)$
 - Insertion Sort, matrix-vector multiplication
- $O(n^3)$
 - All-pair shortest path, matrix-matrix multiplication



Computing Intensive vs. Memory Intensive Problems

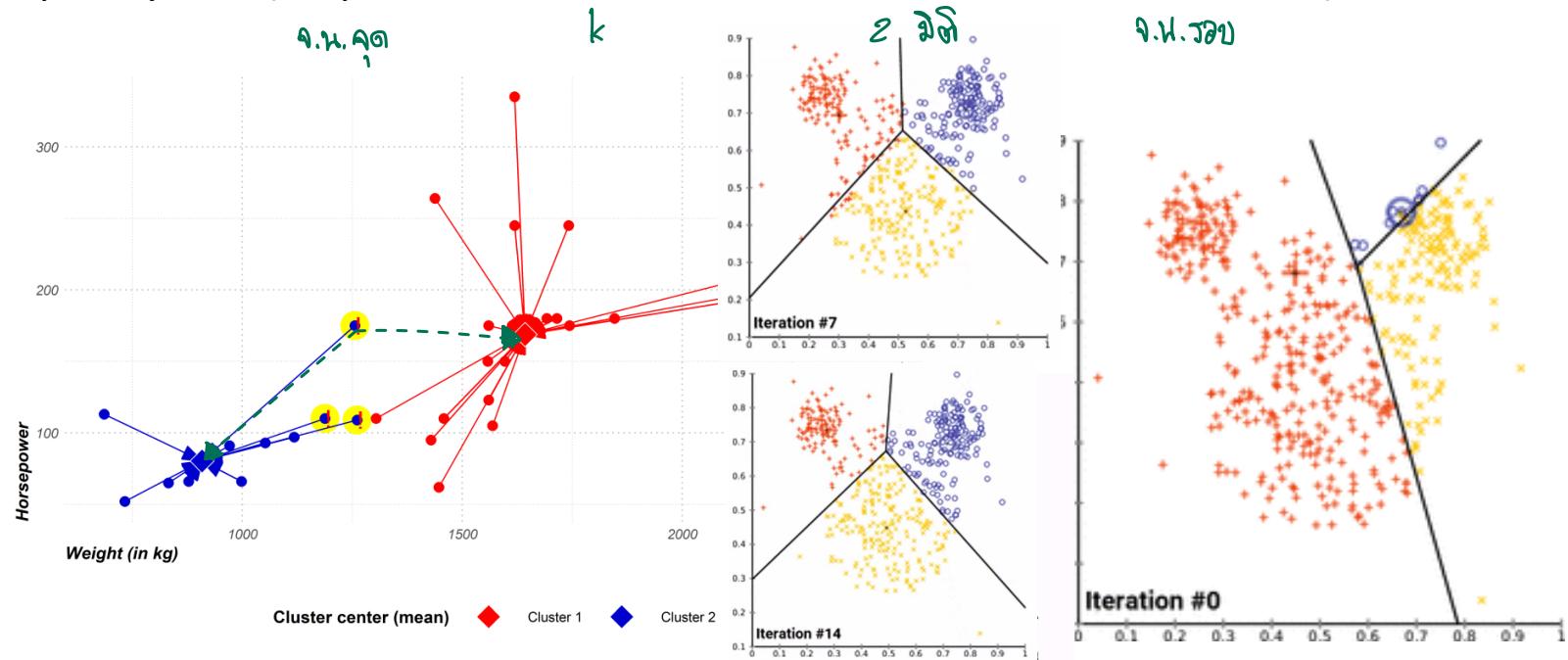
$O(n^3)$ vs. $O(N)$

គម្រោងនៃការប្រើប្រាស់បច្ចេកទេស តួនៅ $\frac{O(N)}{\downarrow}$ វិវាទនូវ $O(n^3)$ យើង
ដែលជា big data

A data science problem can be both computing intensive and memory intensive.

Example: K-Means Clustering

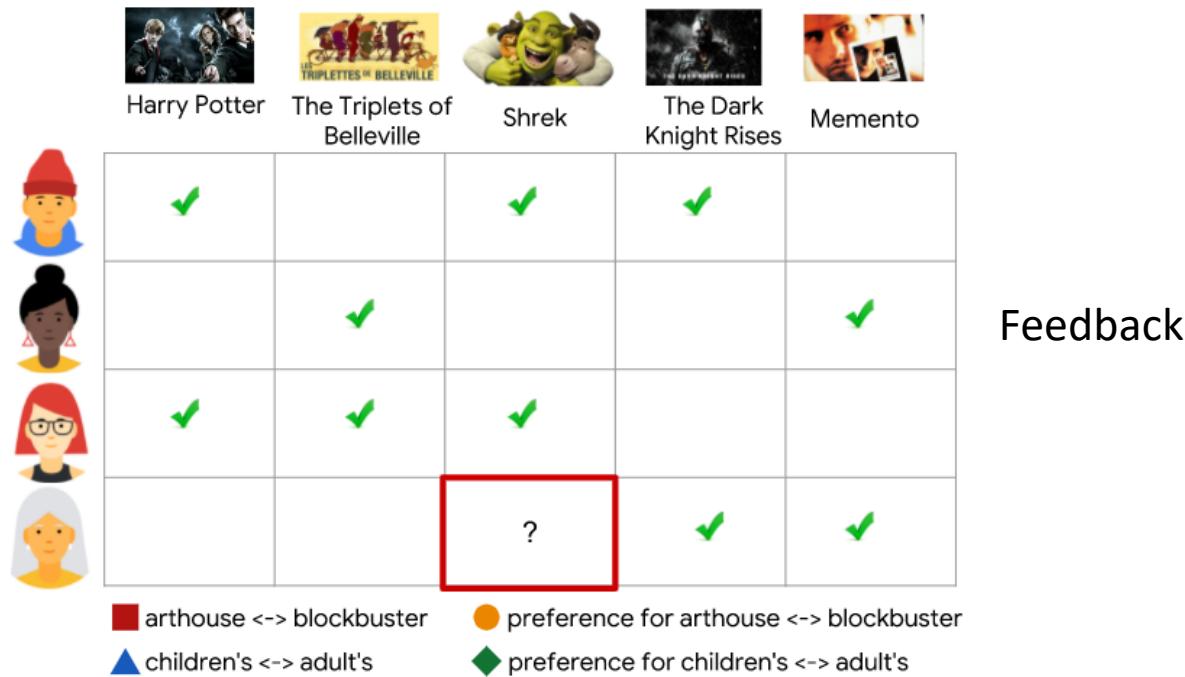
- An algorithm to partition data into clusters (groups of similar data points)
- The difference, or **distance**, between a pair of data points can be calculated from their properties.
- Each cluster has a **centroid** point representing average of each property of its members.
- The algorithm repeats these two steps until the result does not change:
 - Assign each point to the cluster with the nearest centroid.
 - Calculate centroids of each cluster.
- Complexity $\rightarrow O(\text{no.points} * \text{no.clusters} * \text{no.dimensions} * \text{no.iterations})$



Example: Recommendation Systems based on Matrix Factorization

- Feedback matrix contain users' respond to movies. Most cells are empty (sparse matrix).
- Goal: predict users' feedback on movies. In other words, fill the empty cells in *feedback*.

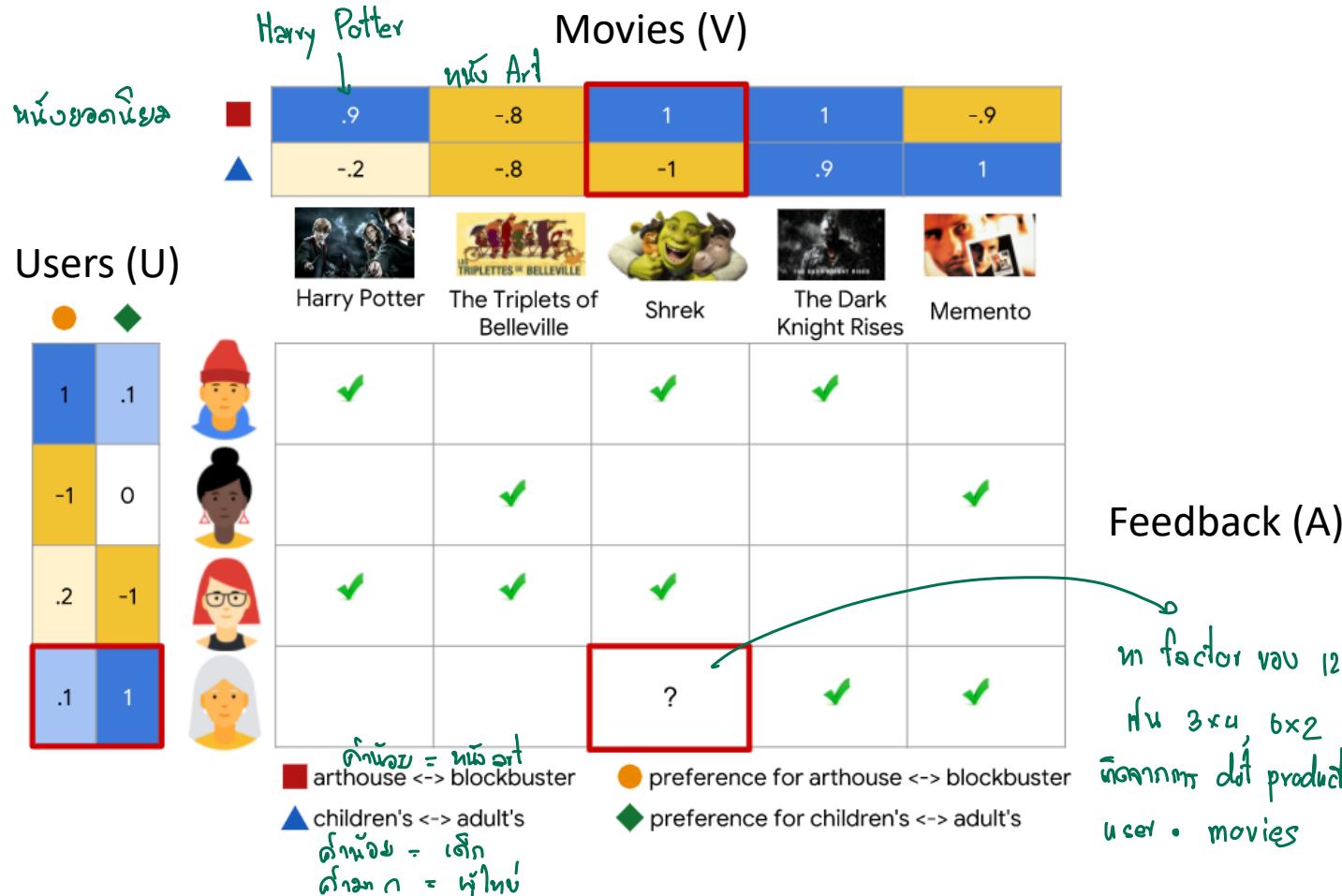
Recommendation គេងប្រើការព័ត៌មានឱ្យរាយ



ແກ້ໄຂຕາງໜູ້

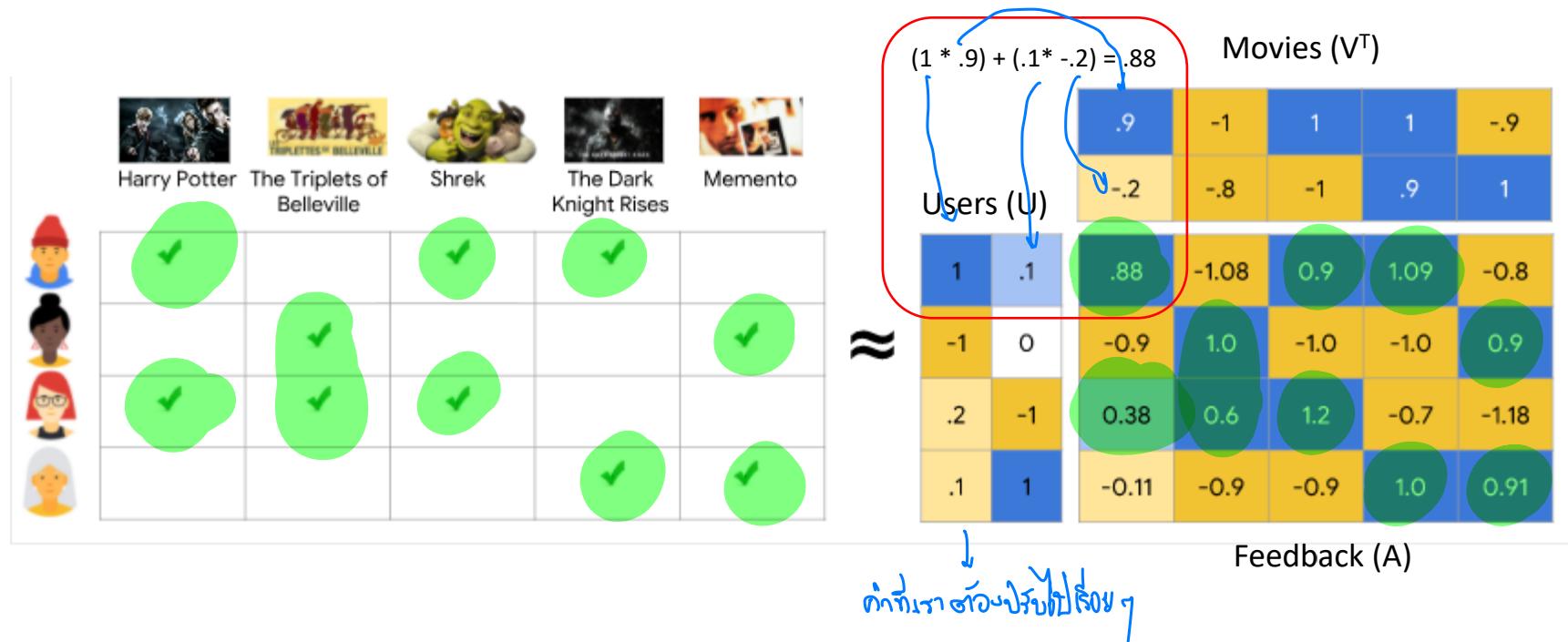
Example: Recommendation Systems based on Matrix Factorization

- *Movies* matrix represents characteristics of each movie.
- *Users* matrix represents preference on each characteristic.



Example: Recommendation Systems based on Matrix Factorization

- Every cell in the *rating* matrix can be calculated from the dot product of *user* and *movie* matrices.
- Matrix Factorization: Find U and V such that $A \approx UV^T$
- Iteratively adjust the value in *user* and *movie* matrices to minimize rating error (objective function).
- Complexity $\rightarrow O(\text{no.rating} * \text{no.features} * \text{no.iteration})$



HPCG Benchmark (used by Top500)

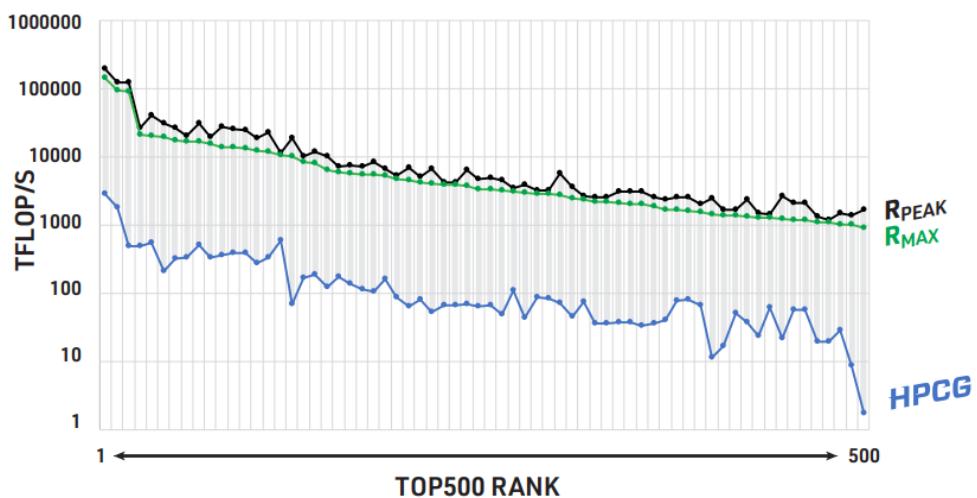
PRECONDITIONED CONJUGATE GRADIENT SOLVER

Code ວິທີ benchmark ໂອຍພາກສ່ວນຈະເຄືອຂັ້ນ
 $p_0 \leftarrow x_0, r_0 \leftarrow b - Ap_0$ dot product ເຊິ່ງ
for $i = 1, 2, \dots, \text{max_iterations}$ **do**

$z_i \leftarrow M^{-1}r_{i-1}$
if $i = 1$ **then**
 $p_i \leftarrow z_i$ dot product
 $\alpha_i \leftarrow \text{dot_prod}(r_{i-1}, z_i)$ ພົບປ່ຽນການດົກໂລຍ
else
 $\alpha_i \leftarrow \text{dot_prod}(r_{i-1}, z_i)$ ການດຳນວາຂອງ matrix
end if
 $\beta_i \leftarrow \alpha_i / \alpha_{i-1}$
 $p_i \leftarrow \beta_i p_{i-1} + z_i$
end if
 $\alpha_i \leftarrow \text{dot_prod}(r_{i-1}, z_i) / \text{dot_prod}(p_i, Ap_i)$
 $x_{i+1} \leftarrow x_i + \alpha_i p_i$
 $r_i \leftarrow r_{i-1} - \alpha_i A p_i$
if $\|r_i\|_2 < \text{tolerance}$ **then**
 STOP
end if

end for

- The benchmark is designed to model data access patterns of real-world applications such as sparse matrix calculations.
- It tests the limitations of memory system and interconnect better than TOP500's Linpack benchmark.



<http://hpcg-benchmark.org/>

<https://www.icl.utk.edu/files/print/2019/hpcg-sc19.pdf>

ສົດຄະພາກວ
SN ສັນບູກທາງທີ່ demand ເຂົາເກົ່າ
HW ເຮົາເຈັນ

Scalability & Parallelism

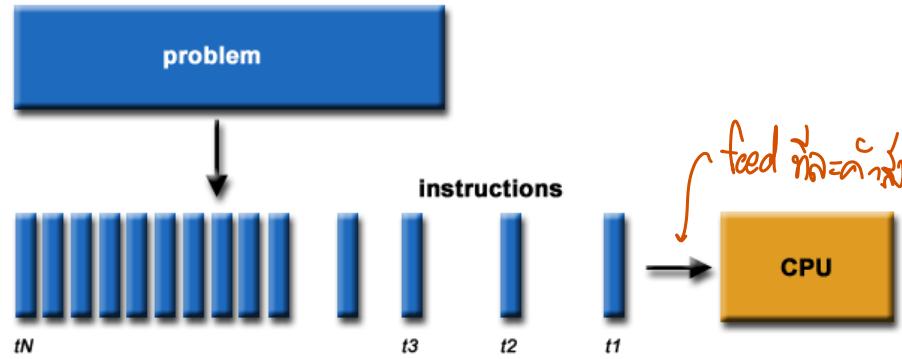
Parallelism

- Parallelism is the potential to divide a task into sub-tasks and perform them in parallel.
HN ห้ามท้อ ใจเด็ดเดี่ยว ใจเด็ดเดี่ยว ใจเด็ดเดี่ยว
SN ทำงานได้ เข้าคิวในนั้น แบบต่อเนื่องกัน
HW ห้ามขาดร่องรอย ลืกงาน
- Parallel computing involves how to exploit parallelism in software with parallelism in hardware (parallel architecture).
ด้วยเดียว 1 core เท่า SN ทำงาน Parallel ก็เท่า Sequential ฉะนั้น
ด้วย 64 core จะ SN ล้วน 1 thread

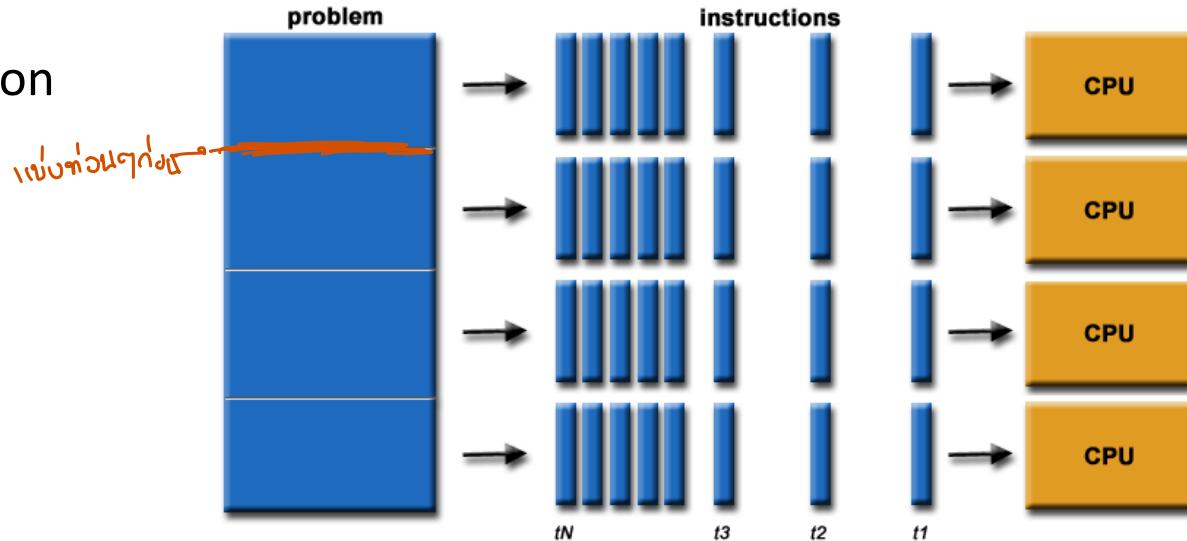
↑
การทำงาน Parallel SN ตัวเดียว HW แบบท่อน

Serial vs. Parallel Computation

Serial Computation



Parallel Computation



Sequential Computation

```
for (i=0; i<6; i++)  
    c[i] = a[i] + b[i];
```

- Primitive operations (i.e. load, add, store) are performed sequentially.

Diagram illustrating a parallel loop over 6 iterations (i=0 to i=5) across 3 time steps (t₀ to t₁₂). The array $a[]$ is updated by adding values from array $b[]$ at each time step, and the result is stored in array $c[]$.

The diagram shows the following sequence of operations:

- Iteration 0:** $a[0]$ is updated at t_0 by adding $b[0]$ and stored in $c[0]$.
- Iteration 1:** $a[1]$ is updated at t_1 by adding $b[1]$ and stored in $c[1]$.
- Iteration 2:** $a[2]$ is updated at t_2 by adding $b[2]$ and stored in $c[2]$.
- Iteration 3:** $a[3]$ is updated at t_3 by adding $b[3]$ and stored in $c[3]$.
- Iteration 4:** $a[4]$ is updated at t_4 by adding $b[4]$ and stored in $c[4]$.
- Iteration 5:** $a[5]$ is updated at t_5 by adding $b[5]$ and stored in $c[5]$.
- Iteration 6:** $a[6]$ is updated at t_6 by adding $b[6]$ and stored in $c[6]$.
- Iteration 7:** $a[7]$ is updated at t_7 by adding $b[7]$ and stored in $c[7]$.
- Iteration 8:** $a[8]$ is updated at t_8 by adding $b[8]$ and stored in $c[8]$.
- Iteration 9:** $a[9]$ is updated at t_9 by adding $b[9]$ and stored in $c[9]$.
- Iteration 10:** $a[10]$ is updated at t_{10} by adding $b[10]$ and stored in $c[10]$.
- Iteration 11:** $a[11]$ is updated at t_{11} by adding $b[11]$ and stored in $c[11]$.

A vertical dashed line at t_2 separates the first two time steps. The label "parallelism" is written in blue at the top right.

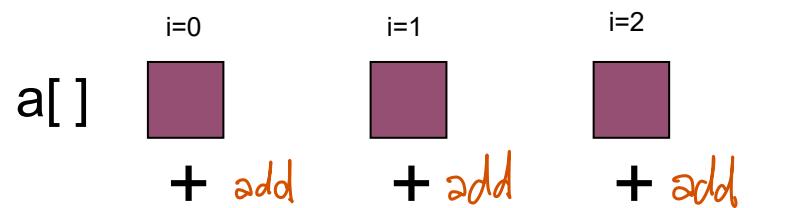
2 cores

Parallel computation with 2 processing elements

- A task is divided and performed in parallel. ແກ້ວມານິ້ນດັບດີ

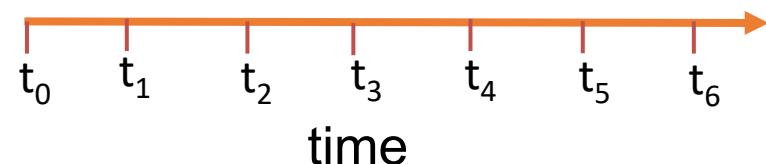
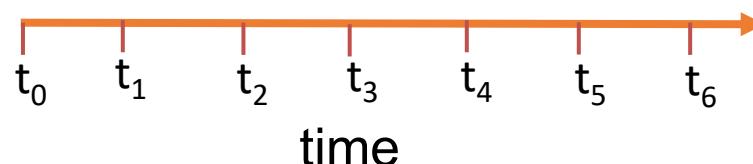
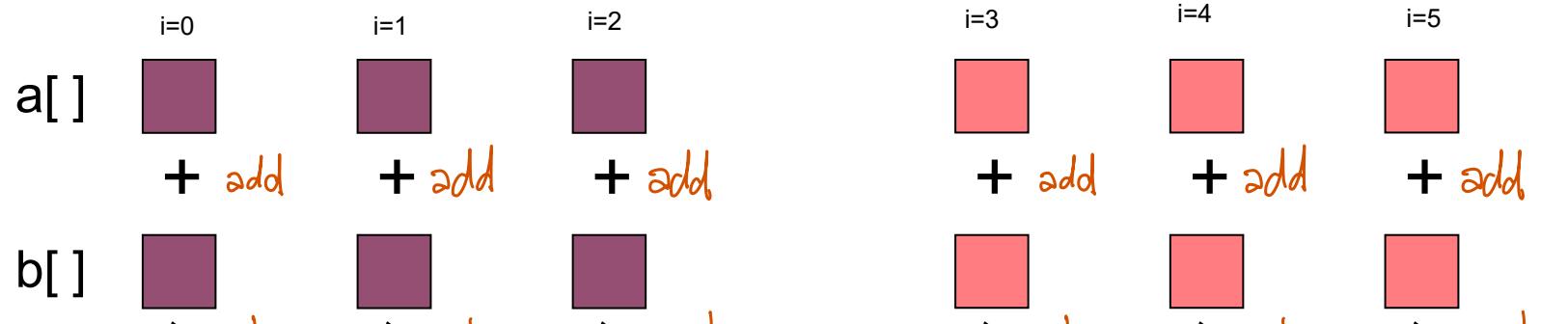
for ($i=0; i<3; i++$)

$$c[i] = a[i] + b[i];$$



for ($i=3; i<6; i++$)

$$c[i] = a[i] + b[i];$$



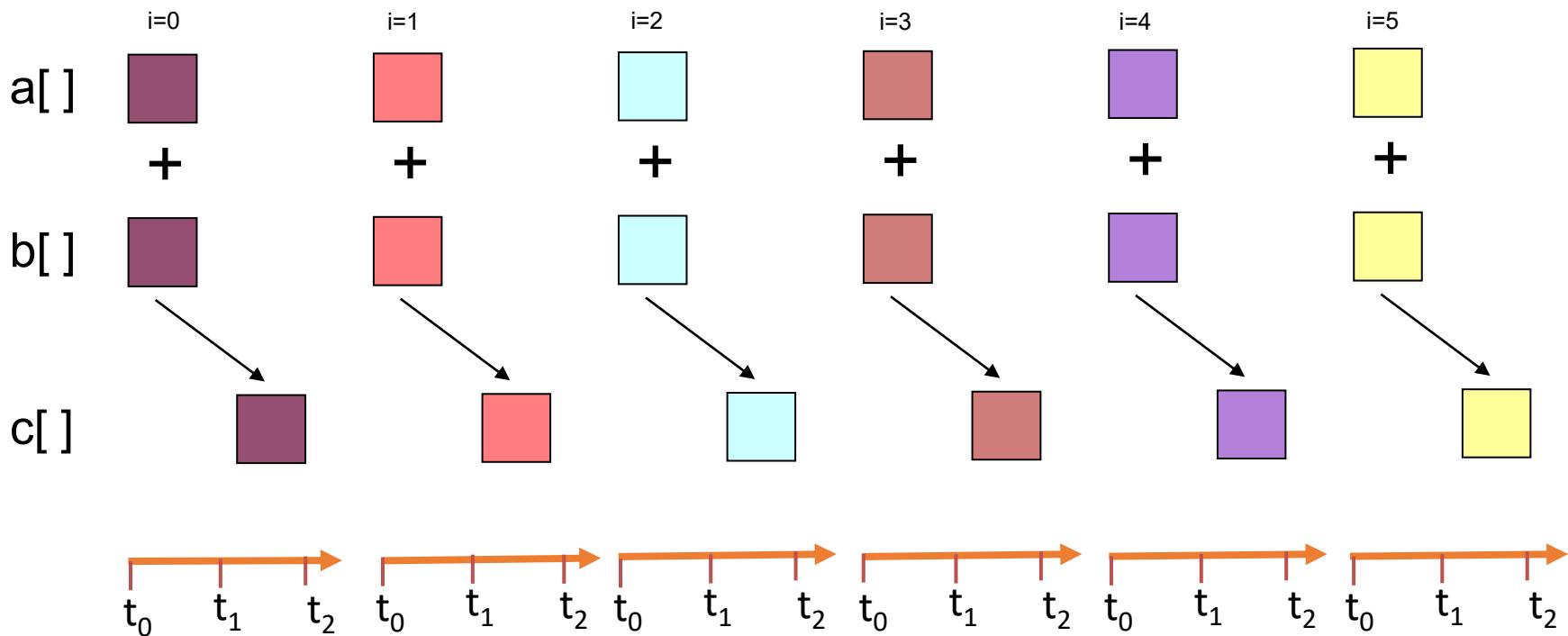
Parallel computation with 6 processing elements

6 cores

- Exploiting all parallelism in software.

¶ parallelism in یهیه

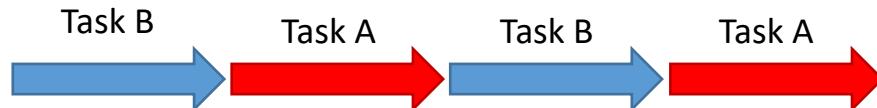
$$c[i] = a[i] + b[i];$$



Parallelism vs. Concurrency

- Concurrency is when multiple tasks logically (actually or appear to) run simultaneously. គំរាលការងារក្នុងក្រឡាយបែងចែង។
- Parallelism is when multiple tasks physically (actually) run simultaneously.
- Concurrency is more general than parallelism.
 - Parallel processes always run concurrently, but concurrent processes do not always run in parallel.

Concurrency but
not parallelism



Concurrency and
parallelism



Scalability

- **Scalability** is the ability of a system or process to handle growing amount of work or its ability to be enlarged to accommodate that growth.

ຄວາມສໍາເລັດໃນການຕີບໂຫຼວງຮັບ ປິ່ນການກົດ ອີ້ນ ສອງດອງຮັບໄດ້ ?

- Scalability of a parallel system is relative to *the ability to exploit parallelism*.
- **Scalable system**
 - The performance improves after adding hardware, proportionally to the added capacity.
 - Support wide range of system sizes
 - Use parallel computing or distributed systems techniques

Levels of parallelism in software

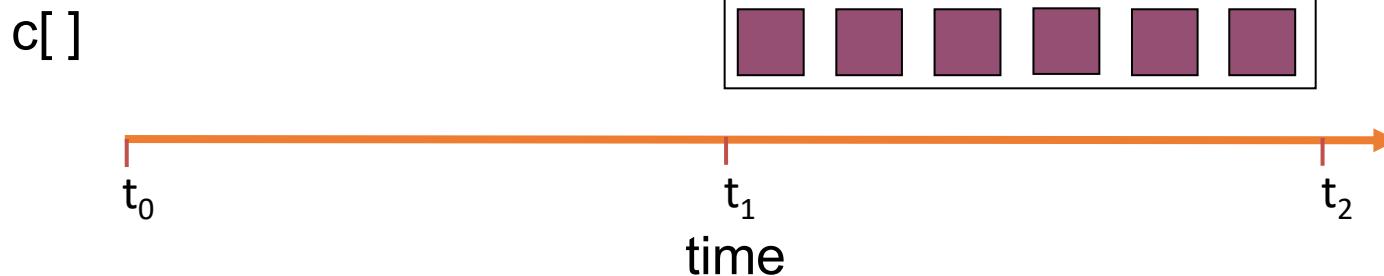
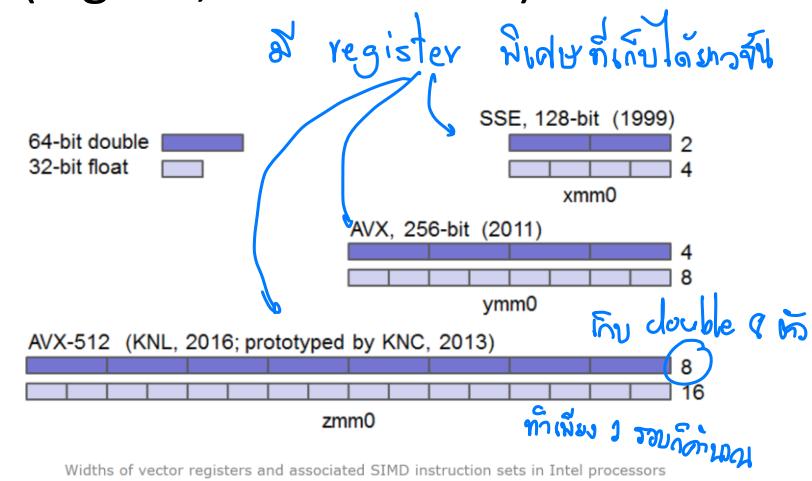
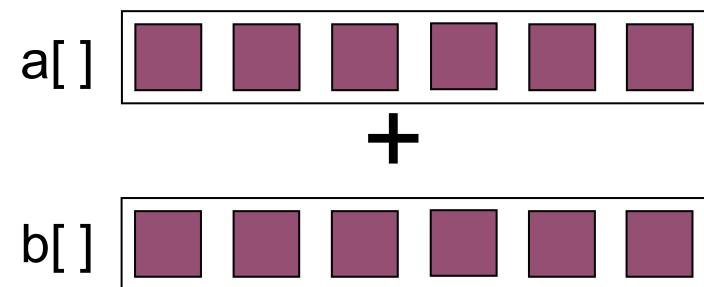
ms match parallel 9_{th}
SW กับ HW

- Instruction-level parallelism $c[ij] = a[ij] + b[ij] + c[ij] + d[ij]$
 - Executing multiple CPU instructions in parallel
- Loop-level parallelism
 - Executing individual iterations of a loop in parallel
- Thread-level
 - Executing multiple threads in parallel
- Process-level $\text{Run 9th hadoop cluster}$
 - Executing multiple processes in parallel

运行第9个 hadoop cluster 通过将任务分派给不同的进程

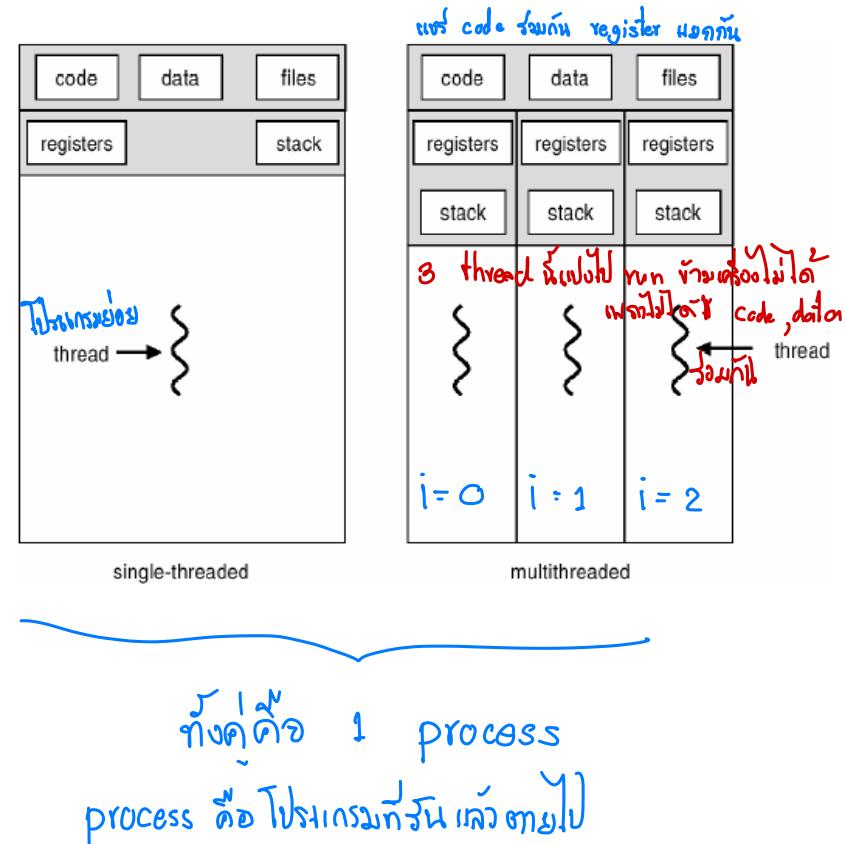
Vector Processing

- Each vector instruction performs on all elements in vector operands at once.
- A vector register has limited size (e.g. 32, 64 elements)

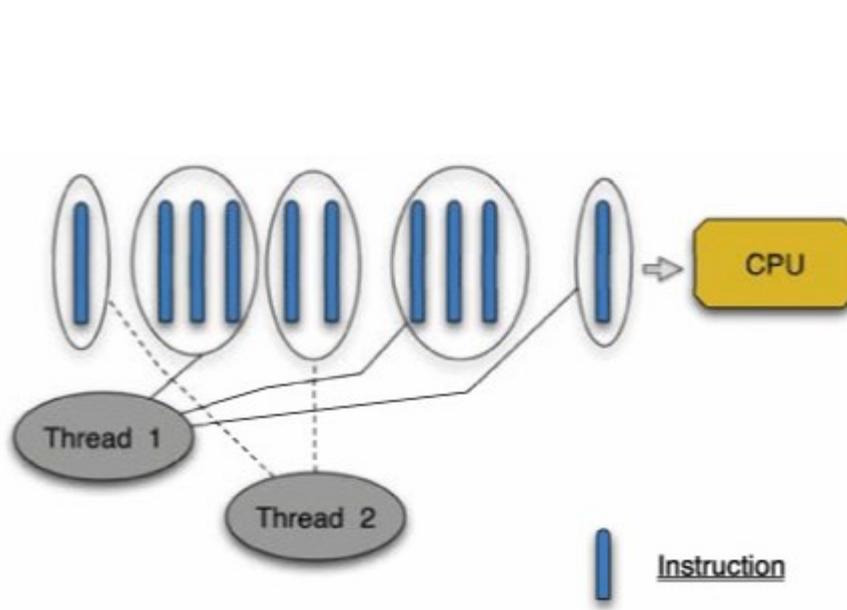


Multi-thread Processing Model

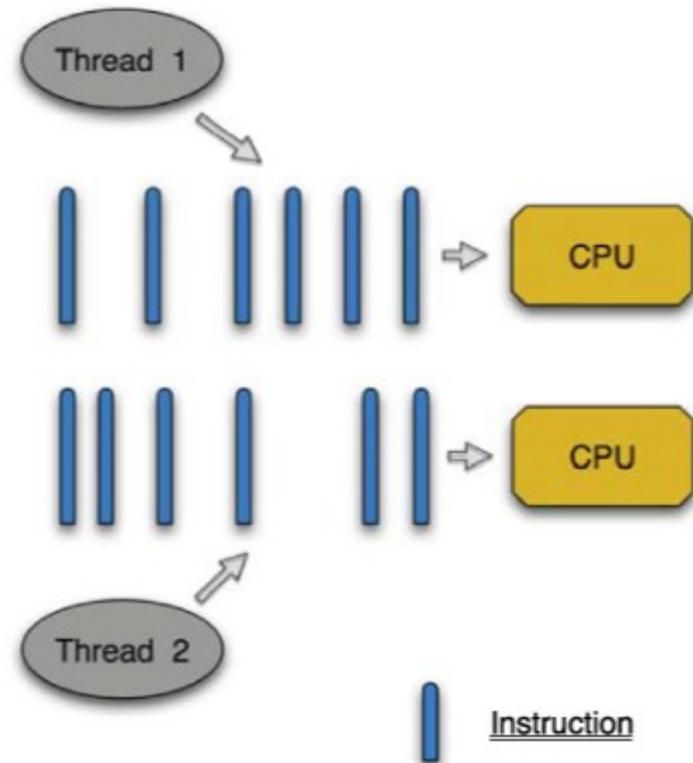
- A process is a running instance of a program
- A process can have multiple threads.
- Processes have separate memory spaces.
- Threads in the same program share memory space.



Threads

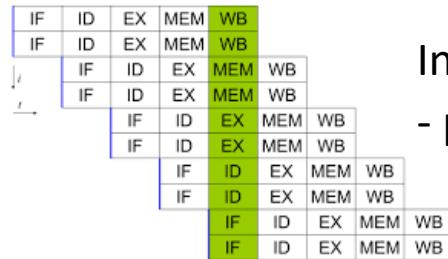


ໄຟຟ້າເຮັດວຽກ
Multiple threads running
concurrently, but not in parallel,
on one CPU core
(no performance gain)



ນີ້ CPU ດລາຍຈຳ ຕົກລາງ parallel ລື່ອງ
Multiple threads running in
parallel on separate CPU cores

Level of parallelism in hardware



งาน data science มักจะ run ที่ high-level language (library)

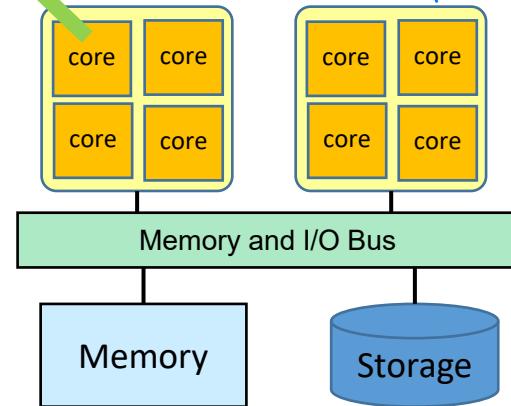
มีจุดหมายรวมกัน ILP ทำให้ต้องมีตัว เรากಡ่ตัว lib ให้ดีๆ และ optimize

Instruction Level Parallelism (ILP) (อยู่ใน หน้า)

- pipelining, superscalar, vector processing

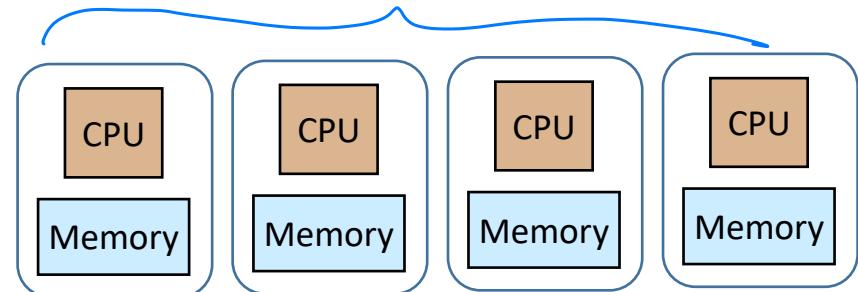
ล็อกอิน run

ทำงานตัวสั่งพร้อมกัน
 executes instruction 同時に実行



Multicore/Multiprocessor

Multi Processor

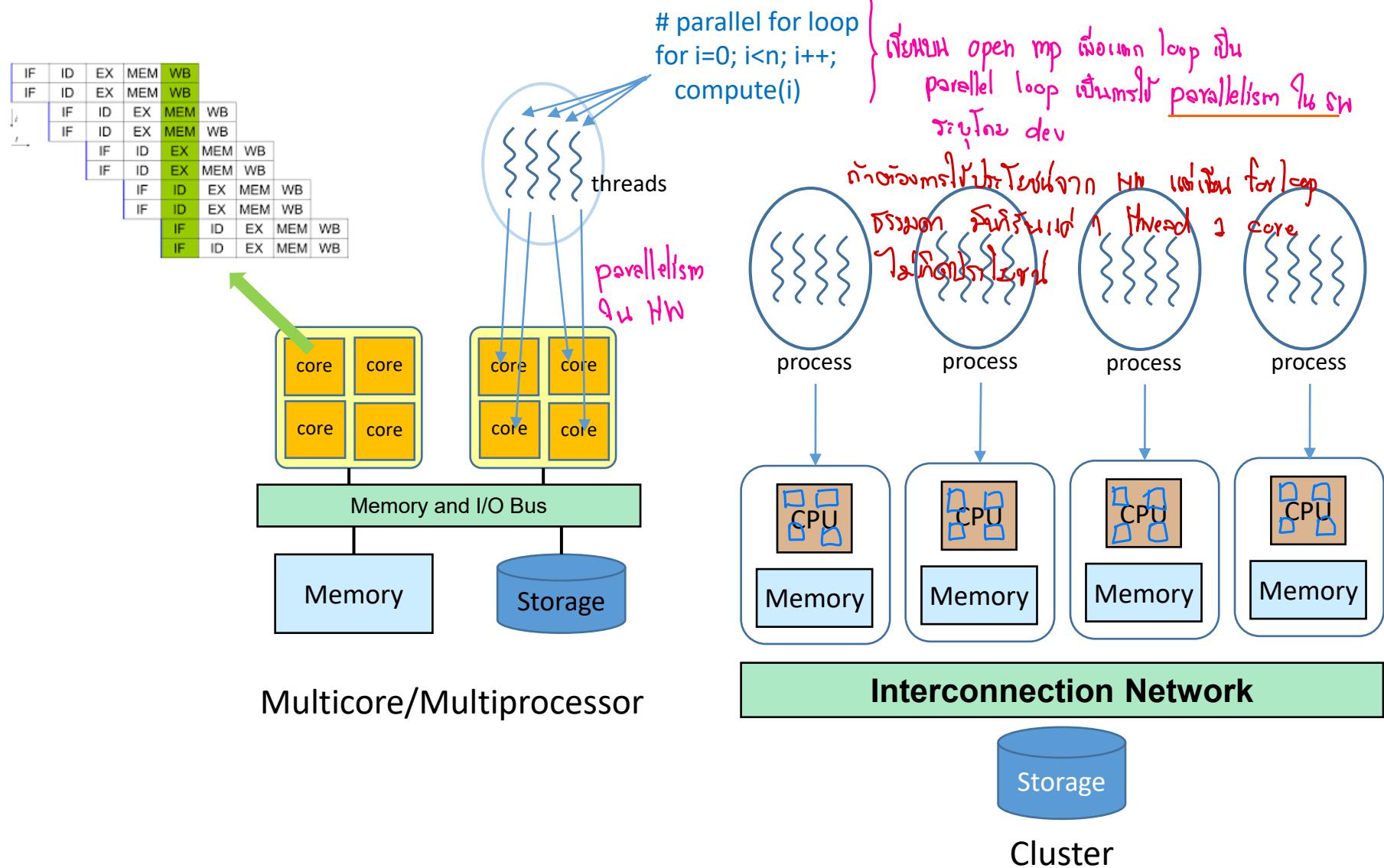


Interconnection Network



Cluster

Exploiting parallelism in software with parallel architecture



Data-Parallelism and Multithreading

Thread#0	Thread#1	Thread#2	Thread#3	Thread#4	Thread#5	Thread#6	Thread#7	Thread#n
a_0	a_1	a_2	a_3	a_4	a_5	a_6	a_7	\dots	\dots	a_n
+	+	+	+	+	+	+	+	+	+	+
b_0	b_1	b_2	b_3	b_4	b_5	b_6	b_7	\dots	\dots	b_n
=	=	=	=	=	=	=	=	=	=	=
c_0	c_1	c_2	c_3	c_4	c_5	c_6	c_7	\dots	\dots	c_n

Exploiting parallelism in software with parallel architecture

□ Parallelism in software

- Instruction-level
- Loop-level
- Thread-level
- Process-level

mapping ไฟล์มูลค่าที่ต้องการ HN

• Parallel architecture

- Single-core (superscalar, vector processing)
- GPU เป็น multi-core มากกว่า
- Multicore/Multiprocessor
1 CPU หลาย core 1 board หลาย processor
กับ 2 เทศบาราชุดในเครือข่ายเดียวกัน

Multicomputer (cluster)

ไม่ต้องการรีเซ็ต iteration 1 รีเซ็ตใหม่
i=2 รีเซ็ตตัวหนึ่งไว้ต่อ HN ไม่ว่าจะการคำนวณแบบ

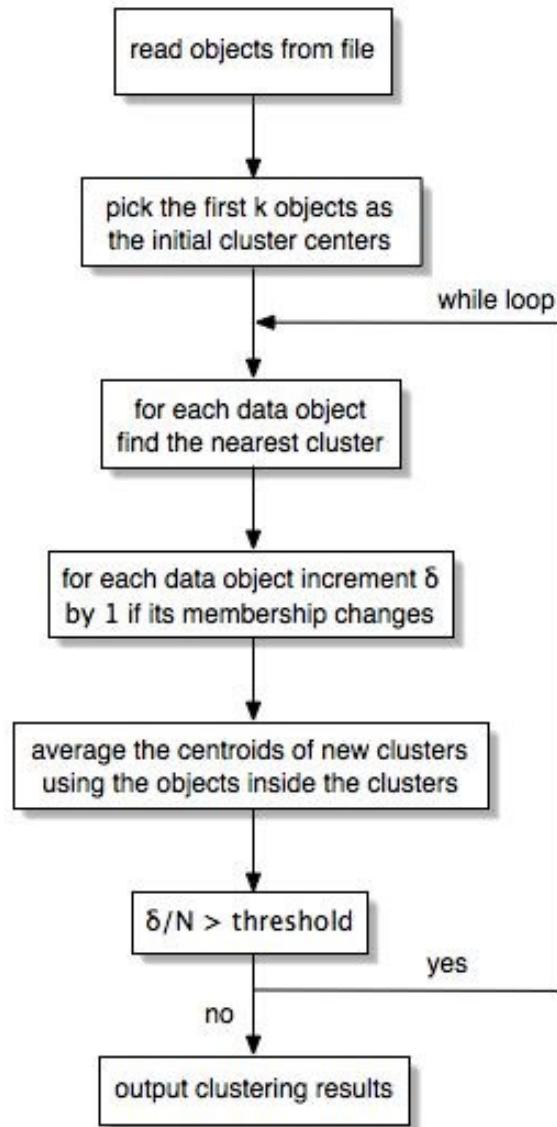
lib หรือ compiler จะปลดการจัดแบ่ง thread แล้วนำไป GPU , Multicore

หรือ vector processing

รันห้องเครื่องในเครือข่าย multi-thread ต้องการรีเซ็ต

1 process แล้ว process ใหม่ thread ใหม่ แล้วรีเซ็ตเครื่องเดิม
concept ของ process ต้องการรีเซ็ตเครื่องเดิม

Sequential K-Means Clustering



N : number of data objects

K : number of clusters

$\text{objects}[N]$: array of data objects

$\text{clusters}[K]$: array of cluster centers

$\text{membership}[N]$: array of object memberships

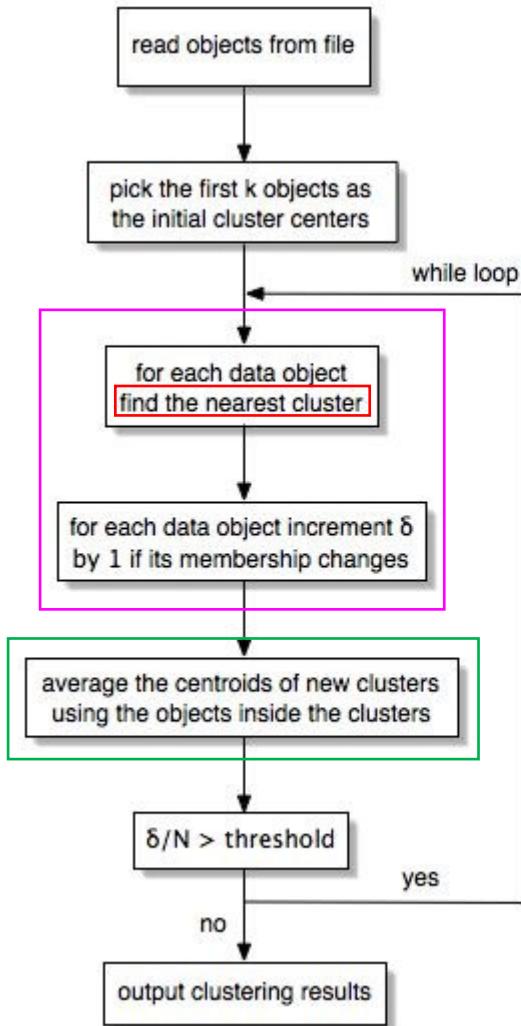
kmeans_clustering()

```
1 while δ/N > threshold
2   δ ← 0
3   for i ← 0 to N-1          loop over each object
4     for j ← 0 to K-1        loop over each cluster
5       distance ← | objects[i] - clusters[j] |
6       if distance < dmin
7         dmin ← distance
8         n ← j
9       if membership[i] ≠ n
10         δ ← δ + 1
11         membership[i] ← n
12         new_clusters[n] ← new_clusters[n] + objects[i]
13         new_cluster_size[n] ← new_cluster_size[n] + 1
14     for j ← 0 to K-1        calculate centroid of new cluster
15       clusters[j][*] ← new_clusters[j][*] / new_cluster_size[j]
16       new_clusters[j][*] ← 0
17       new_cluster_size[j] ← 0
```

Annotations in blue ink:

- loop over each object
- loop over each cluster
- calculate centroid of new cluster

Parallel K-Means Clustering



N : number of data objects

K : number of clusters

$\text{objects}[N]$: array of data objects

$\text{clusters}[K]$: array of cluster centers

$\text{membership}[N]$: array of object memberships

in loop เนื้อหาเรื่อง parallelism

ว่าทำ parallel ยังไง ?

kmeans_clustering()

```

1 while δ/N > threshold
2   δ ← 0
3   for i ← 0 to N-1
4     for j ← 0 to K-1
5       distance ← |objects[i] - clusters[j]|
6       if distance < dmin
7         dmin ← distance
8         n ← j
9     if membership[i] ≠ n
10      δ ← δ + 1
11      membership[i] ← n
12      new_clusters[n] ← new_clusters[n] + objects[i]
13      new_cluster_size[n] ← new_cluster_size[n] + 1
14   for j ← 0 to K-1
15     clusters[j][*] ← new_clusters[j][*] / new_cluster_size[j]
16     new_clusters[j][*] ← 0
17     new_cluster_size[j] ← 0
  
```

Parallelizable loops

2 ขั้นตอนหลักๆ

分け centroid ออก
แล้วคำนวณ distance ของทุกๆ point ที่อยู่ใน cluster นั้นๆ

Parallelism Issues

Perfectly Parallel Problem

ໂປ່ນທຽມຫຼັກພາກສະນູງ, parallel ໄດ້ນຳງ

- Also known as “embarrassingly parallel problem”
- A problem that is easy to divide into parallel subtasks without dependency or communication among subtasks. ເພີ້ນ ການບາກດຳໄຟລ໌ ລາຍລອດ
- Unfortunately, most parallel problems are not in this class.

Issues in Parallelism (1)

- Dependency
 - A task can be executed only after all tasks that it depends on have already been executed.
 - Dependency limits parallelism.

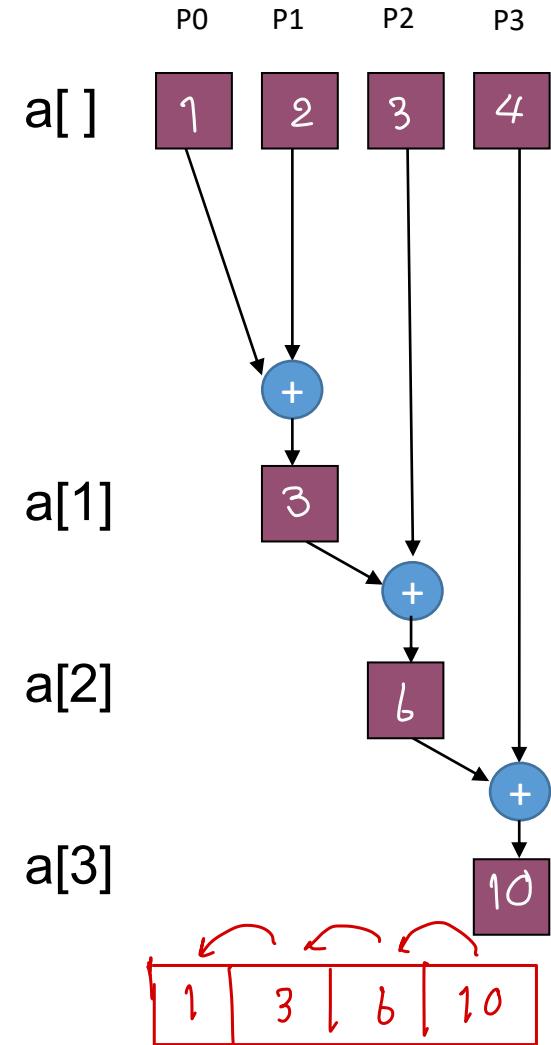
การที่ บางขั้นตอนอาจต้อง ผลลัพธ์จากขั้นตอนหนึ่ง จึงใช้ได้ใน parallelism

Example: Prefix sum

```
for (i=1; i<n; i++)  
    a[i] = a[i] + a[i-1];
```

dependencies

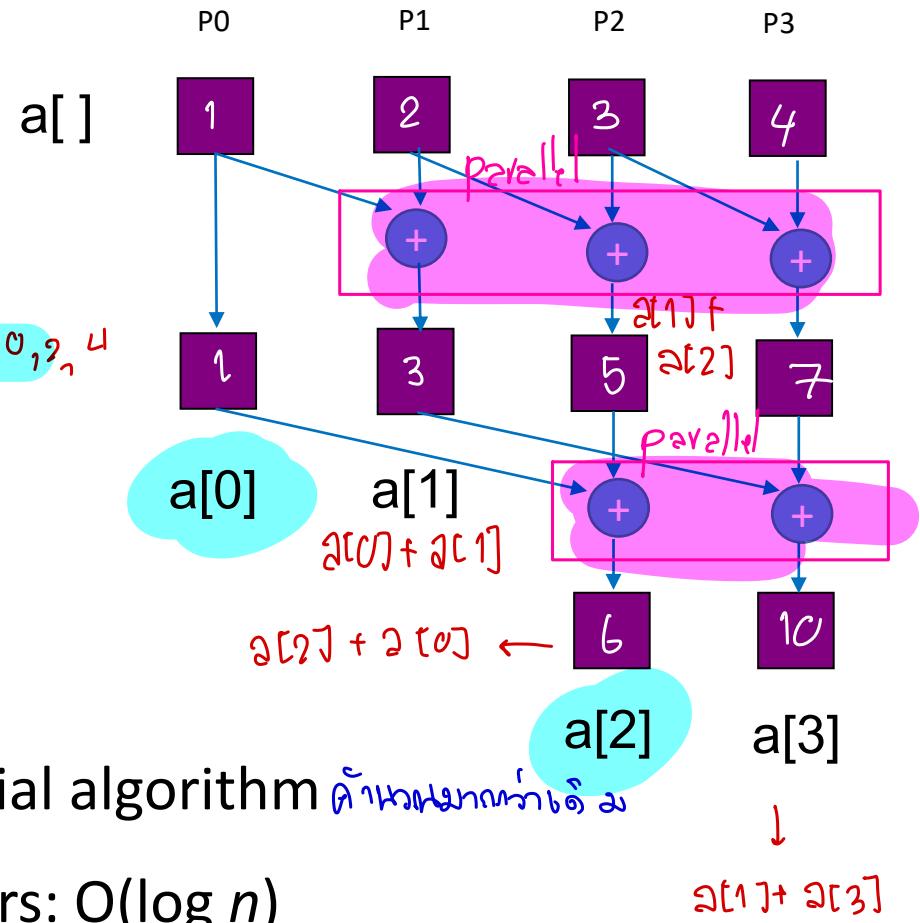
- $a[i]$ depends on $a[i-1]$
- Subtasks cannot run in parallel.
- Time complexity: $O(n)$



Parallel Prefix Sum Algorithm

j : 3rd phase ຖ້ານີ້ມາ
for (j=0; j < log₂(n); j++)
parallel for (i=2^j; i < n; i++)
a[i] = a[i] + a[i-2^j];

ເທົ່ານີ້ມາ phase ພະນັກງານ



- More computation than sequential algorithm ດັ່ງນັ້ນມີຄວາມກຳສົດ ຂັ້ນ
- Time complexity with n processors: $O(\log n)$

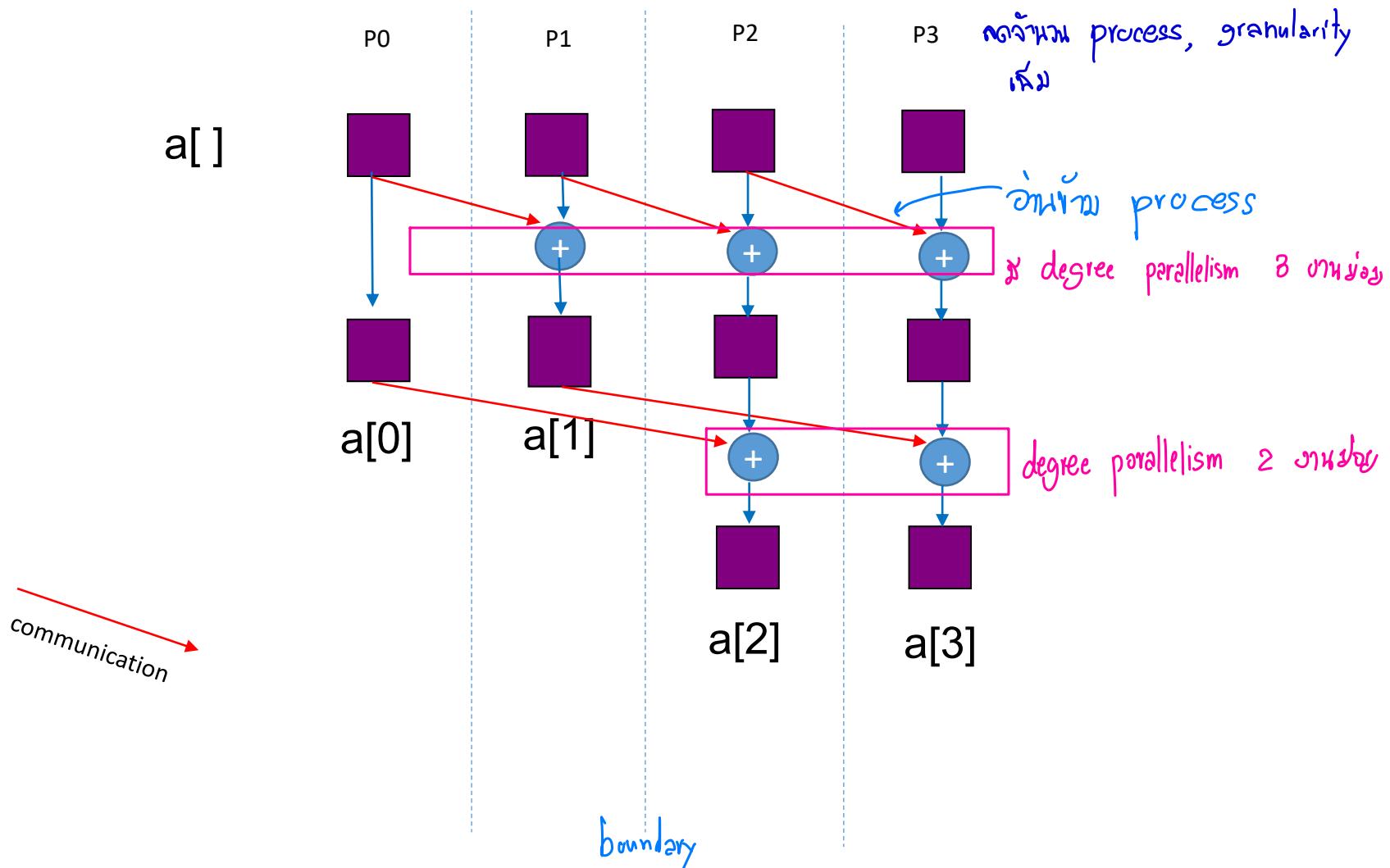
Issues in Parallelism (2)

- Communication and synchronization

- Needs for transferring information among subtasks
- Communication and synchronization cause delay.
*ก่อให้เกิดช้าลง delay
ต้องส่งข้อมูลระหว่างกัน นาน*
- Communication method depends on level of parallelism
 - Threads communicate via shared memory *thread level สื่อสารผ่านหน้าจอ memory ร่วมกัน
กับหน้าจอ*
 - Processes on different machines communicate via messages or remote procedure call (RPC) *สื่อสารผ่านเครือข่ายท้องที่*

Communication

ໃຈ data ໃນກ່າວ process



Issues in Parallelism (3)

- Degree of parallelism

- The number of parallel tasks
- Degree of parallelism in software is usually greater than degree of parallelism in hardware.

ແມ່ນ parallelism ຂັ້ນໄດ້ໄວ້ຕີ ສົນກົວຈະ processing unit ແຕ່ລົມ

ໜຸ່ງປູກໄດ້ໄວ້ ເພື່ອໃຫ້ CPU core ໄດ້ຖຸດາ

- Granularity

- sizes of sub-tasks
- Different levels of parallelism deal with different granularity.

$$\text{degree of parallelism} \propto \frac{1}{\text{Granularity}}$$

Issues in Parallelism (4)

- Other overheads
 - Creating and destroying subtasks / managing worker pool
 - Context switching between tasks that share the same computing element
 - Overheads may relate to degree of parallelism and granularity.

overhead მასში თერთ განვითარებულ