

# システムプログラム概論

## 並列コンピュータシステムとそのOS

2007/4/27

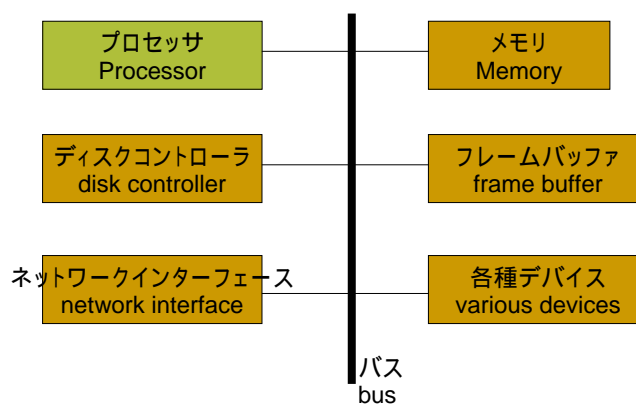
門林 雄基

奈良先端科学技術大学院大学

1

# ユニプロセッサシステム

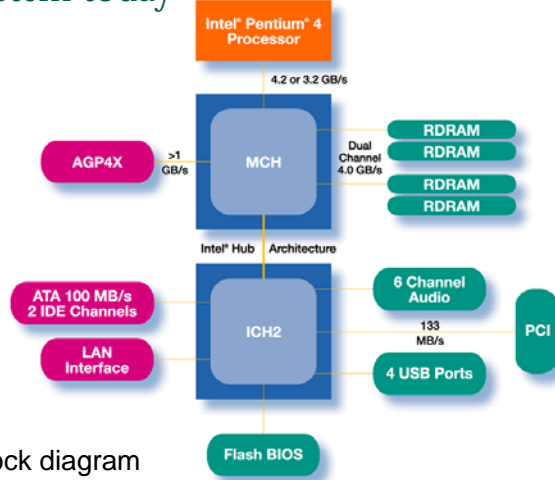
## Uni-processor system revisited



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

2

## 今日のユニプロセッサシステム Uni-processor system today



- Intel i850 chipset block diagram
- Source: intel web site

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

3

## ユニプロセッサシステムの処理モデル Processing model of uni-processor systems



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

4

## より高速な処理を目指して - 並列性の導入 Making systems faster: parallelism

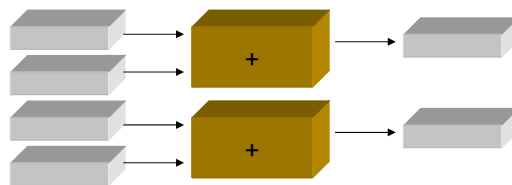
- ビット単位での並列性  
Bit-level parallelism
  - 8 bit, 16 bit, 32 bit, 64 bit...
- 命令単位での並列性  
Instruction-level parallelism
  - Pipelining, Superscalar...  
(本講義の範疇外)
- プロセス・スレッド単位での並列性  
Process-level / Thread-level parallelism

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

5

## 並列システムの処理モデル(1) Processing model of parallel systems (1)

- SIMD (Single Instruction, Multiple Data)
- SIMD today
  - プロセッサ内での実現
  - AltiVec (PowerPC G5), SSE2 (x86) (example)

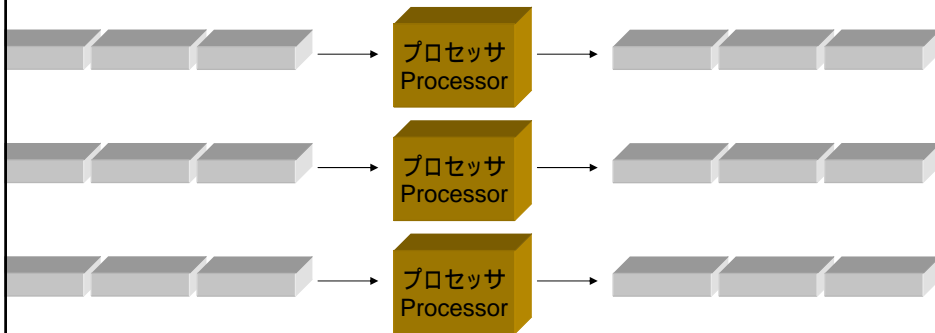


Copyright(C)2007 Youki Kadobayashi. All rights reserved.

6

## 並列システムの処理モデル(2) Processing model of parallel systems (2)

### ■ MIMD (Multiple Instruction, Multiple Data)



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

7

## MIMD systems: part I

- Centralized memory
  - Symmetric Multiprocessor (SMP)
  - Uniform Memory Access
- Distributed memory
  - Non-Uniform Memory Access (NUMA)
  - today: cc-NUMA

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

8

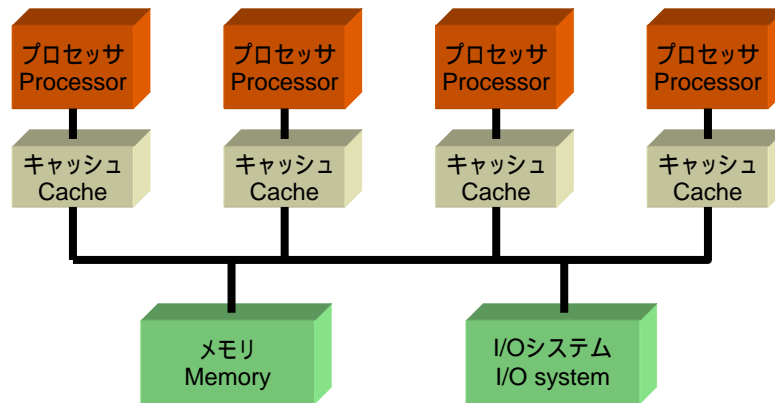
## Dual-processor systems

- PowerMAC G5
- Opteron

## Dual-core systems

- UltraSPARC IV

## Textbook SMP systems



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

11

## SMP today

- Sun V880
- rx8620

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

12

## プログラミング・モデル Programming models

- マルチプログラミング  
Multiprogramming
- 共有メモリ  
Shared memory
  
- メッセージ・パッシング  
Message passing
- データ・パラレル  
Data parallel

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

13

## プログラミング・モデルとOS Programming models and OS

- multiprocessor operating systems
  - Multiprogramming
  - Shared memory
  
- uniprocessor operating systems + library
  - Message passing
  - Data parallel

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

14

## 共有メモリ型システムにおけるプログラミングモデル(1) Programming model in shared memory systems: multiprogramming

プロセスでの並列性 (process-level parallelism)

- メモリ集中型、分散型いづれにも適用可能  
applicable to both NUMA and SMP
- C プログラムから: `fork`, `exec`, `waitpid`, etc.
- シェルから:  

```
#!/bin/bash
for ((i = 0; i < 20; ++i)); do
  ./param-search $i
done
```
- パラメータ探索問題 (parameter-search problems) に向く
- タスク並列型の負荷(task-parallel workloads) に向く

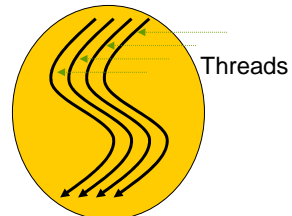
Copyright(C)2007 Youki Kadobayashi. All rights reserved.

15

## 共有メモリ型システムにおけるプログラミングモデル(2) Programming model in shared memory systems: threads

スレッドレベルでの並列性 (thread-level parallelism)

- 共有メモリが必要  
prerequisite: shared memory
  - SMP, cc-NUMA
- 例1) posix threads
  - [example](#)



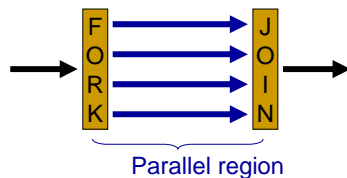
Copyright(C)2007 Youki Kadobayashi. All rights reserved.

16



## スレッドレベルの並列性 (例2): OpenMP Thread-level parallelism by example: 2/2

- OpenMP
- fork-join model
- parallel region
- Example
- [www.openmp.org](http://www.openmp.org)



```
#include <omp.h>
#define N 1000

main()
{
  int i;
  float a[N], b[N], c[N];
  // ...
  #pragma omp parallel shared(a,b,c) private(i)
  {
    #pragma omp for schedule(dynamic,100)
    for (i = 0; i < N; i++)
      c[i] = a[i] + b[i];
  }
}
```

共有変数 (a,b,c)    スレッド固有 (i)

分割サイズ (100)

Parallel region

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

17

## 共有メモリ型システムにおけるプログラミングモデル(3) Programming model in shared memory systems: IPC

- プロセス間通信: inter-process communication
- semaphore, shared memory, message passing, ...
- Cプログラムから:
  - `sem_init`, `sem_wait`, `sem_destroy`
  - `shmget`, `shmat`, `shmdt`
  - `send`, `recv`

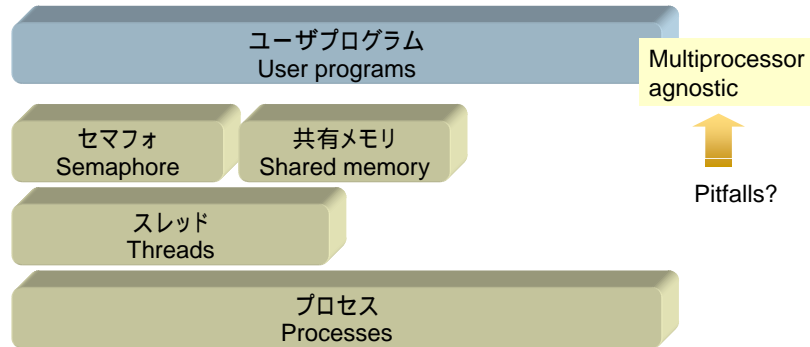
Copyright(C)2007 Youki Kadobayashi. All rights reserved.

18

## マルチプロセッサシステムの抽象化

Abstracting multiprocessor systems

- Multiprogramming, shared memory:
- 同一プログラム、マルチプロセッサ・システムへの可搬性  
same program; scales up to multiprocessor systems



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

19

## 共有メモリシステムにおけるOSの役割

Operating system role in shared-memory systems

- プロセッサの割り当て
  - スレッドへ / プロセスへ
- 通信
  - semaphore, shared memory, message passing
- ライブラリとコンパイラの役割が増大
  - Instruction-level parallelism – optimizing compilers
  - Thread-level parallelism – OpenMP

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

20

## MIMD systems: part II

### ■ Shared memory systems

- Centralized memory Thread-level parallelism with shared state
  - Symmetric Multiprocessor (SMP)
  - Uniform Memory Access
- Distributed memory Process-level parallelism with shared state
  - Non-Uniform Memory Access (NUMA)
  - today: cc-NUMA

### ■ Multicomputers, Clusters (Shared-nothing\*)

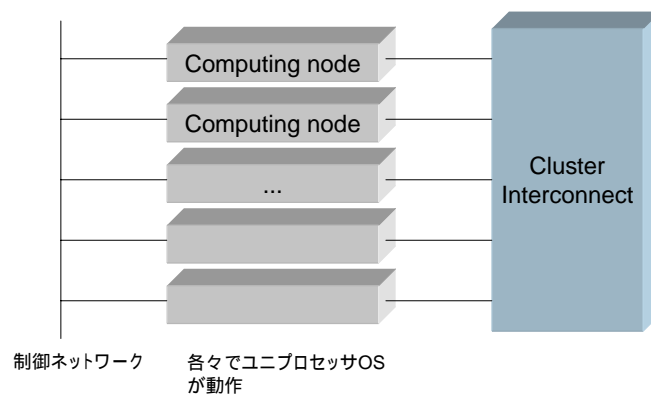
- Research MPPs (dead) Process-level parallelism without shared state\*
- Computing cluster
- Web server cluster, etc.

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

21

## 計算機クラスタ



### Computing cluster



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

22

## 並列計算機 vs 計算機クラスター: 隠れたコスト Hidden cost of computing clusters

- |  |  |  |
|--|--|--|
| <ul style="list-style-type: none"><li>■ SMP, cc-NUMA<ul style="list-style-type: none"><li>□ インターコネクトが高価</li><li>□ メモリアクセス速度の問題</li><li>□ Diminishing returns (IPP p.1)</li><li>□ 独自(商用)OS – 限られた選択肢</li></ul></li><br/><li>■ Cluster<ul style="list-style-type: none"><li>□ インターコネクトは安価</li><li>□ メモリアクセス速度は一定</li><li>□ 高速・安価な市販PCで構成</li><li>□ OS選択の自由</li></ul></li></ul> | <br><br><br><br><br><br><br><br><br><br> | <p><b>Pros</b></p> <ul style="list-style-type: none"><li>Reproducible results</li><li>Reusable software</li><li>Bigger memory</li><li>Better compiler</li><li>Larger MTBF</li><li>Reliable support</li></ul><br><p><b>Cons</b></p> <ul style="list-style-type: none"><li>Reproducibility issues</li><li>Explicit parallelism</li><li>Smaller memory</li><li>Average compiler</li><li>Smaller MTBF</li><li>Unreliable support</li></ul> |
|--|--|--|

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

23

## プログラミング・モデルとOS (再掲) Programming models and OS revisited

- multiprocessor operating systems
  - Multiprogramming
  - Shared memory
  
- uniprocessor operating systems + library
  - Message passing
  - Data parallel

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

24

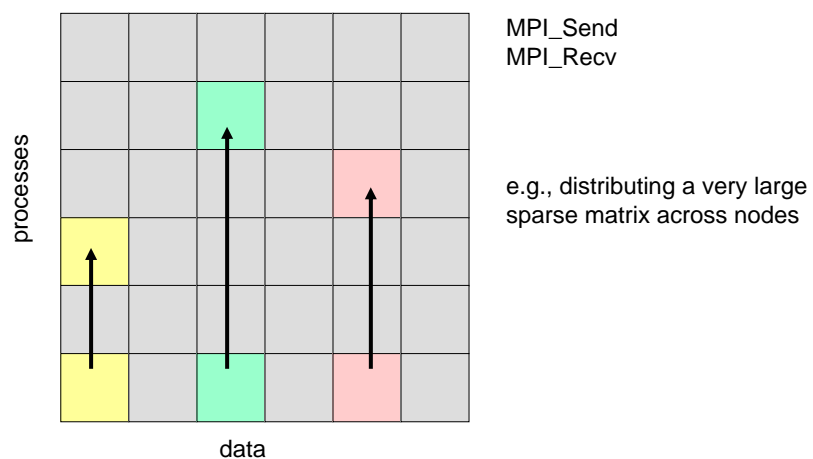
## Message passing programming model

- MPI (Message Passing Interface)
- Send / receive
- Communication within groups
- Scatter / gather
- Reduce
- Barrier synchronization

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

25

## Send/Receive in MPI



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

26

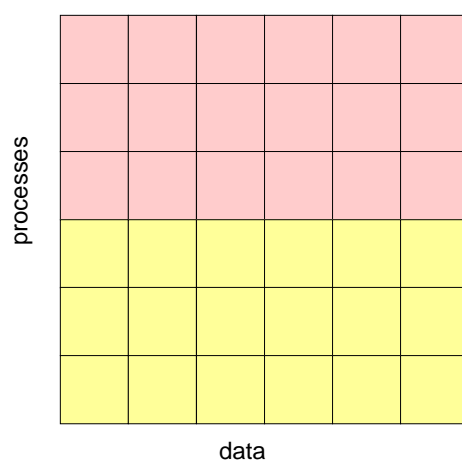
## Blocking send/receive, non-blocking send/receive

- Blocking
  - MPI\_Send
  - MPI\_Recv
  
- Non-blocking
  - MPI\_Isend
  - MPI\_Irecv

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

27

## Communication within groups



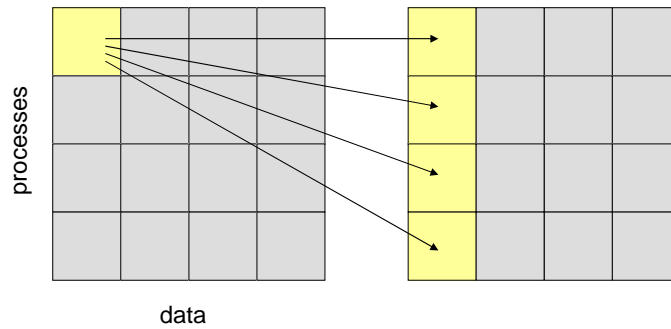
MPI\_Group\_\*

e.g., running the same algorithm  
on two data-sets in parallel

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

28

## Collective communications(1): Broadcast



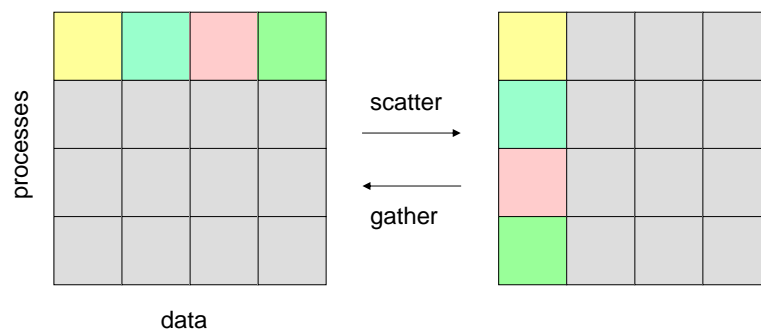
MPI\_Bcast

e.g., initial parameter

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

29

## Collective communications(2): Scatter / gather



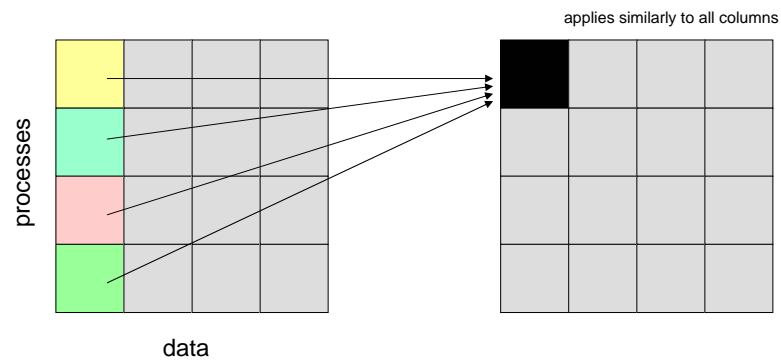
MPI\_Scatter  
MPI\_Gather

e.g., distribute vector across nodes  
matrix multiplication example ([uniprocessor](#), [MPI](#))

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

30

## Global reduction operations: Reduce



MPI\_Reduce  
MPI\_SUM, MPI\_PROD, MPI\_MAX, MPI\_MIN, ...

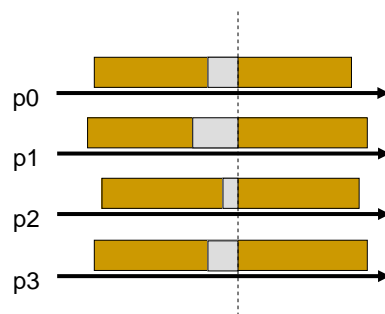
e.g., vector-matrix product

Copyright(C)2007 Youki Kadobayashi. All rights reserved.

31

## Barrier synchronization

- MPI\_Barrier
- Blocks until all group members have called MPI\_Barrier



Copyright(C)2007 Youki Kadobayashi. All rights reserved.

32



## まとめ

- Taxonomy of:
  - Parallelism
  - Multi-processor systems
- Modern systems:
  - Parallelism everywhere
    - system-level: MIMD
    - instruction-level: SIMD
  - Parallelism-friendly
    - programming models
    - libraries
    - compilers