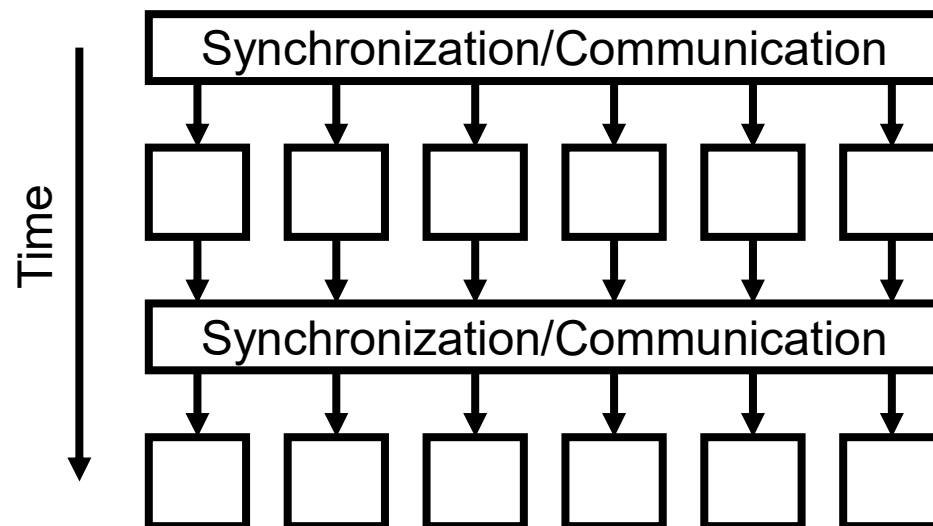# Class Schedule

➢ Review Nov. 7, 2018

➢ Nov. 12, Regular class, Exam Monday, Nov. 14, 2018

➢ The week of Nov. 18, no class, individual prepare term project

➢ Term project presentation: Monday Nov. 25 during and after class hour., Tuesday, Nov. 26 (back up)

➢ Term report due: Nov. 26
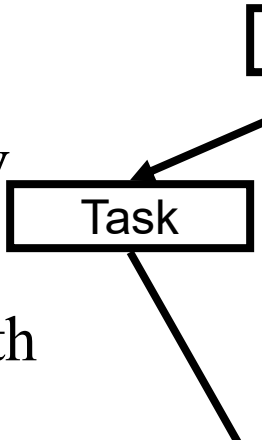
X. Sun (IIT)

**Xian-He Sun**

# Application Structure

➢ Frequently used patterns for parallel applications:

 o Single Program Multiple Data – SPMD (Domain Decomposition)

 o Embarrassingly Parallel

 o Master / Slave

 o Work Pool

 o Divide and Conquer
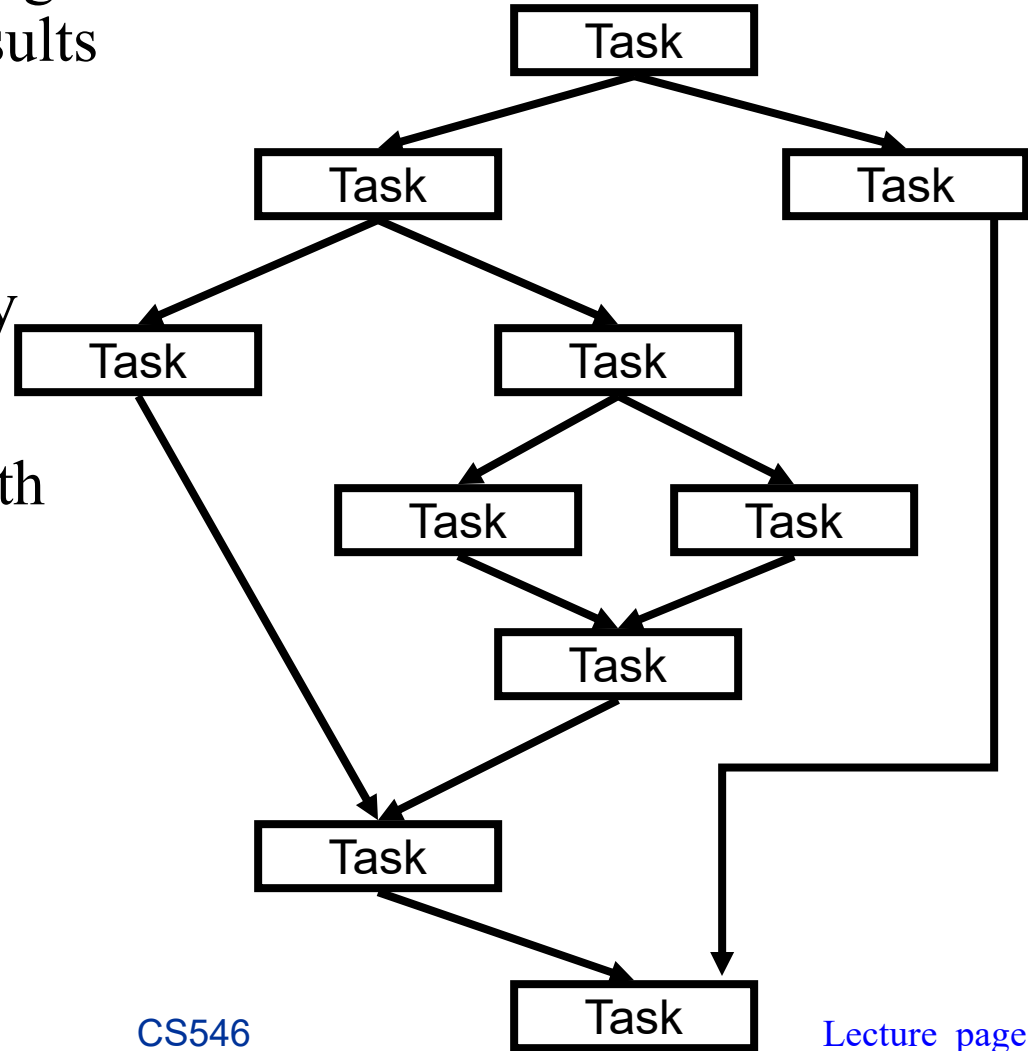
 o Pipeline

 o Competition

# Structure: Single Program Multiple Data

➢ Single program is executed in a replicated fashion.

➢ Processes or threads execute same operations on different data.

➢ Loosely-synchronous: Sequence of phases of computation and communication/synchronization.
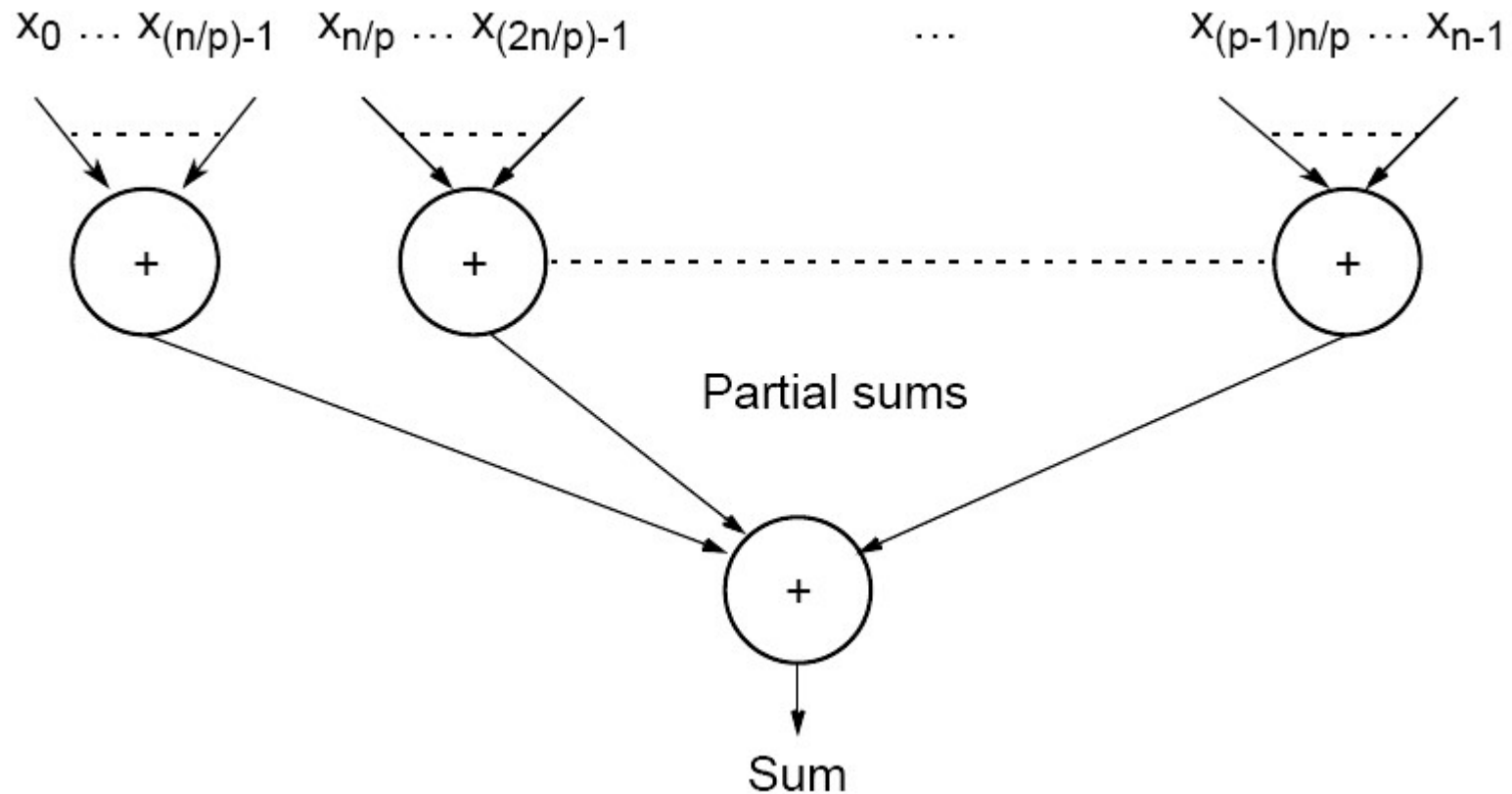
# Structure: Divide and Conquer

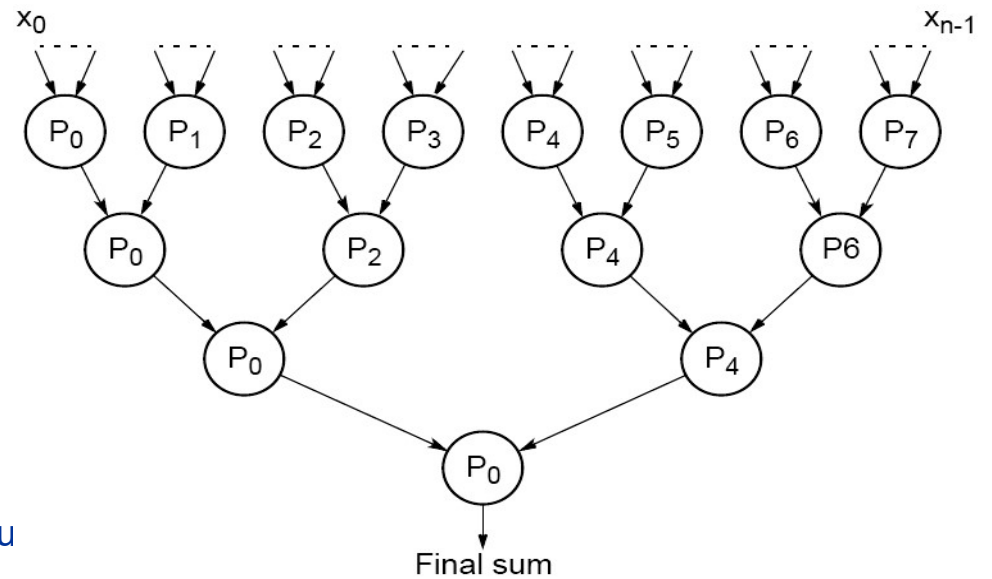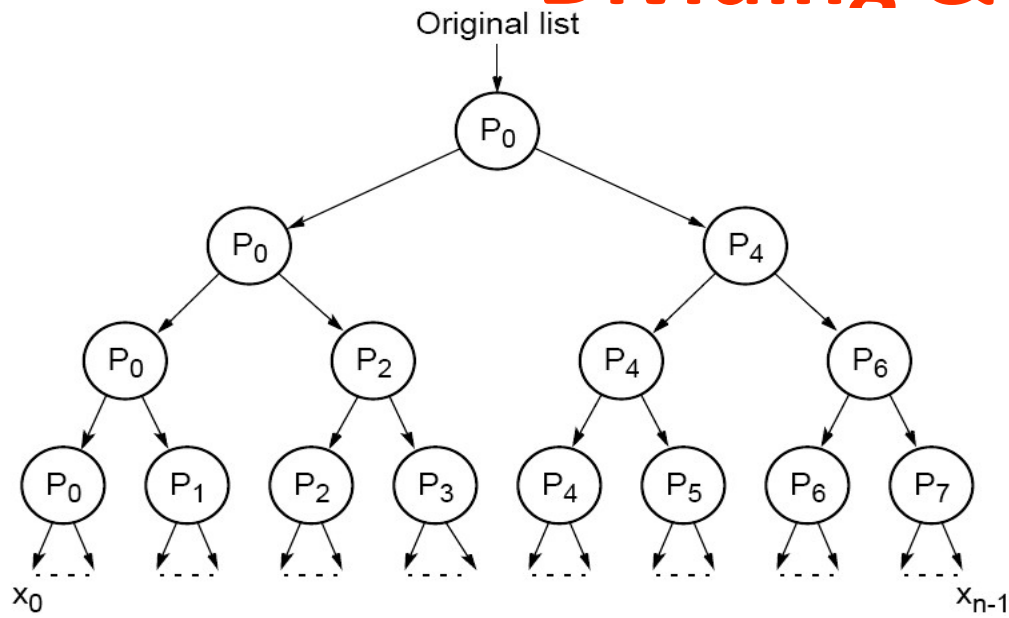- Recursive partitioning of tasks and collection of results
- Problems:
  - load balancing
  - least granularity

- Used on systems with background load

# Adding Numbers

$$x_0 \cdots x_{(n/p)-1} \quad x_{n/p} \cdots x_{(2n/p)-1} \quad \cdots \quad x_{(p-1)n/p} \cdots x_{n-1}$$

Partial sums

Sum

# Dividing & Conquer



Original list

$x_0$         $x_{n-1}$

$x_0$         $x_{n-1}$

Final sum

Lectu
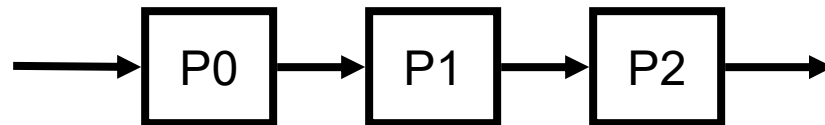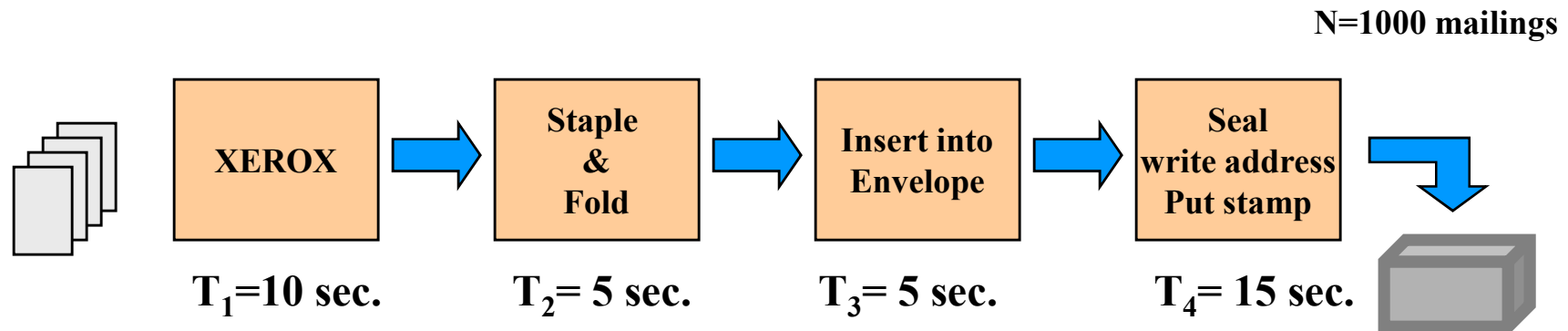
# Structure: Pipeline

➢ Examples

- o Different functions are applied to data: functional decomposition
- o Parallel execution of functions for different data.
- o Signal and image processing
- o Groundwater flow, flow of pollutants, visualization
- o Almost no example of high parallelism

```
  ──────▶ [ P0 ] ──▶ [ P1 ] ──▶ [ P2 ] ──────▶
```

# Pipelining:  Example

We would like to prepare and mail 1000 envelopes each containing
a document of 4 pages to members of an association.

**N=1000 mailings**

| XEROX | Staple & Fold | Insert into Envelope | Seal write address Put stamp |
|---|---|---|---|

$T_1$=10 sec.          $T_2$= 5 sec.          $T_3$= 5 sec.          $T_4$= 15 sec.

- At what intervals, do we see a new envelope prepared for mailing?

  $Max(T_1, T_2, \ldots, T_k) = T_{max} = 15$ sec.

- What is the total time to get N envelopes prepared?

  Time = Cold_Start_Time + $T_{max}$ * (N-1)  $\cong$  N* $T_{max}$

- What is the total time we would have spent if pipelining is not used?

  $N* \Sigma_i\, T_i$

# Pipelining: Example (contd.)

**N=1000 mailings**

| XEROX | Staple & Fold | Insert into Envelope | Seal write address Put stamp |
|---|---|---|---|

$T_1$=10 sec.  $T_2$= 5 sec.  $T_3$= 5 sec.  $T_4$= 15 sec.

How much speedup do we get?

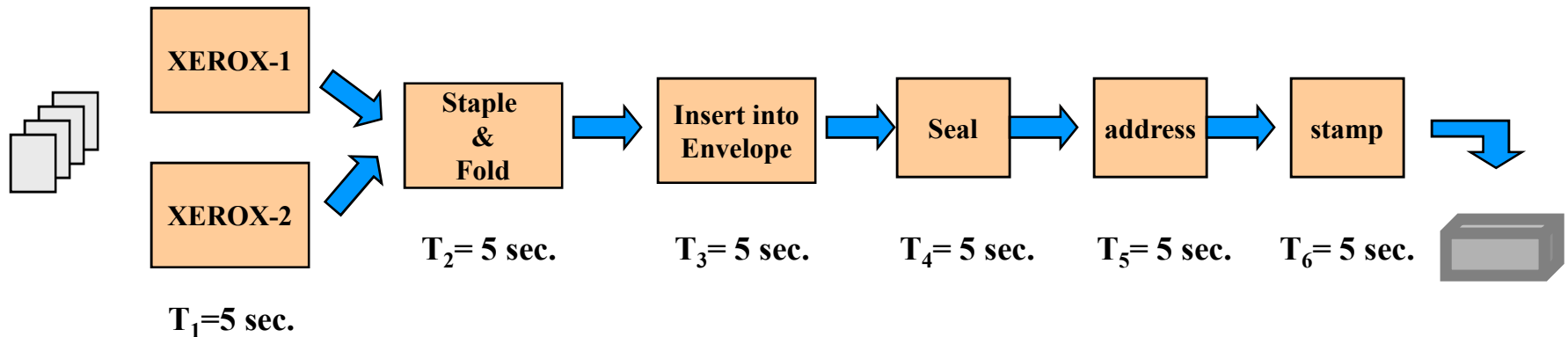Speedup $= T_{seq}/T_{pipe} = [\ N*\sum_i T_i\ ]/\ N* T_{max} = \sum_i T_i\ /\ T_{max}$

Speedup = 35/15

If you can not do much about the completion time for one task (i.e. $\sum_i T_i$ );
what can you do to maximize the speedup?
  (i) Create as many stations (stages) as possible
      and
  (ii) Try to balance the load at each station, i.e. $T_1 = T_2 = \ldots = T_k$

# Pipelining: Example (contd.)

One possible configuration to maximize speedup:



| XEROX-1 | | Staple & Fold | Insert into Envelope | Seal | address | stamp |

$T_1 = 5$ sec.
$T_2 = 5$ sec.
$T_3 = 5$ sec.
$T_4 = 5$ sec.
$T_5 = 5$ sec.
$T_6 = 5$ sec.
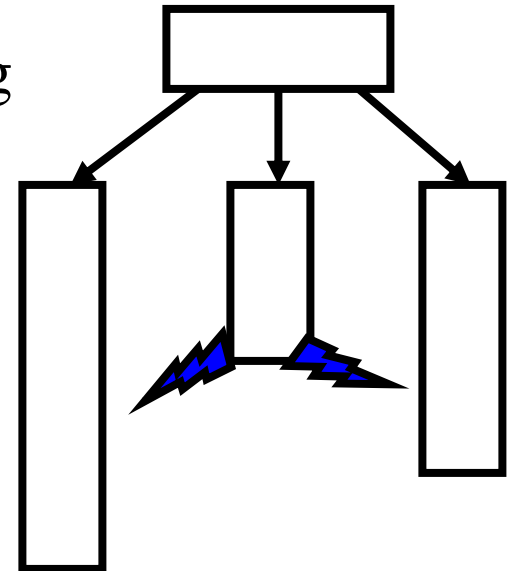
- At what intervals, do we see a new envelope prepared for mailing?

$$Max(T_1, T_2, \ldots, T_k) = T_{max} = 5 \text{ sec.}$$

- What is the speedup now?

Speedup $= 30/5 = 6 =$ number of stages in the pipeline !

# Structure: Competition

➢ Evaluation of multiple solution strategies in parallel.

➢ It might be unknown which strategy is successful or which one is the fastest.

➢ With k processors, k strategies can be tested. If one of the additional strategies - not tested in the sequential program - is very fast, the speedup can be more than k (Superlinear speedup)
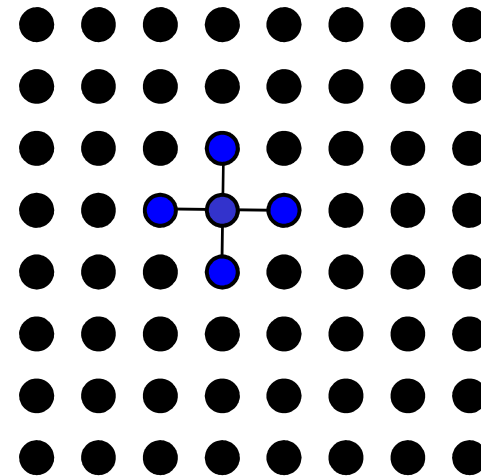
➢ Random search, speculative computing

# Application Structure: Application

➢ Performance Analysis

    o  Pipeline

    o  Based on the parameters to make decision

➢ System Support

    o  MapReduce

➢ Optimization

# Example: Equation Solver Kernel

➢ Solver kernel for a simple partial differential equation.

➢ Finite difference method

➢ Grid $(n+2) \times (n+2)$

➢ Fixed boundaries

➢ Interior points are recomputed

  o Mean value of five points in stencil

  o In place computation

  o Gauss-Seidel method

    • New values of upper and right point

    • Old values of lower and left point

  o Termination if difference between old and new value is below threshold for all points

```
A[i,j]=0.2*(A[i,j]+
  A[i-1,j]+A[i,j-1]+
  A[i,j+1]+A[i+1,j]
```

# Sequential Code

```
int n;                      /*size of matrix: (n + 2-by-n + 2)*/

float: **A, diff = 0;


main ()

begin

read(n);                    /*read input parameter: matrix size*/

A = malloc (a 2-d array of size n + 2 by n + 2 doubles);

initialize(A);              /*initialize the matrix A somehow*/

Solve (A);                  /*call the routine to solve equation*/


end main
```
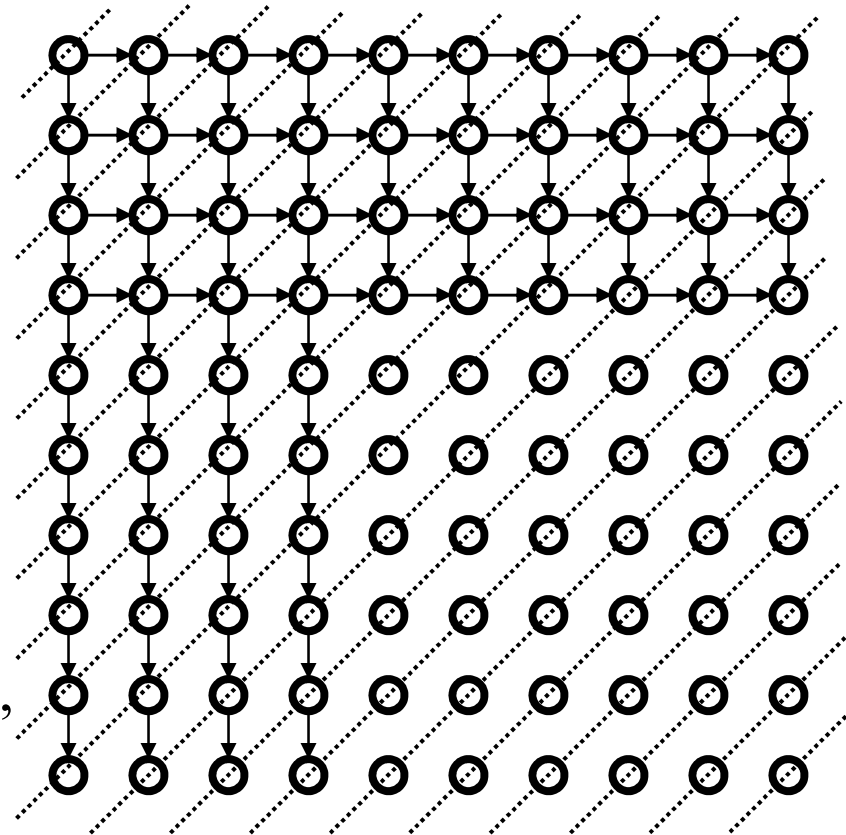
# Routine SOLVE

```
procedure Solve (A)          /*solve the equation system'/
float **A;                   /*A is an (n + 2)-by-(n + 2) array*/
begin
int i, j, done = 0;
float diff = 0, temp;
while (!done) do             /*outermost loop over sweeps*/
  diff = 0;                  /*initialize maximum difference to 0*/
  for i=1 to n do            /*sweep over nonborder points of grid*/
    for j=1 to n do
      temp = A[i,j];         /*save old value of element*/
      A[i,j] = 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
               A[i,j+1] + A[i+1,j]); /*compute average*/
      diff += abs(A[i,j] - temp);
    end for
  end for
  if (diff/(n*n) < TOL) then done = 1;
end while
end procedure
```

# Dependences in Gauss-Seidel

➢ Dependences prohibit row or column wise parallelization

➢ Point-wise synchronization

➢ Parallel execution along anti-diagonals
  o Proportional to n
  o Frequent synchronization, once per anti-diagonal
  o Load imbalance for short anti-diagonals

# Relaxing Ordering Constraints

➤ Jacobi iterations
  - o  Full sweep with old values
  - o  Much slower convergence➔more iterations
  - o  $N^2$ parallel tasks in each iteration

➤ Red-Black iterations:
  - o  Checkerboard like coloring scheme
  - o  Two phases
    - • Computation of red points with old black values
    - • Computation of black points with new red values
  - o  Faster convergence than Jacobi but more iterations than Gauss-Seidel
  - o  Each phase with $n^2/2$ parallel tasks

Red phase



Black phase

# Message Passing Programming (1/7)

- ➢ Processes have private address spaces

- ➢ Values are communicated via send/receive operations

- ➢ Writing local code alike to thread programming

# Message Passing Programming (2/7)

```
Int pid, n, nprocs;              /*process id and number of processes*/
Float **myA;


Main()
Begin
  nprocs=get_proc_num();         /*number of started processes*/
  pid=get_my_rank();             /*pid of the executing process*/
  if (pid==0) read(n);           /*master only*/
  broadcast(0,n,sizeof(int),SIZE); /*communicate to all others*/
  solve();                       /*start computation*/
End
```

# Message Passing Programming (3/7)

```
procedure solve()
begin
int I,j, pid, n1=n/nprocs, done=0;
float temp, tempdiff, mydiff=0;
myA=malloc(n/nprocs+2 by n+2);           /*allocate my rows+ghost rows*/
initialize (myA);                        /*intialize my rows*/


while (!done) do
  mydiff=0;
  if (pid!=0) then
     send(&myA[1,0],n*sizeof(float),pid-1,ROW)
  if (pid!=nprocs-1) then
     send(&myA[n1,0],n*sizeof(float),pid+1,ROW)
  if (pid!=0) then
     receive(&myA[0,0],n*sizeof(float),pid-1,ROW)
  if (pid!=nprocs-1) then
     receive(myA[n1+1,0],n*sizeof(float),pid+1,ROW)
```

# Message Passing Programming (4/7)

```
for i=1 to n1 do              /*sweep over nonborder points of grid*/
    for j=1 to n do
      temp = myA[i,j];        /*save old value of element*/
      myA[i,jl = 0.2 * (myA[i,j] + myA[i,j-1] + myA[i-1,j] +
                myA[i,j+1] + myA[i+1,j]); /*compute average*/
      mydiff += abs(myA[i,j] - temp);
    end for
  end for
```

# Message Passing Programming (5/7)

```
if (pid !=0) then                              /*send local diffs to P0*/
    send (mydiff,sizeof(float),0,DIFF)
    receive(done,sizeof(int),0,DONE)    /*receive done flag from P0*/
else
    for i=1 to nprocs-1 do
       receive(tempdiff,sizeof(float),*,DIFF) /*receive local diffs*/
        mydiff += tempdiff;                    /*accumulate local diffs*/
    endfor
    if (mydiff/(n*n)<TOL) then   done=1 /*check condition*/
    for i=1 to nprocs-1 do
       send(done,sizeof(int),i,done);    /*send done flag to other procs*/
    end for
endif


endwhile
end procedure
```

# Message Passing Programming (6/7)

```
reduce(0,mydiff,sizeof(float),ADD)
if (pid ==0) then
  if (mydiff/(n*n)<TOL) then done=1;
  endif
endif
broadcast(0,done,sizeof(int));
endwhile
end procedure
```

# Message Passing Programming (7/7)

| Syntax | Function |
|---|---|
| Send(src_addr, size, dest, tag) | Send size bytes starting at src_addr to the dest process, with tag identifier. |
| Receive(buffer_addr, size, src, tag) | Recieve a message with the tag identifier from the src process, and put size bytes of it into buffer starting at buffer_addr. |
| Reduce(root_pid, buffer, length, oper) | Compute global value from local values in buffer of the given length of all processes with operation oper. The global value is delivered to root process. |
| Broadcast(root_pid, buffer, length, tag) | The root process sends the value in buffer to the other processes with the tag identifier. |

# The Parallel Partition LU (PPT) Algorithm

*Step* 1. Allocate $A_i, d^{(i)}$ and elements $a_{im}, c_{(i+1)m-1}$ to

the *ith* node, where $0 \le i \le p-1$.

*Step* 2. Use the *LU* decomposition method to solve

$$A_i[\tilde{x}^{(i)}, v^{(i)}, w^{(i)}] = [d^{(i)}, a_{im}e_0, c_{(i+1)m-1}e_{m-1}]$$

*Step* 3. Send $\tilde{x}_0^{(i)}, \tilde{x}_{m-1}^{(i)}, v_0^{(i)}, v_{m-1}^{(i)}, w_0^{(i)}, w_{m-1}^{(i)}$ from the

*ith* node to the other nodes $0 \le i \le p-1$.

*Step* 4. Use the LU method to solve Zy = h on all nodes

*Step* 5. Compute in parallel on $p$ processors

$$\Delta x^{(i)} = [v^{(i)}, w^{(i)}] \begin{bmatrix} y_{2i-1} \\ y_{2i} \end{bmatrix}$$

$$x^{(i)} = \tilde{x}^{(i)} - \Delta x^{(i)}$$

# The Solving process

1. Solve the subsystems in parallel
2. Solve the reduced system
3. Modification

$$\widetilde{\mathbf{A}}^{-1}\mathbf{A}\mathbf{x} = \widetilde{\mathbf{A}}^{-1}\mathbf{d}$$

$$
\begin{bmatrix} A_0^{-1} & & & \\ & A_1^{-1} & & \\ & & & \\ & & & A_{p-1}^{-1} \end{bmatrix}
\begin{bmatrix} A_0 & & & \\ {}^* {}^* & A_1 {}^* & & \\ & {}^* & & \\ & & {}^* {}^* & A_{p-1} \end{bmatrix}
\begin{pmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ \cdot \\ x_{n-1} \end{pmatrix}
=
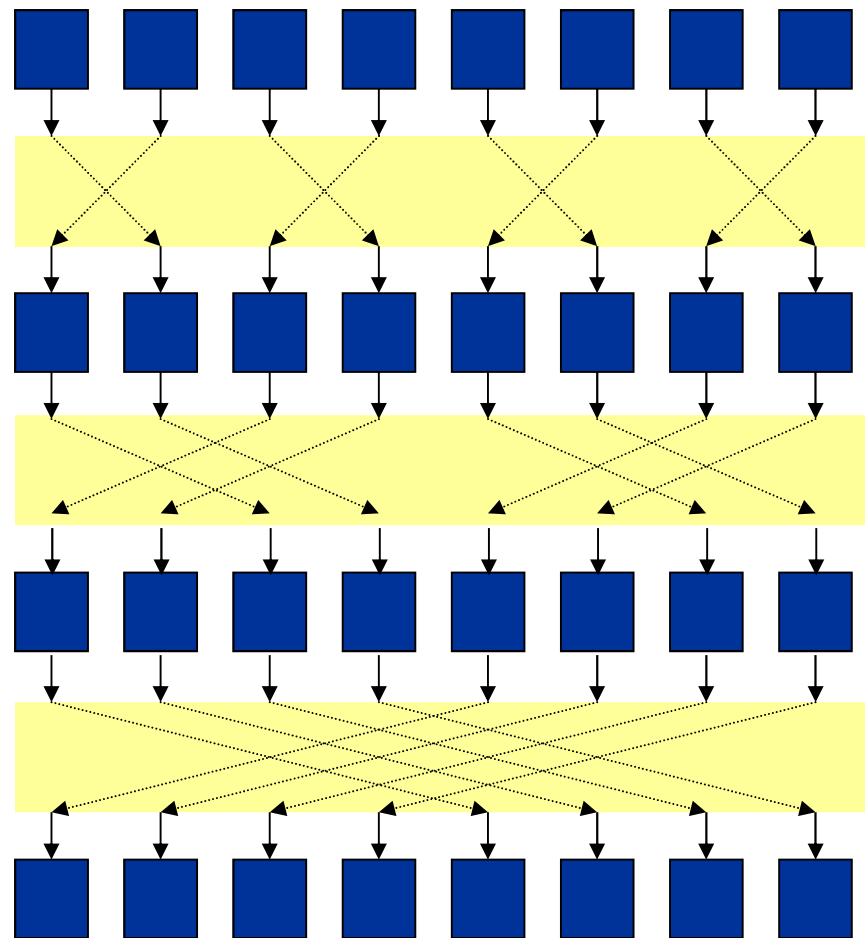\begin{bmatrix} A_0^{-1} & & & \\ & A_1^{-1} & & \\ & & & \\ & & & A_{p-1}^{-1} \end{bmatrix}
\begin{pmatrix} d_0 \\ d_1 \\ \cdot \\ \cdot \\ \cdot \\ d_{n-1} \end{pmatrix}
$$

# The Reduced System  $(Zy=h)$



**Needs global communication**

# All-to-All Total Data Exchange

# Summary Parallel Programming

➢ Steps in the parallelization process

   o  Decomposition, assignment, orchestration, mapping

➢ Variety of programming models for SM and DM systems

   o  Different parallel languages and APIs

   o  Trend towards standards (OpenMP, MPI, Posix)

# Summary Parallel Programming Models

➢ **Global vs local programming**

  o Global: directive-based, data parallel

  o Local: thread-based, remote memory access, message passing

➢ **Data vs functional parallelism**

  o Data parallelism only: data parallel

  o Both: all the others

➢ **Parallelism**

  o High: remote memory access, message passing, thread-based

  o Low: directive-based, high level data parallel

**Figure 4-1    A Dog**



**Figure 4-2    A Pack of Dogs**



**Figure 4-3    A Savage Multiheaded Pooch**

# Concurrent-AMAT: step to optimization

- The traditional AMAT(Average Memory Access Time) :

  $$AMAT = HitCycle + MR \times AMP$$

- MR is the miss rate of cache accesses; and AMP is the average miss penalty

- **Concurrent-AMAT (C-AMAT):**

  $$C\text{-}AMAT = HitCycle/C_H + pMR \times pAMP/C_M = 1/APC$$

- $C_H$ is the hit concurrency; $C_M$ is the *pure* miss concurrency
- *p*MR and *p*AMP are *pure* miss rate and average *pure* miss penalty
- A pure miss is a miss containing at least one cycle which does not have any hit activity

# Recursive in Memory Hierarchy

- **AMAT is recursive**
  - $AMAT = HitCycle_1 + MR_1 \times AMP_1$

  Where $AMP_1 = (HitCycle_2 + MR_2 \times AMP_2)$

- **C-AMAT is also recursive**

$$C\text{-}AMAT_1 = \frac{H_1}{C_{H_1}} + MR_1 \times \kappa_1 \times C\text{-}AMAT_2$$

Where

$$C\text{-}AMAT_2 = \frac{H_2}{C_{H_2}} + pMR_2 \times \frac{pAMP_2}{C_{M_2}}$$

$$\kappa_1 = \frac{pMR_1}{MR_1} \times \frac{pAMP_1}{AMP_1} \times \frac{C_{m_1}}{C_{M_1}}$$

**With Clear Physical Meaning**

# Memory Hierarchy Performance

**1 clk** First-level Cache **300 clks** Main Memory (DRAM)

**Hit Time**

**Miss % * Miss penalty**

- Average Memory Access Time (AMAT)

  = Hit Time + Miss rate * Miss Penalty

  = $T_{hit}(L1)$ + Miss%(L1) * T (memory)

- Example:

  - Cache Hit = 1 cycle

  - Miss rate = 10%

  - Miss penalty = 300 cycles

  - AMAT = 1 + 300*10% = 31 cycles

- To further improve it?

# Example of calculating C-AMAT

Hit Time
1 clk

**First-level Cache**

2
Hit Concurrency

Pure Miss Rate
10%

Pure Miss Penalty
300 clks

3
Pure Miss Concurrency

**Main Memory (DRAM)**

Concurrent Average Memory Access Time (C-AMAT)
= Hit Time/ Hit Concurrency + Pure Miss Rate * (Pure Miss Penalty/ Pure Miss Concurrency)

- Example:
  - Hit Time = 1 cycle
  - Hit Concurrency = 2
  - Pure Miss Rate = 10%
  - Pure Miss Penalty = 300 cycles
  - Miss Concurrency = 3
  - C-AMAT = 1/2 + 10%*(300/3) = 0.5 +10 = 10.5 cycles

# Data Access Time: AMAT



**Hit Time**
**1 clk**

**2**
**Hit Concurrency**

**10 clks**

**3**

**20 clks**

**4**

**300 clks**

**6**

L1 Cache

**L2 Cache**

**L3 Cache**

**Main Memory (DRAM)**

**On-die**

- Average Memory Access Time (AMAT)
  $= T_{hit}(L1) + Miss\%(L1)* (T_{hit}(L2) + Miss\%(L2)* (T_{hit}(L3) + Miss\%(L3)*T(memory) ) )$
- Example: (Latency as shown above)
  - Miss rate: L1=10%, L2=5%, L3=1% (*Be careful miss rate definition*)
  - AMAT
    = 2.115

# Data Access Time:  C-AMAT



- **Concurrent Average Memory Access Time (C-AMAT)**

$$= \frac{H_1}{C_{H_1}} + MR_1 \times \kappa_1 \times \left( \frac{H_2}{C_{H_2}} + MR_2 \times \kappa_2 \times \left( \frac{H_3}{C_{H_3}} + MR_3 \times \kappa_3 \times \frac{H_{Mem}}{C_{H_{Mem}}} \right) \right)$$

- **Example**
  - Miss Rate: L1=10%, L2=5%, L3=1%    pMR, pAMP, AMP, $C_M$, $C_m$:  L1=7%, 10, 10, 5, 4
  - κ: L1=0.56, L2=0.6, L3=0.8                                                      L2=3%, 60, 40, 9, 6
  - C-AMAT≈0.696                                                                          L3=0.8%, 400, 300, 16, 12

# C-AMAT in Multi-Core Environments

# C-AMAT in Multi-Core Environments

# What Does C-AMAT Say?

- C-AMAT is an extension of AMAT to consider concurrency

  - The same as AMAT, if no concurrency present

- C-AMAT introduces the **Pure Miss** concept:

  - Only pure miss causes performance penalty

- **High locality may hurt performance**

  - High locality may lead to pure miss

- **Two contributions of concurrency**

  - Increase bandwidth

  - Latency hiding (overlapping)

# What Does C-AMAT Say?

- C-AMAT also contains the overlapping factor

$$C\text{-}AMAT_1 = \frac{H_1}{C_{H_1}} + MR_1 \times \kappa_1 \times C\text{-}AMAT_2$$

- **Balance** locality, concurrency and Overlapping with C-AMAT

- A good explanation for why
Optimal $\neq$ Optimal Locality + Optimal Concurrence

- C-AMAT uniquely integrates the **joint impact** of locality concurrency, and overlapping for optimization

- Overlapping in connection locality and concurrency is a new issue of research

# Impact of C-AMAT

- New dimensions for optimization: **concurrency, overlapping and balancing**

$$C\text{-}AMAT_1 = \frac{H_1}{C_{H_1}} + MR_1 \times \kappa_1 \times C\text{-}AMAT_2$$

- Can apply at **each layer** of a memory hierarchy

- Existing mechanisms are readily to be extended

  - Every AMAT based optimization has a corresponding C-AMAT extension to include concurrency

- Concurrency as bandwidth increaser and **penalty reducer (overlapping)**: Accurate measure the concurrency contribution

$$\kappa_1 = \frac{pMR_1}{MR_1} \times \frac{pAMP_1}{AMP_1} \times \frac{C_{m_1}}{C_{M_1}}$$

# Misunderstanding of Memory Performance

*Optimal* = ~~Optimal Locality~~ **?**

*Optimal* = Optimal Loca~~l~~ity + Optimal Concurrence **?**

# Practical: Parameters can be measured in hardware



Feedback-based optimization on scheduling and on reconfigurable architecture

# Challenge in Practice

- C-AMAT is general and powerful, but its measurement is environment (hardware) dependent

- Need a good understanding to conduct C-AMAT analysis in Multiccore, GPU, GPGPU, FPGA ASIC, etc. environments

- GPU is a special case of C-AMAT matching

# **Memory stall time:** *the performance we care*

**Traditional AMAT model**

$$CPU\text{-}time = IC \times (CPI_{exe} + f_{mem} \times AMAT) \times Cycle\text{-}time$$

Memory stall time

## **New C-AMAT model**

$$CPU\text{-}time = IC \times (CPI_{exe} + f_{mem} \times C\text{-}AMAT \times (1 - overlapRatio_{c\text{-}m})) \times cycle\text{-}time$$

Memory stall time

$$CPU\text{-}time = IC \times (CPI_{exe} + f_{mem} \times \frac{pMR \times pAMP}{C_M}) \times Cycle\text{-}time$$

Memory stall time

Only pure miss will cause processor stall, and the penalty is formulated here

Y. Liu and X.-H. Sun, "*Reevaluating Data Stall Time with the Consideration of Data Access Concurrency,*" Journal of Computer Science and Technology (JCST), March, 2015

# Application: Layered Performance Matching

Request rates
of computing
components

ALU&FPU

Supply rates
of $L_1$ cache

$APC_1$

$L_1$ cache

Request rates
of $L_1$ cache

Supply rates of
Last level cache

$APC_2$

Request rates of
Last level cache

Last level cache

Supply rates of
main memory

$APC_3$

Main memory

Yu-Hang Liu, Xian-He Sun, "LPM: Concurrency-driven Layered Performance Matching," in ICPP2015, Beijing, China, Sept. 2015.

47

# Idea: Match the Request with Supply

$$LPMR(ALU\,\&\,FPU, L_1) = \frac{Request\ rate\ from\ ALU\&FPU}{Supply\ rate\ by\ L_1\ cache}$$

$$LPMR(L_1, LLC) = \frac{Request\ rate\ from\ L_1\ cache}{Supply\ rate\ by\ LLC}$$

$$LPMR(LLC, MM) = \frac{Request\ rate\ from\ LLC}{Supply\ rate\ by\ main\ memory}$$

- Match at **each memory layer**
- Adjust the supply performance with concurrency

# Quantify Mismatching: with C-AMAT

$$LPMR_1 = \frac{IPC_{exe} \times f_{mem}}{APC_.}$$

$$LPMR_2 = \frac{IPC_{exe} \times f_{mem} \times MR_1}{APC_2}$$

$$LPMR_3 = \frac{IPC_{exe} \times f_{mem} \times MR_1 \times MR_2}{APC_3}$$

- C-AMAT measures the request and supply at each layer
- C-AMAT can increase supply with effective concurrency
- Mismatch ratio directly determines memory stall time

# C-AMAT in Action

**New C-AMAT model**

$$CPU\text{-}time = IC \times (CPI_{exe} + f_{mem} \times \frac{pMR \times pAMP}{C_M}) \times Cycle\text{-}time$$

Memory stall time

Only pure miss will cause processor stall, and the penalty is formulated here

**The Relation of LPMR and Stall time**

$$CPU\text{-}time = IC \times CPI_{exe} \times (1 + \kappa_1 \times LPMR_2) \times Cycle\text{-}time$$

Memory stall time

Y. Liu and X.-H. Sun, "Reevaluating Data Stall Time with the Consideration of Data Access Concurrency," Journal of Computer Science and Technology (JCST), March, 2015

# The LPM Model

**We get the relation between data stall time and LPMR1**

$$Data - stall - time = CPI_{exe} \times (1 - overlapRatio_{c-m}) \times LP.$$

**We get the relation between data stall time and LPMR2**

$$Data - stall -time = CPI_{exe} \times \kappa_1 \times LPMR_2$$

**Because our final goal is to minimize data stall time, the two equations above. provide the baseline for LPM optimizations**

# The Threshold Requirement

$$overlapRatio_{c-m} = \frac{H_1/C_{H_1}}{(H_1/C_{H_1} + pAMP_1/C_{M_1})}$$

$$LPMR_1 \leq \frac{\Delta\%}{1 - overlapRatio_{c-m}}$$

$$LPMR_2 \leq \frac{\Delta\%}{\kappa_1}$$

This is a presentation slide.

# A Matching Example (increase performance)

$$LPMR = \frac{IPC_{exe} \times f_{mem}}{APC} \qquad (1)$$

- **Assume**:

$$Cycle_{CPU} = 2 \text{ ns}$$
$$Cycle_{mem} = 8 \text{ ns}$$
$$f_{mem} = 20\%$$
$$IPC_{exe} = 2.5$$
$$APC = 1$$

- **Then**:
  - In 8 ns, there is 1 memory cycle, and the data supply rate is:
  $$APC * N\_Cycle_{mem} = 1$$
  - In 8 ns, there is 4 cpu cycle, and the data request rate is:
  $$IPC_{exe} * N\_Cycle_{CPU} * f_{mem} = 2$$

- **So**:
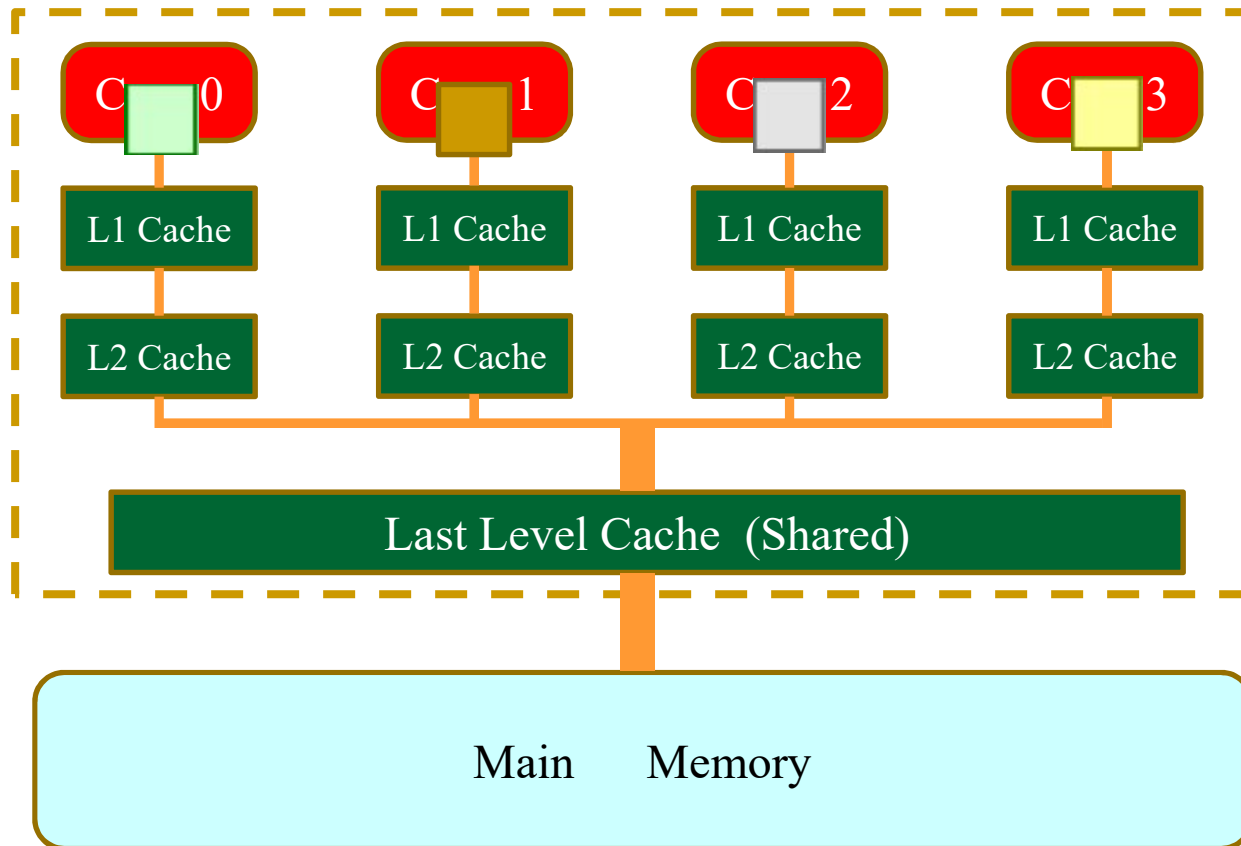  - The data supply rate does not match the data request rate. We can improve memory *concurrency* by adding memory banks or ports to increase data supply ability to adjust the data supply rate. For example, increase *APC* from 1 to 2, so:
  $$APC * N\_Cycle_{mem} = 2$$

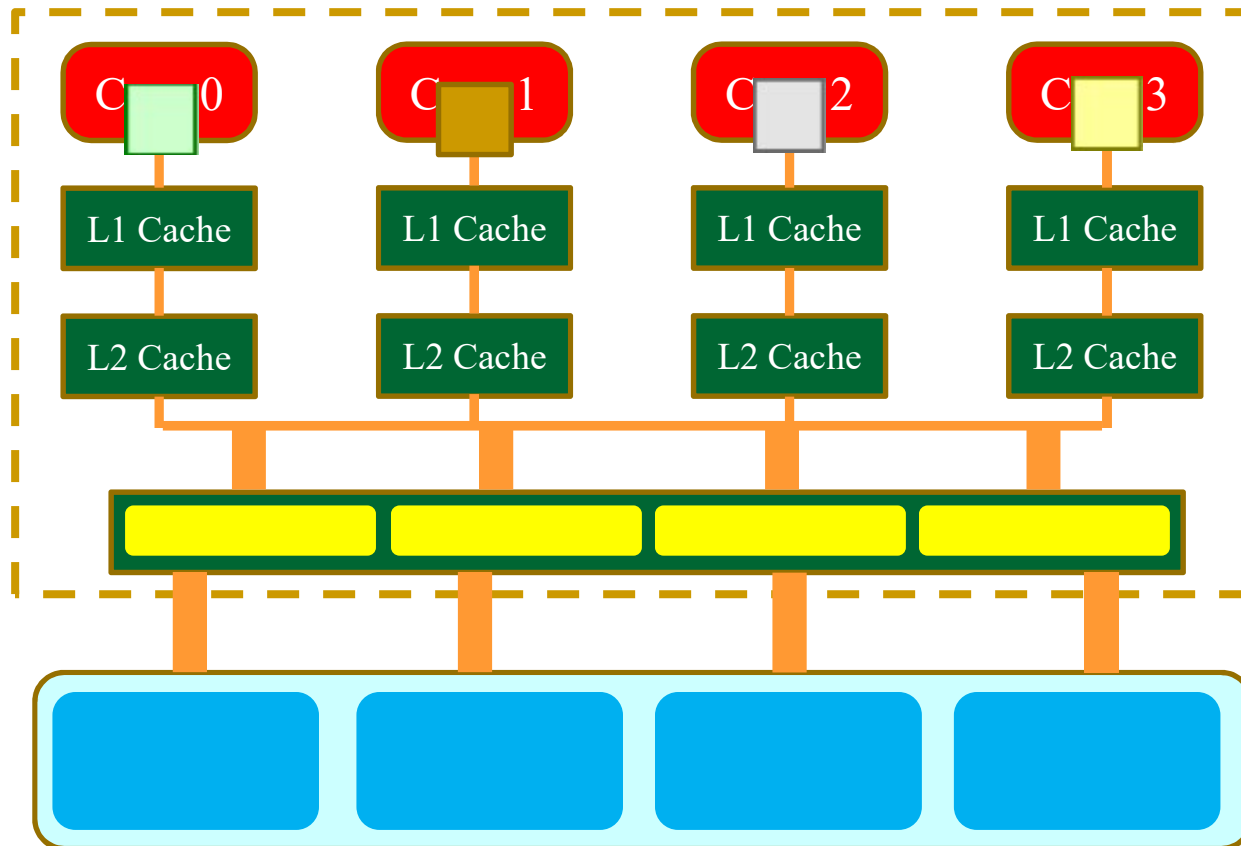  - And the data supply rate matches with the data request rate.

# Memory Access (Concurrency)



**Model and algorithm**

# Memory Access (Concurrency)



**Model and algorithm**

# The LPM Algorithm

BEGIN

Measure $LPMR_1$ and $LPMR_2$

Calculate the thresholds of T1 and T2

$LPMR_1 < T_1$

Yes

No

$LPMR_1 + \Delta < T_1$

Yes

No

$LPMR_2 < T_2$

Yes

No

Reduce hardware overprovision, and update all metrics

Optimize both $L_1$ and $L_2$ layer to reduce $LPMR_1$ and $LPMR_2$, and update all metrics

Optimize only $L_1$ layer to reduce $LPMR_1$, and update all metrics

Stop when stall time is less than 1% of pure execution time

END

**LPMR Reduction Algorithm**

1:       // Initially measure the metrics

2:        For each application or thread, measure the LPMRs in a memory hierarchy

3:        Get the threshold $T_1$ and $T_2$ according to Eq. (22) and (23)

4:     **Begin Do**

5:       // LPM optimization loop

6:         // Case I when both L1 and L2 layer need an optimization

7:           **While** (LPMR$_1$ > $T_1$ and LPMR$_2$ > $T_2$) **Do**

8:              Optimizing at L1 layer and L2 layer,

9:              Update all the metrics (LPMR$_1$, LPMR$_2$, $T_1$ and $T_2$)

10:        **Until** (LPMR$_1$ ≤ $T_1$ or LPMR$_2$ ≤ $T_2$)

11:       // Case II when only L1 layer needs an optimization

12:        **While** (LPMR$_1$ > $T_1$ and LPMR$_2$ ≤ $T_2$) **Do**

13:           Optimizing at L1 layer

14:           Update all the metrics (LPMR$_1$, LPMR$_2$, $T_1$ and $T_2$)

15:        **Until** (LPMR$_1$ ≤ $T_1$)

16:       // Case III when no layer needs to optimize and overprovision may need to reduce

17:       // δ is a positive value

18:        **While** (LPMR$_1$ + δ < $T_1$) **Do**

19:         Reduce hardware overprovision

20:         Update all the metrics (LPMR$_1$, LPMR$_2$, $T_1$ and $T_2$)

21:        **Until** (LPMR$_1$ ≥ $T_1$ − δ )

22:       // Case IV when no layer needs to optimize and no overprovision needs to reduce

23:        **If** ($T_1$ ≥ LPMR$_1$ ≥ $T_1$ − δ)

24:            End the algorithm

25:       **Endif**

26:     **Until (End)**

**Pseudo code**

# Comments and Questions:

- Pipelining can be combined with concurrency for best performance: subarray

- Can you derive the general form of LPMR_i?