

Design of Parallel and High-Performance Computing:  
**Distributed-Memory Models and Algorithms**

Edgar Solomonik

ETH Zurich

December 7, 2015

# Summary

## Lecture overview

- $\alpha$ - $\beta$  communication cost model
- LogP model
- LogGP model (LogP with variable-size messages)
- Algorithms for broadcasts of large messages
- Other collective communication patterns
- Bulk Synchronous Parallel (BSP) model
- PGAS languages / one-sided communication
- Communication-avoiding algorithms
- Overview and final comments

## A simple model for point-to-point messages

The time to send or receive a message of  $s$  bytes is

$$T_{\text{sr}}^{\alpha,\beta}(s) = \alpha + s \cdot \beta$$

- $\alpha$  – **latency/synchronization cost** per message
- $\beta$  – **bandwidth cost** per byte
- each processor can send and/or receive one message at a time

Let  $P$  processors send a message of size  $s$  in a ring,

- the **communication volume** (total amount of data sent) is  $P \cdot s$
- What is the **communication cost** ( $\alpha$ - $\beta$ -model execution time)?
  - if the messages are sent *simultaneously*,

$$T_{\text{sim-ring}}^{\alpha,\beta}(s) = T_{\text{sr}}^{\alpha,\beta}(s) = \alpha + s \cdot \beta$$

- if the messages are sent *in sequence*,

$$T_{\text{seq-ring}}^{\alpha,\beta}(s, P) = P \cdot T_{\text{sr}}^{\alpha,\beta}(s) = P \cdot (\alpha + s \cdot \beta)$$

# Broadcasts in the $\alpha$ - $\beta$ model

The execution time of a broadcast of a message of size  $s$  to  $P$  processors is

- using a **binary tree** of height

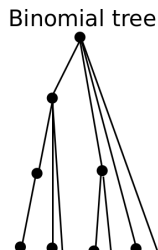
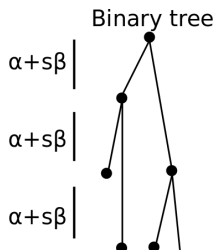
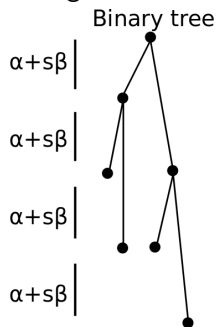
$$h = 2(\log_2(P + 1) - 1),$$

$$\begin{aligned} T_{\text{bcast-bin}}^{\alpha,\beta}(s, P) &= h \cdot T_{\text{sr}}^{\alpha,\beta}(s) \\ &= 2(\log_2(P+1) - 1) \cdot (\alpha + s \cdot \beta) \end{aligned}$$

- using a **binomial tree** of height

$$h' = \log_2(P + 1),$$

$$\begin{aligned} T_{\text{bcast-bnm}}^{\alpha,\beta}(s, P) &= h' \cdot T_{\text{sr}}^{\alpha,\beta}(s) \\ &= \log_2(P+1) \cdot (\alpha + s \cdot \beta) \end{aligned}$$



# The LogP model

Limitations of the  $\alpha$ - $\beta$  messaging model:

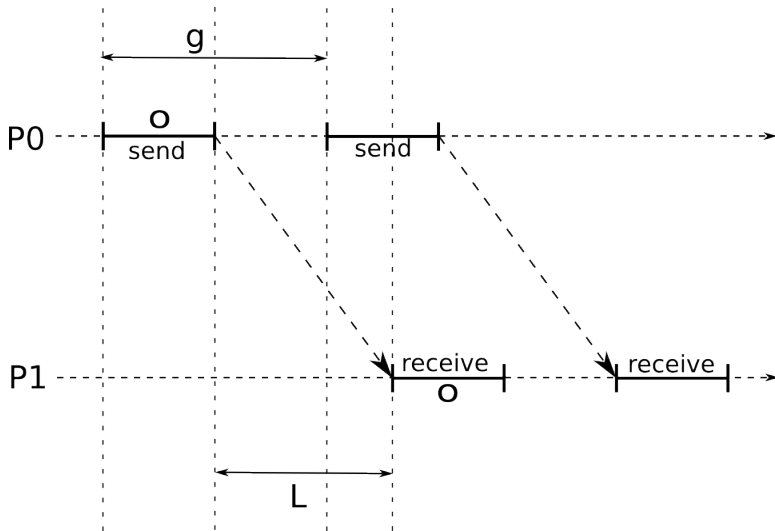
- both sender and receiver block until completion
- a processor cannot send multiple messages simultaneously
- no overlap between communication and computation

The **LogP model** (Culler et al. 1996) enables modelling of overlap by modelling the cost of sending a message of one 'datum' in terms of

- $L$  – network **latency** cost (processor free)
- $o$  – sender/receiver sequential **overhead** (processor occupied)
- $g \geq o$  – **gap** between two sends or two receives (processor free)
- $P$  – number of **processors**
- the LogP communication cost for sending a message of  $s$  datums is

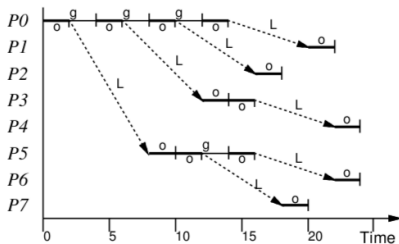
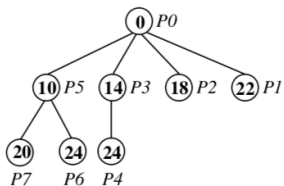
$$T_{sr}^{\text{LogP}}(s) = 2o + L + (s - 1) \cdot g$$

# Messaging in the LogP model



# Broadcasts in the LogP model

Same idea as binomial tree, forward message as soon as it is received, keep forwarding until all nodes obtain it (Karp et al. 1993)



## The LogGP model

The LogP model parameter  $g$  is associated with the datum size

- this injection rate implies a fixed-sized packet (datum) can be sent anywhere after a time interval of  $g$
- modern computer networks do not have a small fixed packet size and achieve higher bandwidth for large messages

The **LogGP model** (Alexandrov et al. 1997) introduces another bandwidth parameter  $G$ , which dictates the large-message bandwidth

- $G$  – **Gap per byte**; time per byte (processor free)
- $g \geq o$  – **gap** between injection/retrieval of bytes of two messages
- the LogGP time for sending a message of  $s$  bytes is

$$T_{sr}^{\text{LogGP}}(s) = 2o + L + (s - 1) \cdot G$$



# The LogGP model

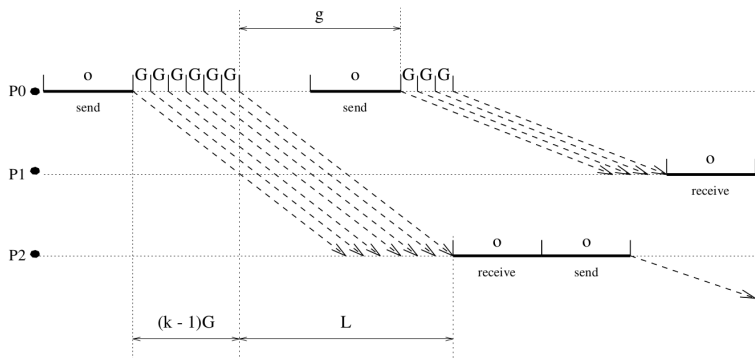


Diagram taken from: Alexandrov, A., Ionescu, M. F., Schauser, K. E., and Scheiman, C. LogGP: incorporating long messages into the LogP model—one step closer towards a realistic model for parallel computation. ACM SPAA, July 1995.

## Large-message broadcasts

Lets now consider broadcasts of a message of a size  $s \geq P$  bytes

- recall binomial tree broadcast cost:

$$T_{\text{bcast-bnm}}^{\alpha-\beta}(s, P) = \log_2(P+1) \cdot (\alpha + s \cdot \beta)$$

- consider instead the following broadcast schedule
  - the root sends a different segment of the message to each processor
  - all processors exchange segments in  $P - 1$  near-neighbor ring exchanges
- the cost of this broadcast schedule is

$$\begin{aligned} T_{\text{bcast-ring}}^{\alpha-\beta}(s, P) &= (P - 1)(T_{\text{sr}}^{\alpha-\beta}(s/P) + T_{\text{sim-ring}}^{\alpha-\beta}(s/P)) \\ &= 2(P - 1)(\alpha + s/P \cdot \beta) \end{aligned}$$

- for sufficiently large message sizes, the new schedule is faster,

$$\lim_{s \rightarrow \infty} \left( \frac{T_{\text{bcast-bnm}}^{\alpha-\beta}(s, P)}{T_{\text{bcast-ring}}^{\alpha-\beta}(s, P)} \right) \approx \log_2(P)/2$$

## Pipelined binary tree broadcast

Send a fixed-size packet to left child then to right child (entire message of size  $s$ )

- if the LogP model datum size is  $k_{\text{LogP}}$  bytes, the LogP cost is

$$T_{\text{PBT}}^{\text{LogP}}(s, P) \approx \log(P) \cdot (L + 2g + o) + 2(s/k_{\text{LogP}}) \cdot g$$

- in the LogGP model, we can select a packet size  $k$  and obtain the cost

$$T_{\text{PBT}}^{\text{LogGP}}(s, P, k) \approx \log(P) \cdot (L + 2g + o + 2k \cdot G) + 2(s/k) \cdot (g + k \cdot G)$$

minimizing the packet size  $k$ ,

$$k_{\text{opt}}^{\text{LogGP}}(s, P) = \underset{k}{\operatorname{argmin}}(T_{\text{PBT}}^{\text{LogGP}}(s, P, k))$$

(via e.g. differentiation by  $k$ ) we obtain the optimal packet size

$$k_{\text{opt}}^{\text{LogGP}}(s, P) = \sqrt{s/\log(P)} \cdot \sqrt{\frac{g}{G}}$$

so the best packet size, depends not only on architectural parameters, but also on dynamic parameters: the number of processors and message size

## Pipelined binary tree broadcast contd.

In LogP we obtained

$$T_{\text{PBT}}^{\text{LogP}}(s, P) \approx \log(P) \cdot (L + 2g + o) + 2(s/k_{\text{LogP}}) \cdot g$$

In LogGP we obtained,

$$T_{\text{PBT}}^{\text{LogGP}}(s, P, k) \approx \log(P) \cdot (L + 2g + o + 2k \cdot G) + 2(s/k) \cdot (g + k \cdot G)$$

$$k_{\text{opt}}^{\text{LogGP}}(s, P) = \sqrt{s/\log(P)} \cdot \sqrt{\frac{g}{G}}$$

- in the  $\alpha$ - $\beta$  model for a packet size of  $k$ , we obtain the cost

$$T_{\text{PBT}}^{\alpha, \beta}(s, P, k) \approx 2(\log(P) + s/k)(\alpha + k \cdot \beta)$$

with a minimal-cost packet size of

$$k_{\text{opt}}^{\alpha, \beta}(s, P) = \sqrt{s/\log(P)} \cdot \sqrt{\frac{\alpha}{\beta}}$$

## Pipelined binary tree broadcast conclusions

The LogP model is inflexible, while the LogGP and the  $\alpha$ - $\beta$  models capture the key **input** and **architectural** scaling dependence

$$T_{\text{PBT}}^{\alpha,\beta}(s, P, k) \approx 2(\log(P) + s/k)(\alpha + k \cdot \beta)$$

$$k_{\text{opt}}^{\alpha,\beta}(s, P) = \sqrt{s/\log(P)} \cdot \sqrt{\frac{\alpha}{\beta}}$$

The minimized cost in the  $\alpha$ - $\beta$  model is

$$\begin{aligned} T_{\text{oPBT}}^{\alpha,\beta}(s, P) &= T_{\text{PBT}}^{\alpha,\beta}(s, P, k_{\text{opt}}^{\alpha,\beta}(s, P)) \\ &\approx 2 \left( \log(P) + \sqrt{s \cdot \log(P)} \cdot \sqrt{\frac{\beta}{\alpha}} \right) \cdot \left( \alpha + \sqrt{\frac{s}{\log(P)}} \cdot \sqrt{\alpha \cdot \beta} \right) \\ &= 2 \log(P) \cdot \alpha + 4 \sqrt{s \cdot \log(P)} \cdot \sqrt{\alpha \cdot \beta} + 2s \cdot \beta \end{aligned}$$

**Q:** Could we get rid of the factor of two constant in the  $O(s \cdot \beta)$  cost?

**A:** Not so long as the root sends two copies of the whole message...

## Double Tree

**Observation:** *the leaves of a binary tree,  $(P - 1)/2$  processors, send nothing, while the internal nodes do all the work.*

### Double Pipelined Binary Tree Broadcast

- define two pipelined binary trees with a shared root
- non-root processors act as a leaf in one and as an internal node in the second
- send half of the message down each tree

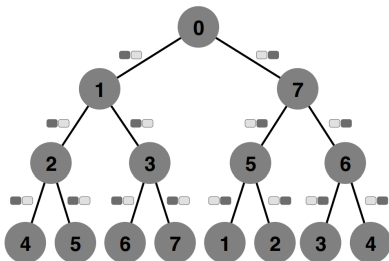


Diagram taken from: Hoefler, Torsten, and Dmitry Moor. "Energy, Memory, and Runtime Tradeoffs for Implementing Collective Communication Operations."

## Double pipelined binary tree

The cost of the double pipelined binary tree is essentially the same as the cost of a single pipelined binary tree with half the message size,

$$T_{\text{DPBT}}^{\alpha, \beta}(s, P) \approx 2 \log(P) \cdot \alpha + 2\sqrt{2s \cdot \log(P)} \cdot \sqrt{\alpha \cdot \beta} + s \cdot \beta$$

for a sufficiently large message size ( $s$ ) this is twice as fast as a single pipelined binary tree.

How close is the double pipelined binary tree to optimum?

- for fixed-size packets, lower bound (Johnsson and Ho 1989) is

$$T_{\text{broadcast-lb}}^{\alpha, \beta}(s, P) \approx \log(P) \cdot \alpha + 2\sqrt{s \cdot \log(P)} \cdot \sqrt{\alpha \cdot \beta} + s \cdot \beta$$

- attained by algorithm of Träff and Ripke 1995,

$$T_{\text{broadcast}}^{\alpha, \beta}(s, P) = T_{\text{broadcast-lb}}^{\alpha, \beta}(s, P).$$

- showing optimality for variable-size packets is an open question

## Other types of collective communication

We can classify collectives into four categories

- **One-to-All:** Broadcast, Scatter
- **All-to-One:** Reduce, Gather
- **All-to-One + One-to-All:** Allreduce (Reduce+Broadcast), Allgather (Gather+Broadcast), Reduce-Scatter (Reduce+Scatter), Scan
- **All-to-All:** All-to-all

MPI (Message-Passing Interface) provides all of these as well as variable size versions (e.g. (All)Gatherv, All-to-allv), see online for specification of each routine.

We now present algorithms for and their cost in the  $\alpha - \beta$  model, with

$$s = \begin{cases} \text{input size} & : \text{one-to-all collectives} \\ \text{output size} & : \text{all-to-one collectives} \\ \text{per-processor input/output size} & : \text{all-to-all collectives} \end{cases}$$



## Tree collectives

We have demonstrated how (double/pipelined) binary trees and binomial trees can be used for broadcasts

- *A reduction may be done via any broadcast tree with the same communication cost, with reverse data flow*

$$T_{\text{reduce}} = T_{\text{broadcast}} + \text{cost of local reduction work}$$

Scatter is strictly easier than broadcast, pipeline half message to each child in a binary tree

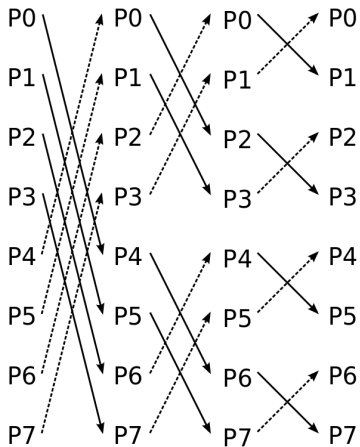
$$T_{\text{scatter}}^{\alpha, \beta}(s, P) \approx 2 \log(P) \cdot \alpha + s \cdot \beta$$

- *A gather may be done via the reverse of any scatter algorithm:*

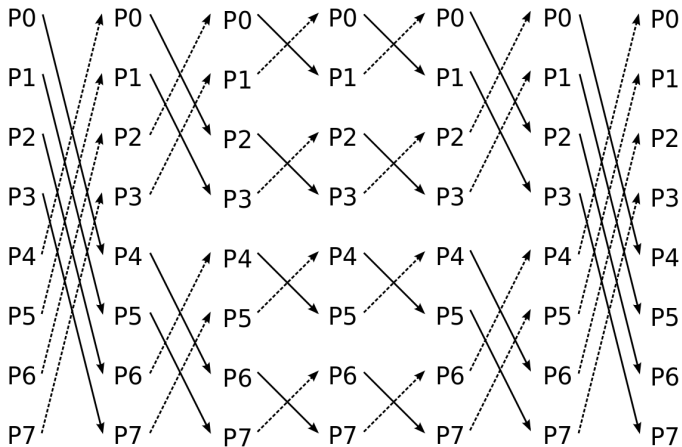
$$T_{\text{gather}} = T_{\text{scatter}}$$

**All-to-One + One-to-All** collectives can be done via two trees, but is this most efficient? What about **All-to-All** collectives?

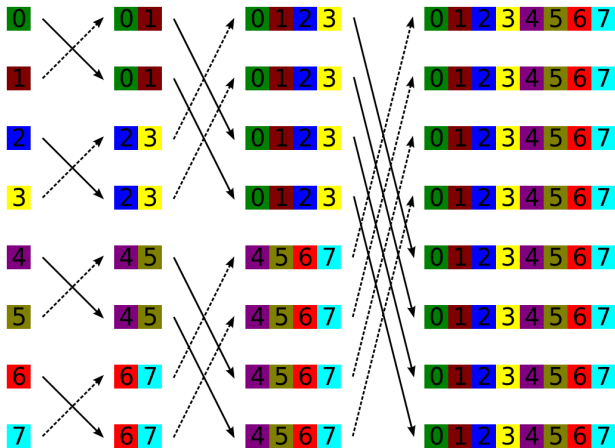
# Butterfly network



# Butterfly network



# Butterfly Allgather (recursive doubling)



## Cost of butterfly Allgather

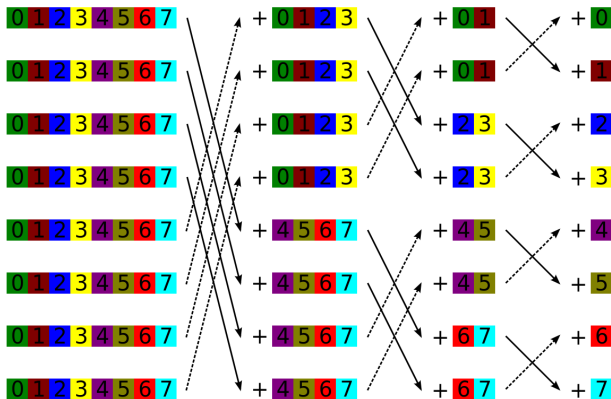
The butterfly has  $\log(P)$  levels. The size of the message doubles at each level until all  $s$  elements are gathered, so the total cost is

$$\begin{aligned}
 T_{\text{allgather}}^{\alpha,\beta}(s, P) &= \begin{cases} 0 & : P = 1 \\ T_{\text{allgather}}^{\alpha,\beta}(s/2, P/2) + \alpha + (s/2) \cdot \beta & : P > 1 \end{cases} \\
 &\approx \log(P) \cdot \alpha + \sum_{i=1}^{\log(P)} s/2^i \cdot \beta \\
 &\approx \log(P) \cdot \alpha + s \cdot \beta
 \end{aligned}$$

The geometric summation in the cost is characteristic of one-to-all, all-to-one, and all-to-one-to-all butterfly algorithms

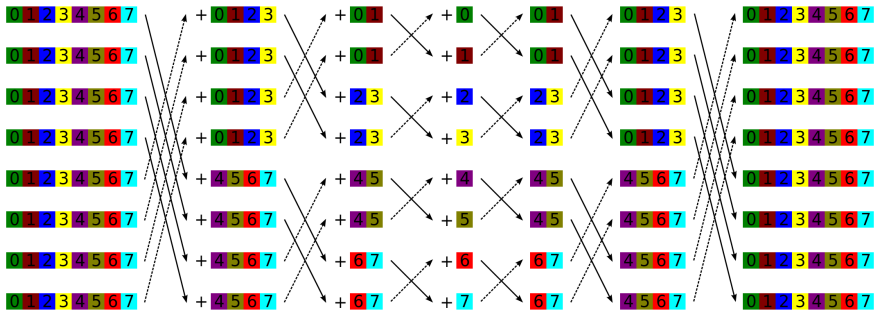
- no pipelining necessary to achieve linear bandwidth cost

# Butterfly Reduce-Scatter (recursive halving)

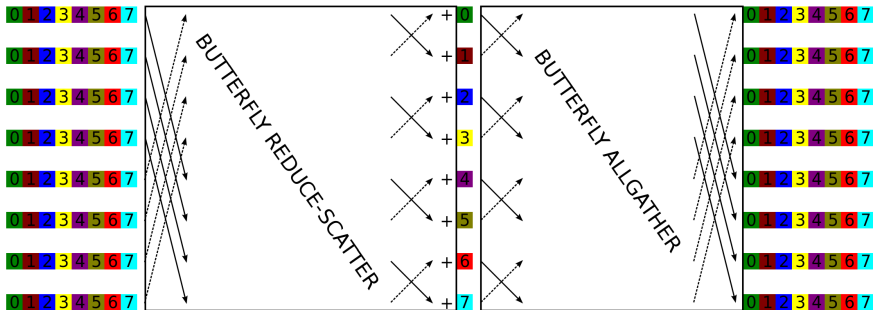


$$T_{\text{reduce-scatter}} = T_{\text{allgather}} + \text{cost of local reduction work}$$

# Butterfly Allreduce



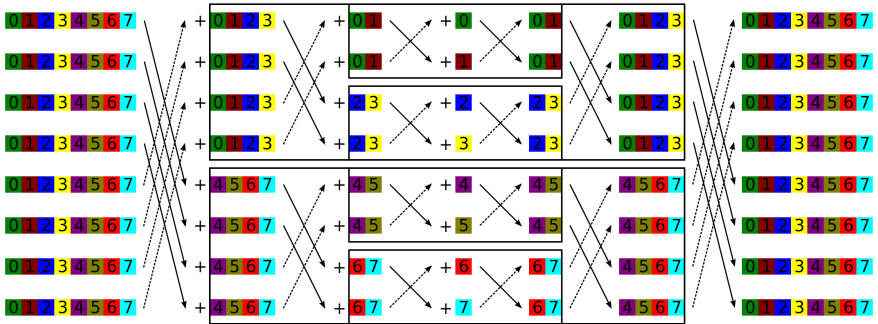
# Butterfly Allreduce



$$T_{\text{allreduce}} = T_{\text{reduce-scatter}} + T_{\text{allgather}}$$

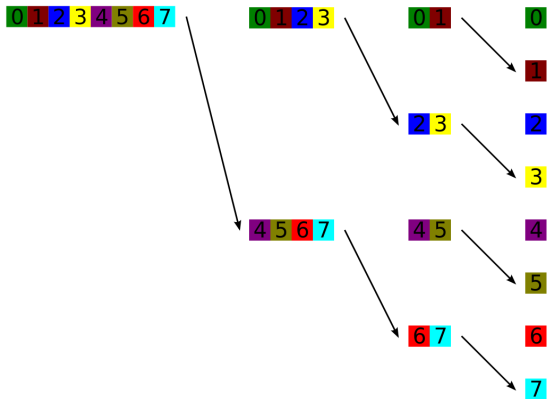


# Butterfly Allreduce: note recursive structure of butterfly



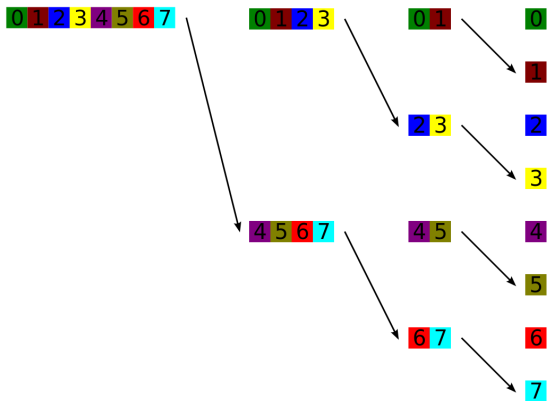
Its possible to do Scan (each processor ends up with a unique value of a prefix sum rather than the full sum) in a similar fashion, but also with operator application done additionally during recursive doubling (Allgather)

# Butterfly Scatter



Question: Which tree is this equivalent to?

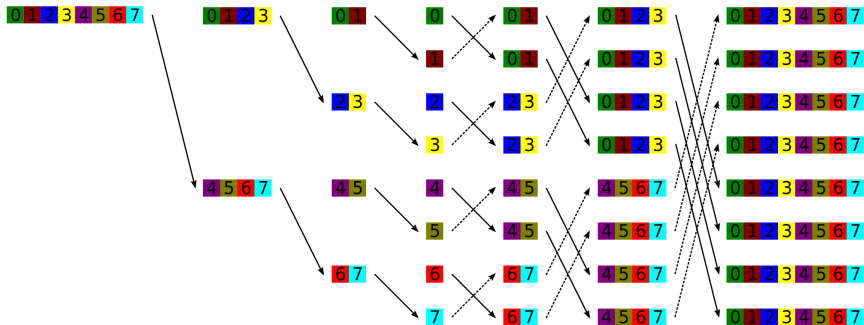
# Butterfly Scatter



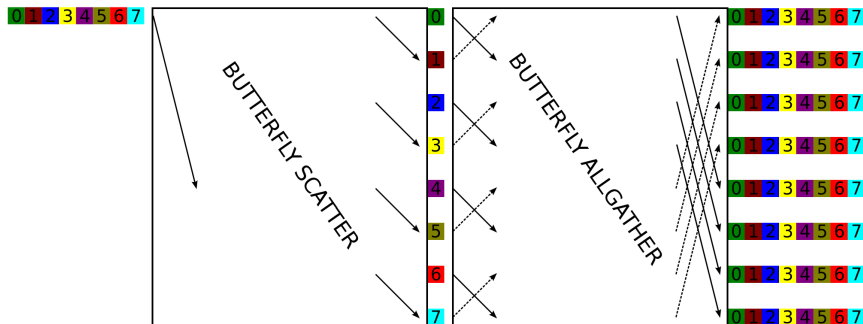
Question: Which tree is this equivalent to?

Answer: Binomial tree.

# Butterfly Broadcast

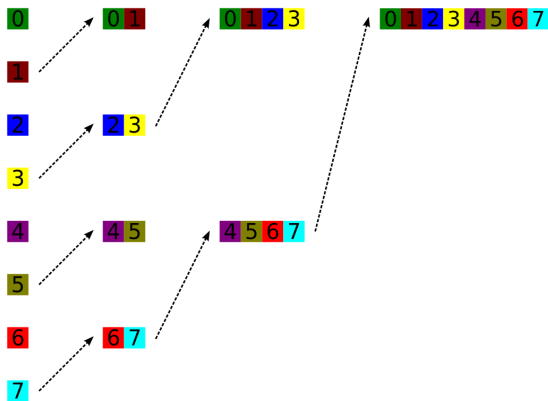


# Butterfly Broadcast



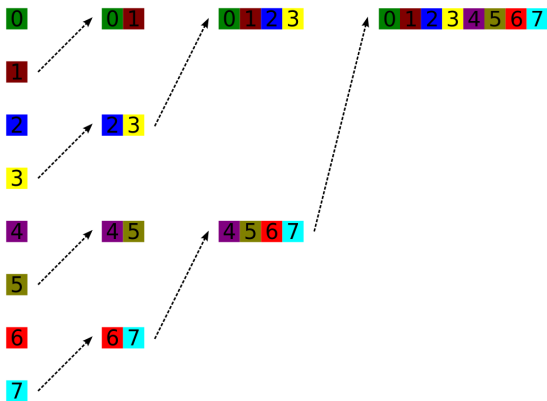
$$T_{\text{broadcast}} = T_{\text{scatter}} + T_{\text{allgather}}$$

# Butterfly Gather



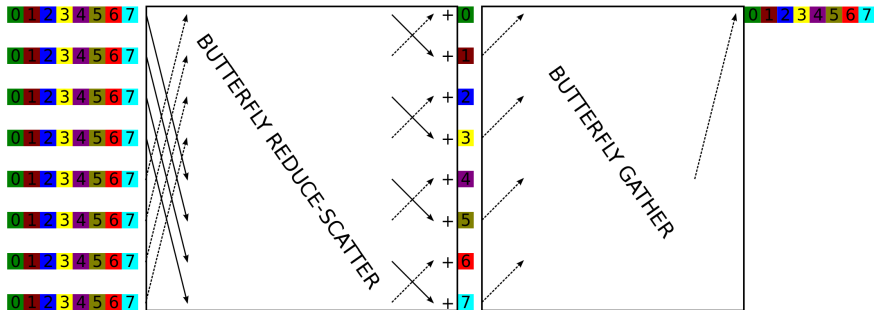
Question: Which other collective could use Gather as a subroutine?

# Butterfly Gather



Question: Which other collective could use Gather as a subroutine?  
Answer: Reduction.

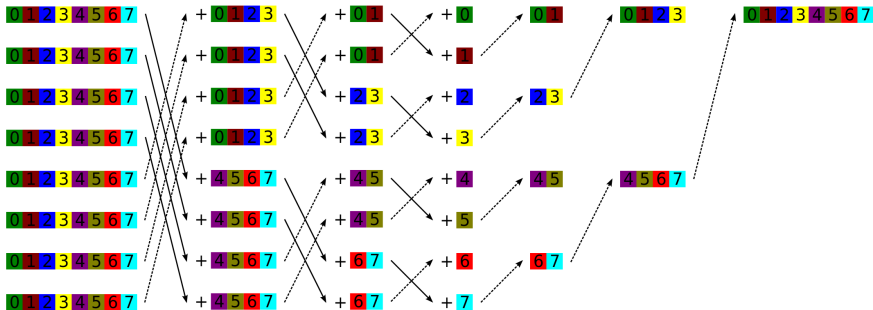
# Butterfly Reduce



$$T_{\text{reduce}} = T_{\text{reduce-scatter}} + T_{\text{gather}}$$

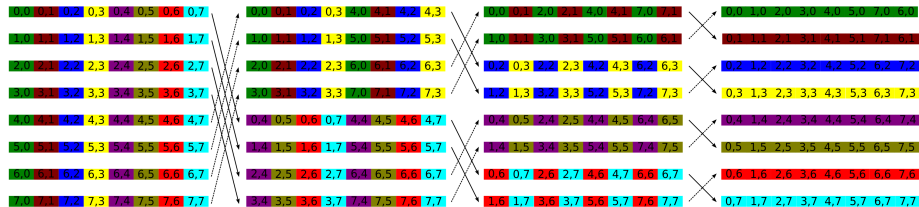


# Butterfly Reduce



$$T_{\text{reduce}} = T_{\text{reduce-scatter}} + T_{\text{gather}}$$

# Butterfly All-to-All



Note that the size of the message stays the same at each level

$$\begin{aligned}
 T_{\text{all-to-all}}^{\alpha, \beta}(s, P) &= \begin{cases} 0 & : P = 1 \\ T_{\text{all-to-all}}^{\alpha, \beta}(s, P/2) + \alpha + (s/2) \cdot \beta & : P > 1 \end{cases} \\
 &= \alpha \cdot \log(P) + \beta \cdot \sum_{i=1}^{\log(P)} s/2 = \alpha \cdot \log(P) + \beta \cdot s/2 \cdot \log(P)
 \end{aligned}$$

Its possible to do All-to-All in less bandwidth cost (as low as  $\beta \cdot s$  by sending directly to targets) at the cost of more messages (as high as  $\alpha \cdot P$  if sending directly)

## BSP model definition

The **Bulk Synchronous Parallel (BSP) model** (Valiant 1990) is a theoretical execution/cost model for parallel algorithms

- execution is subdivided into **supersteps**, each associated with a global synchronization
- within each superstep each processor can send and receive up to  $h$  messages (called an  **$h$ -relation**)
- the cost of sending or receiving  $h$  messages of size  $m$  is  $h \cdot m \cdot \hat{g}$
- the total cost of a superstep is the max over all processors at that superstep
- when  $h = 1$  the BSP model is closely related to the  $\alpha$ - $\beta$  model with  $\beta = \hat{g}$  and LogGP mode with  $G = \hat{g}$
- we will focus on a variant of BSP with  $h = P$  and for consistency refer to  $\hat{g}$  as  $\beta$  and the cost of a synchronization as  $\alpha$

## Synchronization vs latency

By picking  $h = P$ , we allow a global barrier to execute in the same time as the point-to-point latency

- this abstraction is good if the algorithm's performance is not expected to be latency-sensitive
- messages become non-blocking, but progress must be guaranteed by barrier
- collectives can be done in linear bandwidth cost with  $O(1)$  supersteps
- enables high-level algorithm development: how many collective protocols does the algorithm need to execute?
- global barrier may be a barrier of a subset of processors, if BSP is used recursively

## Nonblocking communication

The paradigm of sending non-blocking messages then synchronizing later is sensible

- MPI provides non-blocking 'I(send/rcv)' primitives that may be 'Wait'ed on in bulk (these are slightly slower than blocking primitives, due to buffering)
- MPI and other communication frameworks also provide **one-sided** messaging primitives which are *non-blocking and zero-copy* (no buffering)
- one-sided communication progress must be guaranteed by a barrier on all or a subset of processors (or MPI Win Flush between a pair)

## (Reduce-)Scatter and (All)Gather in BSP

When  $h = P$  all discussed collectives that require a single butterfly can be done in time  $T_{\text{butterfly}} = \alpha + s \cdot \beta$  i.e. they can all be done in one superstep

- Scatter: root sends each message to its target (root incurs  $s \cdot \beta$  send bandwidth)
- Reduce-Scatter: each processor sends its portion to every other processor (every processor incurs  $s \cdot \beta$  send and receive bandwidth)
- Gather: send each message to root (root incurs  $s \cdot \beta$  receive bandwidth)
- Allgather: each processor sends its portion to every other processor (every processor incurs  $s \cdot \beta$  send and receive bandwidth)

when  $h < P$ , we could perform the above algorithms using a butterfly with 'radix'= $h$  (number of neighbors at each butterfly level) in time

$$T_{\text{butterfly}} = \log_h(P) \cdot \alpha + s \cdot \beta$$

## Other collectives in BSP

The Broadcast, Reduce, and Allreduce collectives may be done as combinations of collectives in the same way as with Butterfly algorithms, using two supersteps

- Broadcast done by Scatter then Allgather
- Reduce done by Reduce-Scatter then Gather
- Allreduce done by Reduce-Scatter then Allgather

BSP preserves this hierarchical algorithmic structure and costs.

However, BSP with  $h = P$  can do all-to-all in  $O(s)$  bandwidth and  $O(1)$  supersteps (as cheap as other collectives), when  $h < P$ , the logarithmic factor on the bandwidth is recovered.

## Systems for one-sided communication

BSP employs the concept of non-blocking communication, which presents practical challenges

- to avoid buffering or additional latency overhead, the communicating processor must know be aware of the desired buffer location of the remote processor
- if the location of the remote buffer is known, the communication is called 'one-sided'
- with network hardware known as Remote Direct Memory Access (RDMA) one-sided communication can be accomplished without disturbing the work of the remote processor

One-sided communication transfers are commonly be formulated as

- **Put** – send a message to a remote buffer
- **Get** – receive a message from a remote buffer



# Partitioned Global Address Space (PGAS)

**PGAS** programming models facilitate non-blocking remote memory access

- they allow declaration of buffers in a globally-addressable space, which other processors can access remotely
- **Unified Parallel C (UPC)** is a compiler-based PGAS language that allows direct indexing into globally-distributed arrays (Carlson et al. 1999)
- **Global Arrays** (Nieplocha et al. 1994) is a library that supports a global address space via a one-sided communication layer (e.g. ARMCI, Nieplocha et al. 1999)
- MPI supports one-sided communication via declaration of **windows** that declare remotely-accessible buffers

## Matrix multiplication

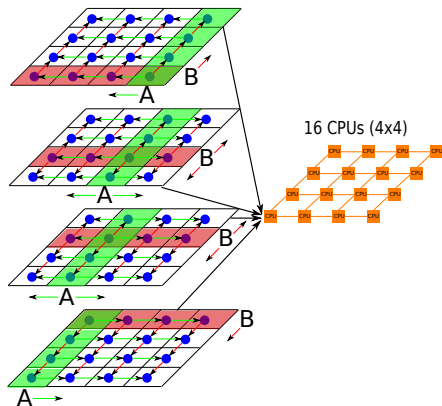
Matrix multiplication of  $n$ -by- $n$  matrices  $A$  and  $B$  into  $C$ ,  $C = A \cdot B$  is defined as, for all  $i, j$ ,

$$C[i, j] = \sum_k A[i, k] \cdot B[k, j]$$

A standard approach to parallelization of matrix multiplication is commonly referred to as **SUMMA** (Agarwal et al. 1995, Van De Geijn et al. 1997), which uses a 2D processor grid, so blocks  $A_{lm}$ ,  $B_{lm}$ , and  $C_{lm}$  are owned by processor  $\Pi[l, m]$

- SUMMA variant 1: iterate for  $k = 1$  to  $\sqrt{P}$  and for all  $i, j \in [1, \sqrt{P}]$ 
  - broadcast  $A_{ik}$  to  $\Pi[i, :]$
  - broadcast  $B_{kj}$  to  $\Pi[:, j]$
  - compute  $C_{ij} = C_{ij} + A_{ik} \cdot B_{kj}$  with processor  $\Pi[i, j]$

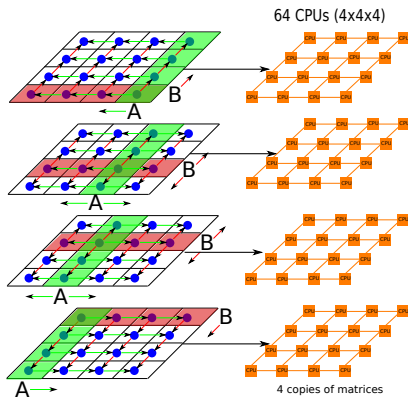
## SUMMA algorithm



$$T_{\text{SUMMA}}^{\alpha, \beta} = 2\sqrt{P} \cdot T_{\text{broadcast}}^{\alpha, \beta}(n^2/p, \sqrt{P}) \leq 2\sqrt{P} \cdot \log(P) \cdot \alpha + \frac{4n^2}{\sqrt{P}} \cdot \beta$$

# 3D Matrix multiplication algorithm

Reference: Agarwal et al. 1995 and others



$$\begin{aligned}
 T_{3D-MM}^{\alpha,\beta} &= 2T_{\text{broadcast}}^{\alpha,\beta}(n^2/p^{2/3}, p^{1/3}) + T_{\text{reduce}}^{\alpha,\beta}(n^2/p^{2/3}, p^{1/3}) \\
 &\leq 2 \log(P) \cdot \alpha + \frac{6n^2}{p^{2/3}} \cdot \beta
 \end{aligned}$$

## LU factorization

The LU factorization algorithm provides a stable (when combined with pivoting) replacement for computing the inverse of a  $n$ -by- $n$  matrix  $A$ ,

$$A = L \cdot U$$

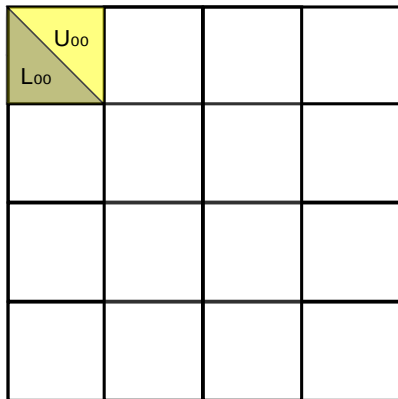
where  $L$  is lower-triangular and  $U$  is upper-triangular is computed via Gaussian elimination: for  $k = 1$  to  $n$ ,

- set  $L[k, k] = 1$  and  $U[k, k : n] = A[k, k : n]$
- divide  $L[k+1 : n, k] = A[k+1 : n, k] / U[k, k]$
- update Schur complement

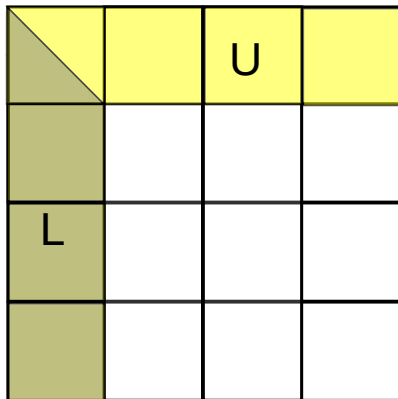
$$A[k+1 : n, k+1 : n] = A[k+1 : n, k+1 : n] - L[k+1 : n, k] \cdot U[k, k+1 : n]$$

this algorithm can be blocked analogously to matrix multiplication

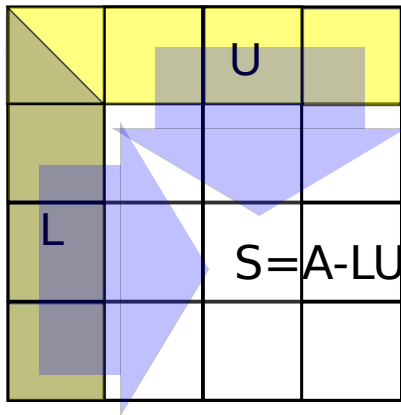
# Blocked LU factorization



# Blocked LU factorization

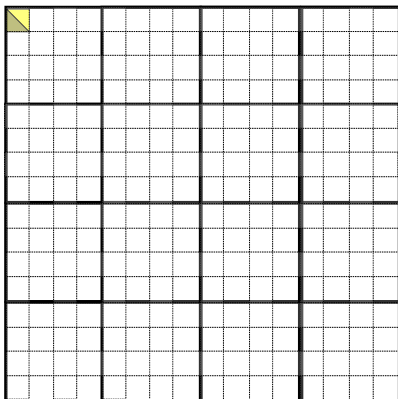


# Blocked LU factorization

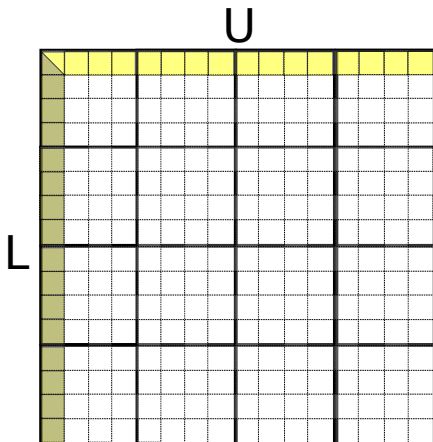




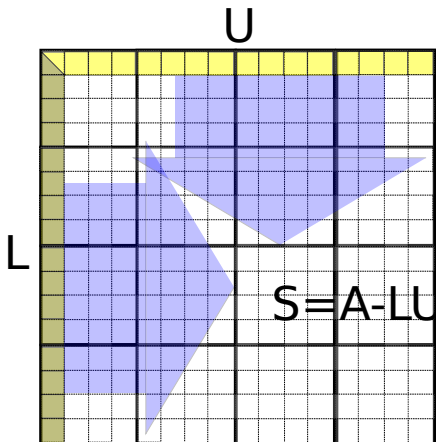
# Block-cyclic LU factorization



# Block-cyclic LU factorization



# Block-cyclic LU factorization



## Recursive matrix multiplication

Now lets consider a recursive parallel algorithm for matrix multiplication

$$\begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \cdot \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

$$C_{11} = A_{11} \cdot B_{11} + A_{12} \cdot B_{21}$$

$$C_{21} = A_{21} \cdot B_{11} + A_{22} \cdot B_{21}$$

$$C_{12} = A_{11} \cdot B_{12} + A_{12} \cdot B_{22}$$

$$C_{22} = A_{21} \cdot B_{12} + A_{22} \cdot B_{22}$$

This requires 8 recursive calls to matrix multiplication of  $n/2$ -by- $n/2$  matrices, as well as matrix additions at each level, which can be done in linear time

## Recursive matrix multiplication: analysis

If we execute all 8 recursive multiplies in parallel with  $P/8$  processors, we obtain a cost recurrence of

$$T_{\text{MM}}^{\alpha,\beta}(n, P) = T_{\text{MM}}^{\alpha,\beta}(n/2, P/8) + O(\alpha) + O\left(\frac{n^2}{P} \cdot \beta\right)$$

The bandwidth cost is dominated by the base cases, where it is proportionate to

$$\left(n/2^{\log_8(P)}\right)^2 = \left(n/P^{\log_8(2)}\right)^2 = \left(n/P^{1/3}\right)^2 = n^2/P^{2/3}$$

for a total that we have seen before (3D algorithm)

$$T_{\text{MM}}^{\alpha,\beta}(n, P) = O(\log(P) \cdot \alpha) + O\left(\frac{n^2}{P^{2/3}} \cdot \beta\right)$$

## Recursive LU factorization

LU factorization has the form

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \cdot \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix}$$

and can be computed recursively via

$$\begin{aligned} [L_{11}, U_{11}] &= \text{LU}(A_{11}) \\ L_{21} &= A_{21} \cdot U_{11}^{-1} \\ U_{12} &= L_{11}^{-1} \cdot A_{12} \\ [L_{22}, U_{22}] &= \text{LU}(A_{22} - L_{21} \cdot U_{12}) \end{aligned}$$

The inverses  $L_{11}^{-1}$  and  $U_{11}^{-1}$  may be obtained as part of the recursion in the first step (see Tiskin 2002 for details). There are two recursive calls to LU and 3 matrix multiplications needed at each step

## Recursive LU factorization: analysis

The two recursive calls within LU factorization must be done in sequence, so we perform them with all the processors. We have to also pay for the cost of matrix multiplications at each level

$$\begin{aligned}T_{\text{LU}}^{\alpha,\beta}(n, P) &= 2T_{\text{LU}}^{\alpha,\beta}(n/2, P) + O(T_{\text{MM}}^{\alpha,\beta}(n, P)) \\ &= 2T_{\text{LU}}^{\alpha,\beta}(n/2, P) + O\left(\log(P) \cdot \alpha + \frac{n^2}{P^{2/3}} \cdot \beta\right)\end{aligned}$$

with base-case cost (sequential execution)

$$T_{\text{LU}}^{\alpha,\beta}(n_0, P) = O(\log(P) \cdot \alpha) + n_0^2 \cdot \beta$$

the bandwidth cost goes down at each level and we can execute the base-case sequentially when  $n_0 = n/P^{2/3}$ , with a total cost of

$$T_{\text{LU}}^{\alpha,\beta}(n, P) = O(P^{2/3} \cdot \log(P) \cdot \alpha) + O\left(\frac{n^2}{P^{2/3}} \cdot \beta\right)$$

## Conclusion and summary

### Summary:

- important parallel communication models:  $\alpha$ - $\beta$ , LogP, LogGP, BSP
- collective communication: binomial trees are good for small-messages, pipelining and/or butterfly needed for large-messages
- collective protocols provide good building blocks for parallel algorithms
- recursion is a thematic approach in communication-efficient algorithms



## Backup slides