



Department of Mathematical
& Statistical Sciences

UNIVERSITY OF COLORADO **DENVER**



Greedy Trees for MPI Reductions

Julien Langou

University of Colorado Denver

MARCH 28, 2013



Table of Contents

Introduction and Model

Unidirectional

Nonuniform Segmentation

Bidirectional

Conclusion



Table of Contents

Introduction and Model

Unidirectional

Nonuniform Segmentation

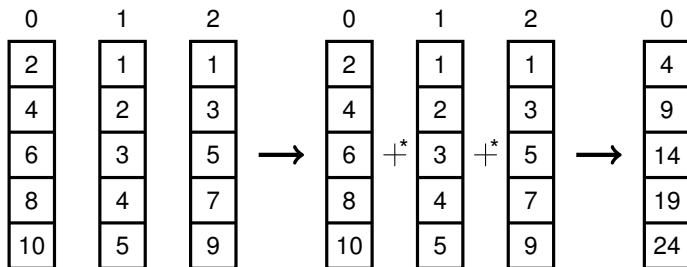
Bidirectional

Conclusion



What is a reduction?

- We will consider a set of p processors with distributed memory and each processor has a message of size m .
- A reduction combines the messages entry-wise, and returns the value on one specified processor.
- Example: $p = 3$, $m = 5$



* Can be any associative operation.



Communication Model

- Unidirectional system - At any given time a processor is allowed to send a message to another processor or receive a message from another processor, but not both.
- Communication time between two processors is given by the linear model, $\alpha + \beta m$, where α is the latency (start up time) and β is the inverse bandwidth.
- The time for the computation is given γm .
- Cannot overlap communication and computation.
- In practice we have the relationship, $\alpha \gg \beta > \gamma$.
- The message m can be split into q segments of size s_i .
- Uniform segmentation: $s_i = s$ for all i .



Table of Contents

Introduction and Model

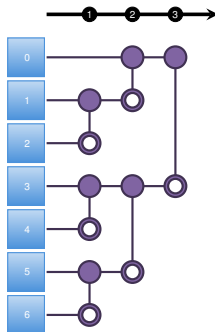
Unidirectional

Nonuniform Segmentation

Bidirectional

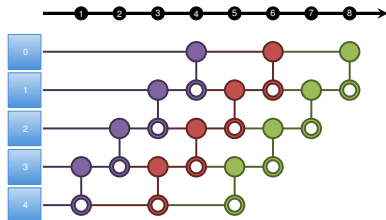
Conclusion

Binomial Tree



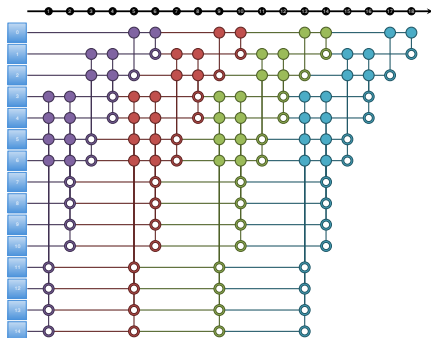
- Best for small messages, $\alpha \gg \beta m$.
- Minimizes the number of communications started.
- No segmentation. Only increases latency.
- $time = \lceil \log_2 p \rceil (\alpha + \beta m + \gamma m)$

Pipeline Tree



- Best for long messages, $\alpha \ll \beta m$.
- Poor startup, but optimal overhead for trailing segments.
- $time = ((p - 1) + 2(q - 1))(\alpha + \beta s + \gamma s)$

Binary Tree



- At each step, two different processors send to the same receiving processor.
- An iteration is therefore twice as long as compared to the other trees.
- Good for medium sized messages, $\alpha \approx \beta m$.
- $time = (2(\lceil \log_2 p + 1 \rceil - 1) + 4(q - 1))(\alpha + \beta s + \gamma s)$

Time Complexity for Binomial, Pipeline, and Binary

Binomial	Time	$\lceil \log_2 p \rceil (\alpha + \beta m + \gamma m)$
Pipeline	Time	$(p-1)(\alpha + \beta s + \gamma s) + 2(q-1)(\alpha + \beta s + \gamma s)$
	s_{opt}	$\left(\frac{2m\alpha}{(p-3)(\beta + \gamma)} \right)^{1/2}$
	T_{opt}	$\left[((p-3)\alpha)^{1/2} + (2m(\beta + \gamma))^{1/2} \right]^2$
Binary	Time	$2(\lceil \log_2(p+1) \rceil - 1)(\alpha + \beta s + \gamma s) + 4(q-1)(\alpha + \beta s + \gamma s)$
	s_{opt}	$\left(\frac{2m\alpha}{(N-3)(\beta + \gamma)} \right)^{1/2}$
	T_{opt}	$2 \left[((N-3)\alpha)^{1/2} + (2m(\beta + \gamma))^{1/2} \right]^2$

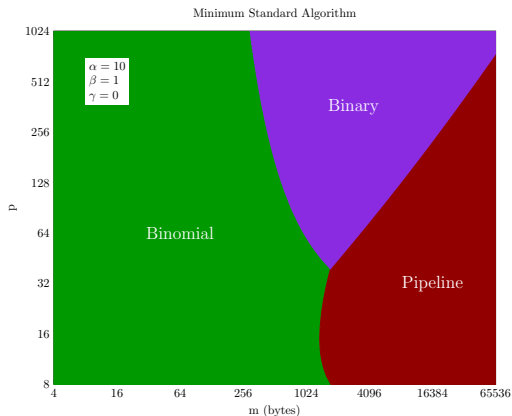
$N = \lceil \log_2(p+1) \rceil$, s_{opt} is the optimal equi-segment size, and T_{opt} is the time for the algorithm at s_{opt} . Formulae are valid for $p > 3$.



Lower Bounds for each term in communication time

	Latency	Bandwidth	Computation
Reduce	$\lceil \log_2 p \rceil \alpha$	$2m\beta$	$\frac{p-1}{p} m\gamma$
Binomial	$\lceil \log_2 p \rceil \alpha$	$\lceil \log_2 p \rceil m\beta$	$\lceil \log_2 p \rceil m\gamma$
Pipeline	$(p-1)\alpha$	$(p-3+2m)\beta$	$(p+2m-3)\gamma$
Binary	$2(N-1)\alpha$	$2(N-3+2m)\beta$	$2(N-3+2m)\gamma$

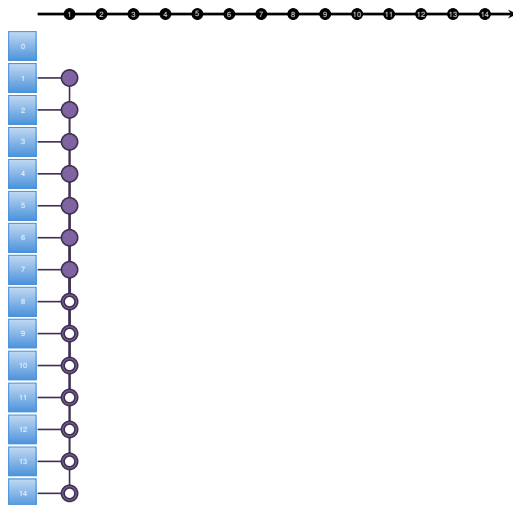
$$*N = \lceil \log_2(p+1) \rceil$$



Regions where binomial or pipeline or binary is better in term of the number of processors (p) and the message size (m). For each algorithm, each p and each m , the optimal segment size is used. The machine parameters are $\alpha = 10$, $\beta = 1$, $\gamma = 0$.



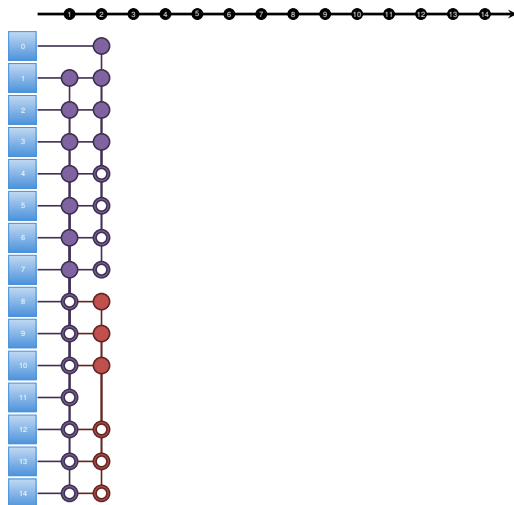
Greedy Tree - Unidirectional, No Computation



- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

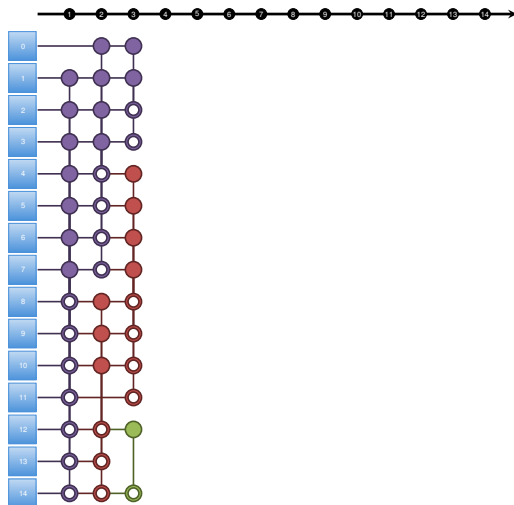


Greedy Tree - Unidirectional, No Computation



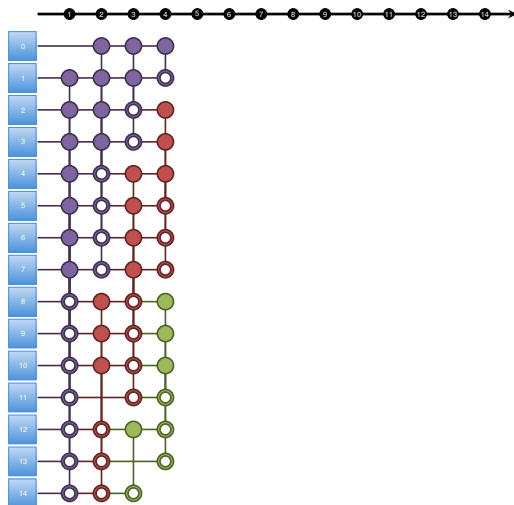
- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional, No Computation



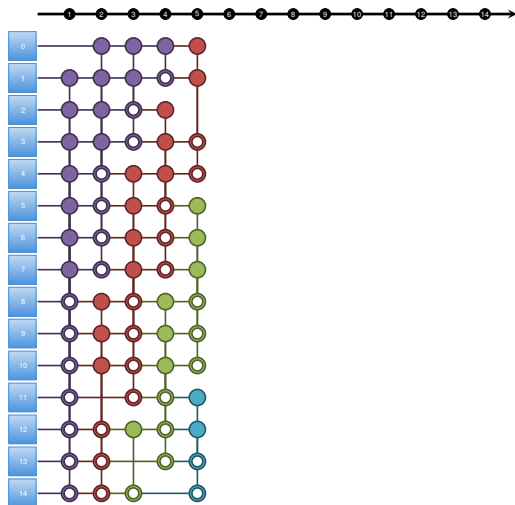
- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional, No Computation



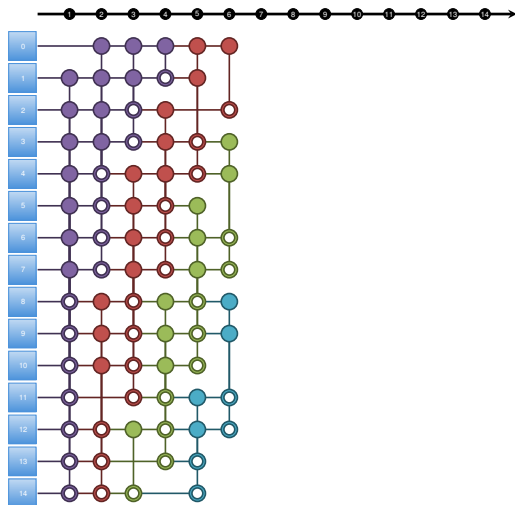
- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional, No Computation



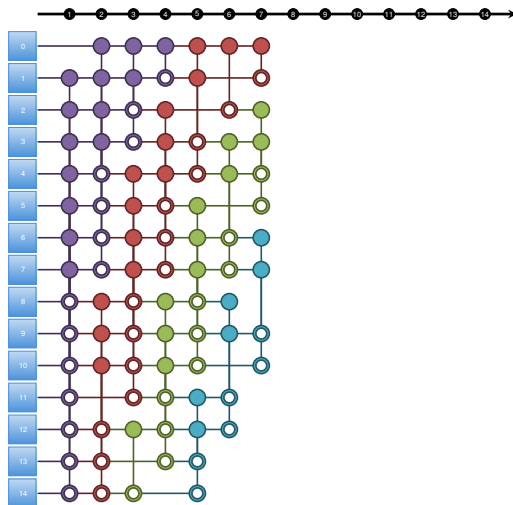
- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional, No Computation



- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

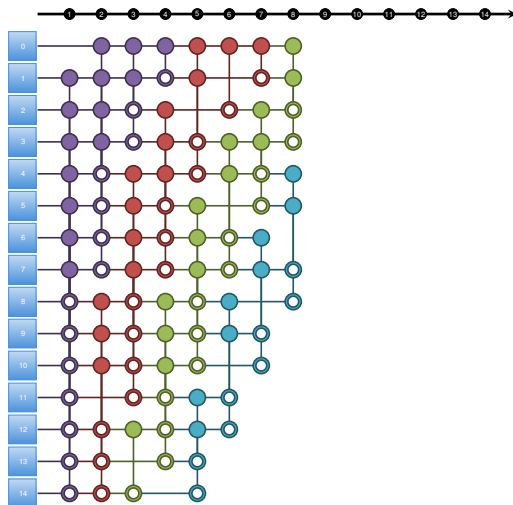
Greedy Tree - Unidirectional, No Computation



- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).



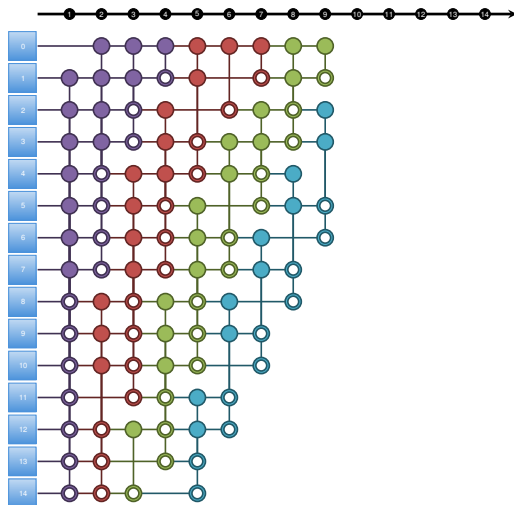
Greedy Tree - Unidirectional, No Computation



- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).



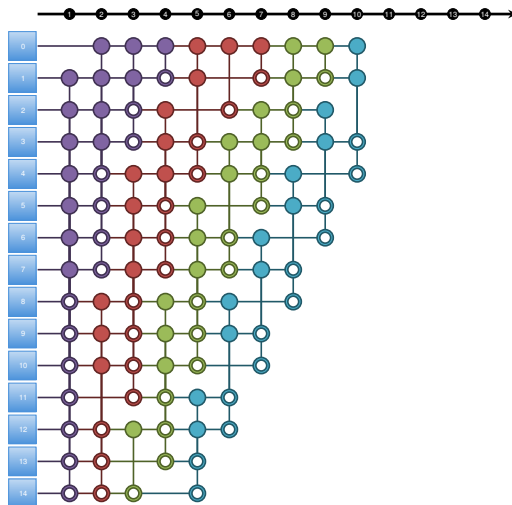
Greedy Tree - Unidirectional, No Computation



- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

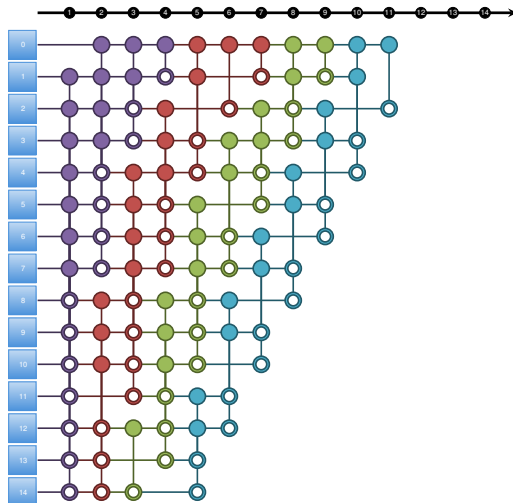


Greedy Tree - Unidirectional, No Computation



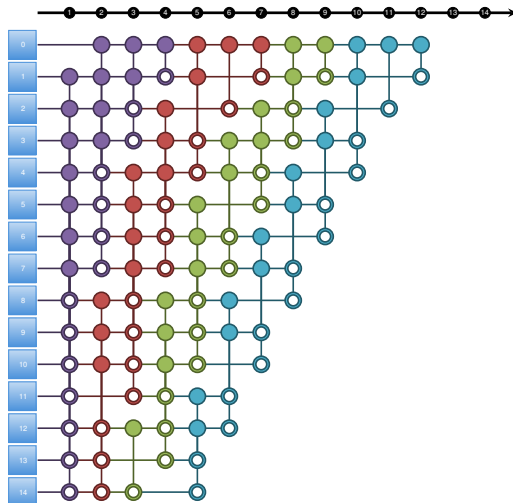
- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional, No Computation



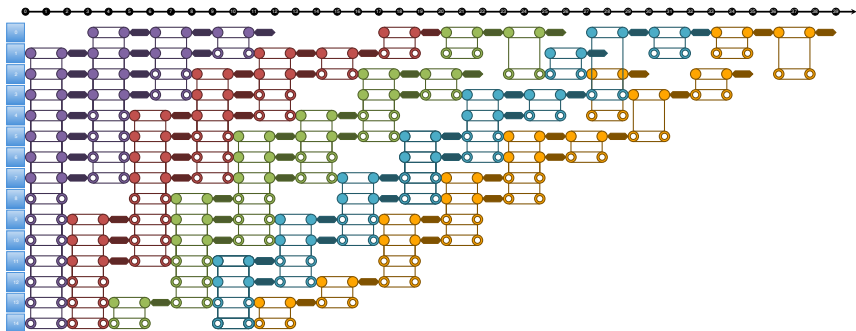
- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional, No Computation



- Optimal for uniform segmentation.
- Motivated by greedy QR factorization scheme. [Cosnard and Robert '86]
- Different tree for each segment.
- We assume the operation is commutative (and associative as well, of course).

Greedy Tree - Unidirectional with Computation



15 processors, 5 equi-segments, $t_{comm} = 2$ and $t_{comp} = 1$.

Filled in circles represent receiving processors, open circles represent sending processors, and hexagons represent computation.



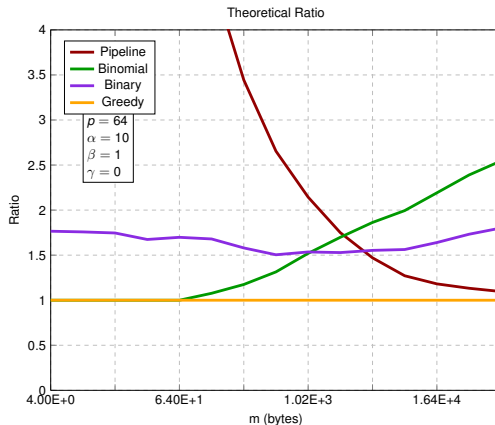
Theoretical Results

Theorem

In an unidirectional system assuming that segments are reduced in order the time complexity of the greedy algorithm is no worse than any reduction algorithm.

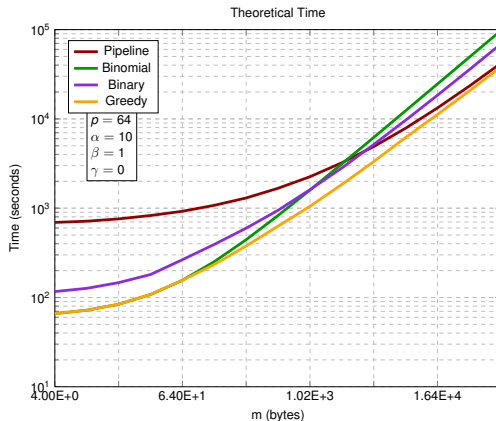
- Segments do not have to be equal size.
- Reducing segments in order allows for less work space (buffer size is $2 \times$ segment size), also this will work for non-commutative operations.
- Each algorithm is tuned for optimal uniform segmentation.
- For given parameters, $p, m, \alpha, \beta, \gamma$, which algorithm is the best?

Theoretical Results (ratio)



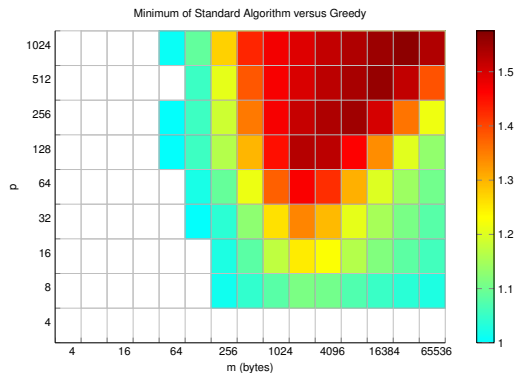
- Each algorithm is tuned for optimal uniform segmentation.
- 50% speed up for “medium sized” messages.

Theoretical Results (time)



- Each algorithm is tuned for optimal uniform segmentation.
- 50% speed up for “medium sized” messages.

Theoretical Results (ratio versus best of “standard algorithms”)



- Each algorithm is tuned for optimal uniform segmentation.
- 50% speed up for “medium sized” messages.

C code for the unidirectional greedy reduction algorithm.

```
#include "reduce.h"
int global_s;
int Reduce_greedy( void *sendbuf_notype, void *recvbuf_notype, int n,
MPI_Datatype mpi_datatype, MPI_Op mpi_op, int root, MPI_Comm mpi_comm){
    int pool_size;
    int my_rank;
    int j, i, q, z, s2, s3, qstart, qend, snext;
    int s=global_s;
    MPI_Status status;
    int *tempbuf, *tempbuf2;
    int *sendbuf, *recvbuf;
    int *hist;

    sendbuf = (int *) sendbuf_notype;
    recvbuf = (int *) recvbuf_notype;

    if (root != 0 ) MPI_Abort( mpi_comm, 512);

    MPI_Comm_size(mpi_comm, &pool_size); MPI_Comm_rank(mpi_comm, &my_rank);

    if(my_rank != 0 ) tempbuf = (int*)malloc(s*sizeof(int));
    tempbuf2 = (int*)malloc(s*sizeof(int));

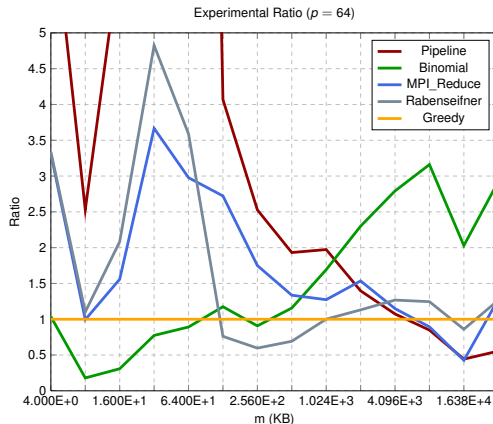
    if ( my_rank == 0 )
        for(i=0;i<n;i++) recvbuf[i] = sendbuf[i];
    else
        for(i=0;i<n;i++) tempbuf[i] = sendbuf[i];

    q = n/s;
    if( (n % s) != 0 ){ s2 = n % s; q++; }
    else s2 = s;

    hist = (int *)malloc(q*sizeof(int));
    for(i = 0; i < q; i++) hist[i] = pool_size;
    qstart = 0;
    qend = 0;
    while( hist[q-1] > 1 ){
        if (qstart != q-1 && hist[qstart+1]-hist[qstart] > 1){ qstart++;
        }
        for(i = qstart; i >= qend; i--){
            z = (i == 0) ? hist[0]/2 : (hist[i]-hist[i-1])/2;
            s3 = ( i == q-1 ) ? s2 : s;
            if( my_rank < hist[i] && my_rank >= hist[i] - z ){
                MPI_Send(tempbuf, s3, mpi_datatype, my_rank-z, 512, mpi_comm );
                if( i < q-1 ){
                    snext = (i == q-1) ? s2 : s;
                    for(j=0;j<snext;j++) tempbuf[j] = sendbuf[(i+1)*s + j];
                }
            }
            if( my_rank < hist[i] - z && my_rank >= hist[i] - 2*z ){
                MPI_Recv( tempbuf2, s3, mpi_datatype, my_rank+z, 512, mpi_comm, &status );
                if ( my_rank == 0 ){
                    for( j = 0; j < s3; j++ ){ recvbuf[(i*s + j) += tempbuf2[j];
                    }
                }
                else
                    for( j = 0; j < s3; j++) tempbuf[j] += tempbuf2[j];
            }
            hist[i] -= z;
            if (hist[i] == 1){
                hist[i] = 0; qend++;
            }
        }
    }
    free( hist ); if(my_rank != 0) free( tempbuf ); free( tempbuf2 );
    return 0;
}
```

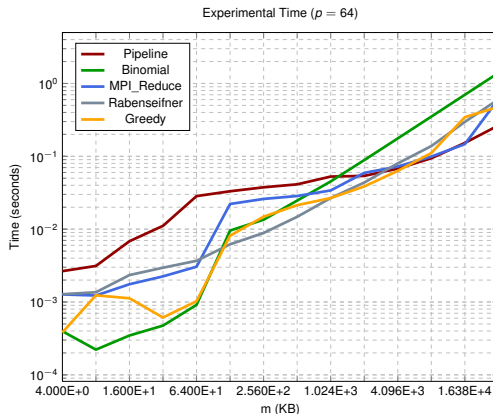
- `global_s` is the size of a segment and needs to be initialized (if possible “tuned”) in advance.
- The implementation is restricted to `root` being 0, `MPI_Datatype` being `int`, and `MPI_Op` being `+`. These restrictions are not a consequence of the algorithm and can be removed.

Experimental Results (ratio)



- Implemented with OpenMPI v1.4.3 point-to-point functions MPI_Send and MPI_Recv.

Experimental Results (time)



- Implemented with OpenMPI v1.4.3 point-to-point functions MPI_Send and MPI_Recv.



Table of Contents

Introduction and Model

Unidirectional

Nonuniform Segmentation

Bidirectional

Conclusion



Nonuniform Segmentation

(Note: our algorithm is optimal no matter the segmentation, uniform or not.)
Why does segmentation have to be uniform?

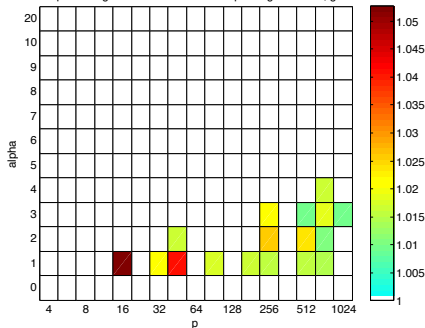
- Experiment: Fix the message size to $m = 10$ and check all possible segmentations.
- Results for greedy
 - 61 of the 986 total trials were optimized by a nonuniform segmentation.
 - The maximum improvement of nonuniform versus uniform segmentation was 7.3%.
 - Of the 61 trials optimized by a nonuniform segmentation the average improvement was 2%.
- For pipeline, all trials were optimized by uniform segmentation.

Nonuniform Segmentation

Sample Segmentations:

Parameters	Percent	Best Uniform	Optimal
$p = 12, \alpha = 0, \beta = 1, \gamma = 1$	7.3%	(1,1,1,1,1,1,1,1,1)	(2,2,1,1,1,1,1,1)
$p = 48, \alpha = 1, \beta = 1, \gamma = 1$	2.6%	(2,2,2,2,2)	(3,2,2,2,1) (2,2,1,2,2,1)
$p = 256, \alpha = 5, \beta = 1, \gamma = 1$	3.0%	(4,4,2)	(5,3,2)

Ratio for optimal segmentation versus best equal segmentation, gamma = 0



Ratio for optimal segmentation versus best equal segmentation, gamma = 1

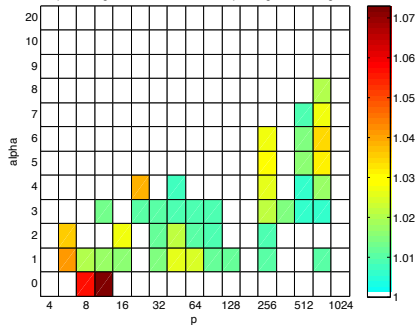




Table of Contents

Introduction and Model

Unidirectional

Nonuniform Segmentation

Bidirectional

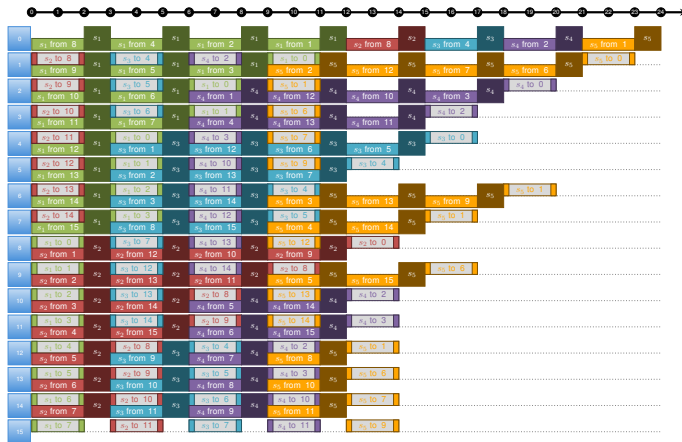
Conclusion



Bidirectional

- Adapt the unidirectional greedy algorithm to the bidirectional context.
- Reduction and broadcast are essentially the same in bidirectional context. (All processors are performing computation.)
- Optimal broadcast algorithm [Traff and Ripke '08, Bar-Noy and Kipnis '94] scheduled in reverse has time $(\lceil \log_2 p \rceil + q - 1)(t_{comm} + t_{comp})$.
- Bidirectional greedy has same time complexity (conjecture).
- Rabenseifner ['04] provides algorithm for optimal computation.

Bidirectional Greedy Tree with computation



16 processors, 5 equi-segments, $t_{comm} = 2$ and $t_{comp} = 1$.

Solid rectangles represent receiving processors and rectangles with end strips represent sending processors. The darker rectangles represent computation.

Pseudo-code for bidirectional greedy algorithm

Algorithm 2: Bidirectional Greedy Algorithm

```

 $S = \text{zeros}(p, 1);$ 
 $R = \text{zeros}(p, 1);$ 
 $C = \text{zeros}(p, 2);$ 
 $t = 0;$ 
 $M = \text{zeros}(p, q);$ 
while  $\min\{M(i, j) \mid 1 \leq i \leq p \text{ and } 1 \leq j \leq q\} = 0$  do
     $\text{segStart} = \min_j\{M(i, j) = 0\};$ 
     $\text{stop} = 1;$ 
     $j = \text{segStart} - 1;$ 
    while  $\text{stop} \neq 0$  do
         $j \leftarrow j + 1;$ 
         $\text{COMP} = \{i \mid C(i, 1) \leq t < C(i, 2) \text{ or } t + t_{\text{comm}}(j) > C(i)\};$ 
         $I = \{i \mid M(i, j) = 0\};$ 
         $\text{sendProc} = I \setminus \text{COMP} \setminus \{i \mid S(i) > t\};$ 
         $\text{recvProc} = I \setminus \text{COMP} \setminus \{i \mid R(i) > t\};$ 
         $\text{freeProc} = \text{sendProc} \cup \text{recvProc};$ 
         $\text{sendProc} \leftarrow \text{sendProc} \setminus \text{freeProc};$ 
         $\text{recvProc} \leftarrow \text{recvProc} \setminus \text{freeProc};$ 
         $s = |\text{sendProc}|;$ 
         $r = |\text{recvProc}|;$ 
         $f = |\text{freeProc}|;$ 
        if  $s = r$  then
             $y = \lfloor f/2 \rfloor;$ 
             $\text{sendProc} \leftarrow \text{sendProc} \cup \{\text{freeProc}(i) \mid f - y + 1 \leq i \leq f\};$ 
             $\text{recvProc} \leftarrow \text{recvProc} \cup \{\text{freeProc}(i) \mid 1 \leq i \leq y\};$ 
        else if  $s < r$  then
             $y = r - s;$ 
             $m = \min\{f, y\};$ 
             $x = \lfloor (f - m)/2 \rfloor;$ 
            if  $m > 0$  then
                 $\text{sendProc} \leftarrow \text{sendProc} \cup \{\text{freeProc}(i) \mid (f - (m + x) + 1) \leq i \leq f\};$ 
            if  $x > 0$  then
                 $\text{recvProc} \leftarrow \text{recvProc} \cup \{\text{freeProc}(i) \mid 1 \leq i \leq x\};$ 
        else if  $r < s$  then
             $y = s - r;$ 
             $m = \min\{f, y\};$ 
             $x = \lfloor (f - m)/2 \rfloor;$ 
            if  $m > 0$  then
                 $\text{recvProc} \leftarrow \text{recvProc} \cup \{\text{freeProc}(i) \mid 1 \leq i \leq m + x\};$ 
            if  $x > 0$  then
                 $\text{sendProc} \leftarrow \text{sendProc} \cup \{\text{freeProc}(i) \mid f - x + 1 \leq i \leq f\};$ 
         $I = \min(|\text{sendProc}|, |\text{recvProc}|);$ 
        if  $I = 0$  then
             $\text{sendProc} \leftarrow \emptyset;$ 
             $\text{recvProc} \leftarrow \emptyset;$ 
        else
             $\text{sendProc} = \{\text{sendProc}(i) \mid 1 \leq i \leq I\};$ 
             $\text{recvProc} = \{\text{recvProc}(i) \mid 1 \leq i \leq I\};$ 
             $M(i, j) = t + t_{\text{comm}}(j), \forall i \text{ s.t. } i \in \text{sendProc};$ 
             $S(i) = t + t_{\text{comm}}(j), \forall i \text{ s.t. } i \in \text{sendProc};$ 
             $R(i) = t + t_{\text{comm}}(j) + t_{\text{comp}}(j), \forall i \text{ s.t. } i \in \text{recvProc};$ 
             $C(i, 1) = t + t_{\text{comm}}(j), \forall i \text{ s.t. } i \in \text{recvProc};$ 
             $C(i, 2) = t + t_{\text{comm}}(j) + t_{\text{comp}}(j), \forall i \text{ s.t. } i \in \text{recvProc};$ 
            if  $|\{i \mid M(i, j) = 0\}| = 1$  then
                 $M(1, j) = \max\{M(i, j) \mid 1 \leq i \leq \text{stop}\} + t_{\text{comp}}(j);$ 
            if  $|\{i \mid S(i) \leq t\}| + |\{i \mid R(i) \leq t\}| < 2$  then
                 $\text{stop} = 0;$ 
            else if  $j \geq q$  then
                 $\text{stop} = 0;$ 
         $t = t + 1;$ 

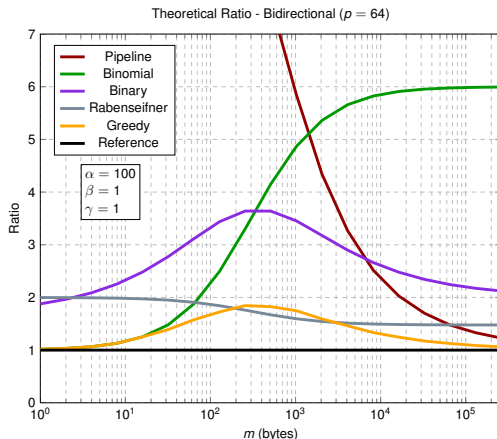
```

Lower bounds for each term in communication time

	Latency	Bandwidth	Computation
Reduce	$\lceil \log_2 p \rceil \alpha$	$m\beta$	$\frac{p-1}{p} m\gamma$
Binomial	$\lceil \log_2 p \rceil \alpha$	$\lceil \log_2 p \rceil m\beta$	$\lceil \log_2 p \rceil m\gamma$
Pipeline	$(p-1)\alpha$	$(p-3+m)\beta$	$(p-3+m)\gamma$
Binary	$2(N-1)\alpha$	$(N-3+2m)\beta$	$(N-3+2m)\gamma$

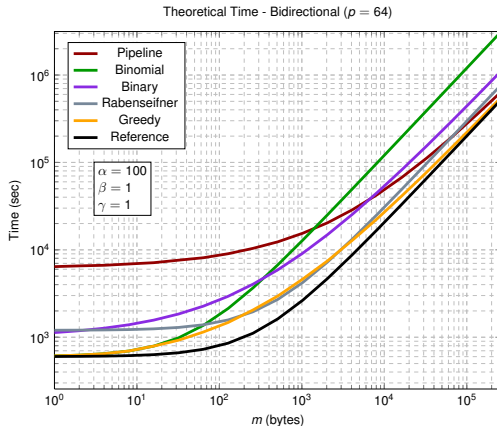
$$*N = \lceil \log_2(p+1) \rceil$$

Theoretical Results (ratio)



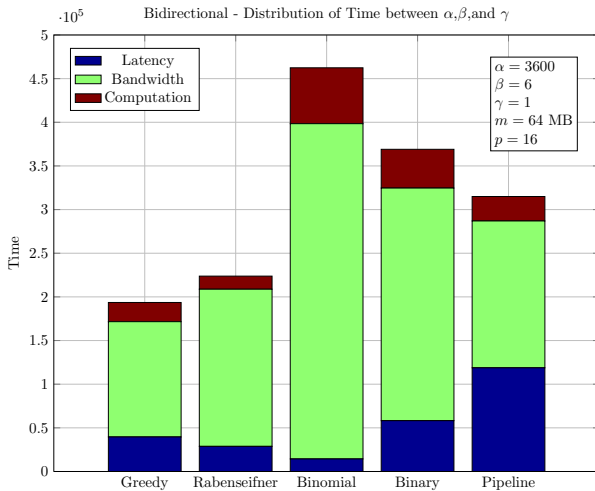
Reference Line: $\lceil \log_2 p \rceil \alpha + m\beta + \frac{p-1}{p} \gamma$

Theoretical Results (time)

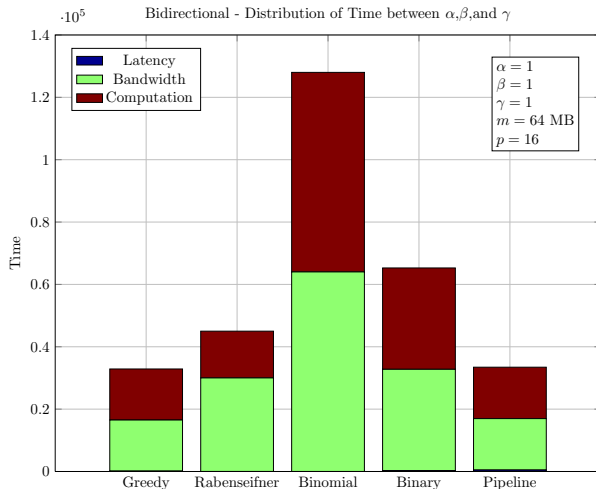


Reference Line: $\lceil \log_2 p \rceil \alpha + m\beta + \frac{p-1}{p} \gamma$

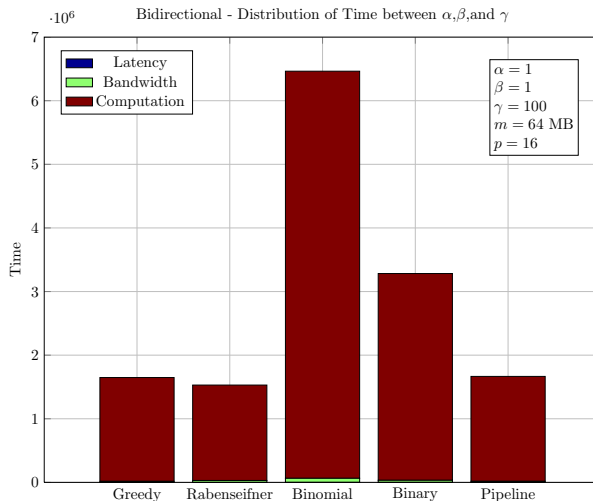
Distribution of time between α, β, γ



Distribution of time between α, β, γ



Distribution of time between α, β, γ





Conclusion

Unidirectional:

- Compared the greedy reduction with three standard algorithms.
- Greedy was the best theoretically.
- Most improvement is for medium sized messages (1Kb - 1Mb).
- Nonuniform segmentation is considered.
- Greedy is optimized by nonuniform segmentation for some machine parameters.

Bidirectional:

- Adapt unidirectional greedy algorithm for a bidirectional system
- Same time complexity of optimal broadcast algorithm.
- Implementation coming soon...

Questions?