

Agenda for today:

- Using Multi-threaded programming
- Message passing:
 - Parallel primitives
 - Reduce
- Map/Reduce

Upcoming Deadlines



Python Programming 2

Not available until Apr 1 | Due Apr 15 at 10am | -/5 pts



Preparation for Lecture 4/15

Not available until Apr 1 | Due Apr 15 at 10am



Assignment 2

Available until Apr 20 | Due Apr 17 at 10am | -/14 pts

DIT 873 & DAT346

Techniques for Large-Scale Data

Lecture 4

Recap

Monte Carlo: Computing Pi

```
p = multiprocessing.Pool(args.workers)
s = p.map(sample_pi, [n]*args.workers)
```

- List of inputs
- Apply the same function to each list element
- Replace list elements by functions output

Here:
only communication
providing argument &
returning result!

Threads writing to memory

Multi-threaded character counting:

Thread #1

the quick brown fox

count['i'] += 1

	...	c	d	e	f	g	h	i	j	k	l	...
Shared count array		9	12	32	2	7	5	6	4	11	9	

count['i'] += 1

Thread #2

multi-threading

Possible Execution

1. Thread 1 reads current value 6 from count array into register
2. Thread 1 increments register
3. Thread 2 reads current value 6 from count array
4. Thread 1 writes incremented value 7
5. Thread 2 increments register
6. Thread 2 writes incremented value 7

⇒ count['i'] incorrect

See Sec. 3.8.3.1 Rauber & Runger (2013)

Threads writing to memory

Multi-threaded character counting:

Thread #1

```
acquire lock count['i']  
count['i'] += 1  
release lock count['i']
```

	...	c	d	e	f	g	h	i	j	k	l	...
Shared count array		9	12	32	2	7	5	6	4	11	9	

Thread #2

```
acquire lock count['i']  
count['i'] += 1  
release lock count['i']
```

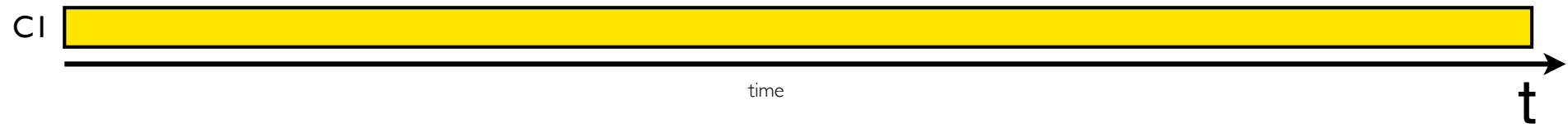
Possible Execution

1. Thread 1 locks count[i]
2. Thread 2 tries to lock count[i] and waits
3. Thread 1 reads current value 6 from count array into register
4. Thread 1 increments register
5. Thread 1 writes incremented value 7
6. Thread 1 releases lock
7. Thread 2 locks count[i]
8. Thread 2 reads current value 7 from count array
9. Thread 2 increments register
10. Thread 2 writes incremented value 8
11. Thread 2 releases lock

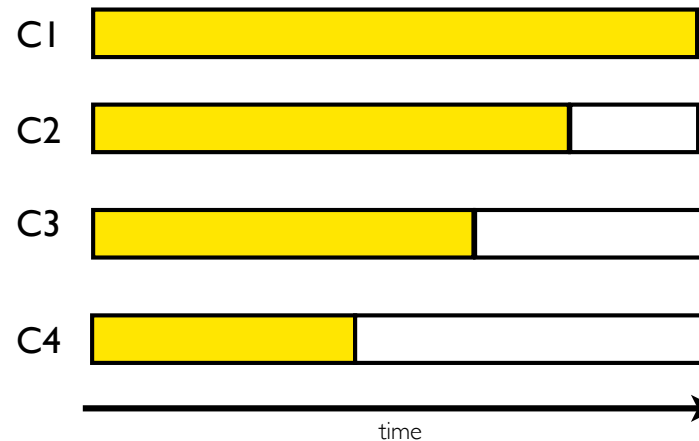
⇒ count['i'] correct

Load Balancing: k-means

Single core



Serial and parallel sections on 4 cores



 Busy

 Idle

Under the hood: Python arrays

- Numpy, Python array module expose the underlying C-arrays
- This allows sharing them in memory
- See `multiprocessing.RawArray`

<https://docs.python.org/2/library/array.html>

Continuing ...

Threads in Python for the lazy

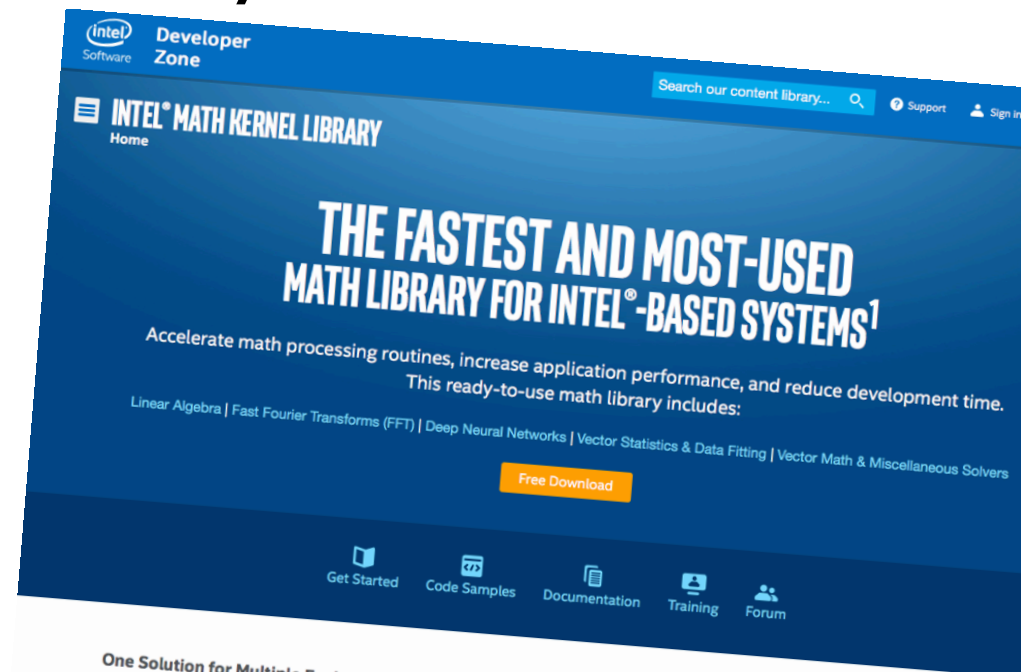
- Intel Math Kernel Library (MKL) is multi-threaded
- used by conda install of numpy
- mostly for linalg, some basic array operations including sum

```
import mkl  
import numpy as np
```

```
mkl.set_num_threads(4)  
x = np.ones((N,N))  
y = np.ones((N,N))  
np.dot(x,y)
```

<https://software.intel.com/en-us/mkl>

<https://software.intel.com/en-us/articles/numpy scipy-with-intel-mkl>



Pitfalls

- Multiprocessing + threaded MKL could lead to reduced performance

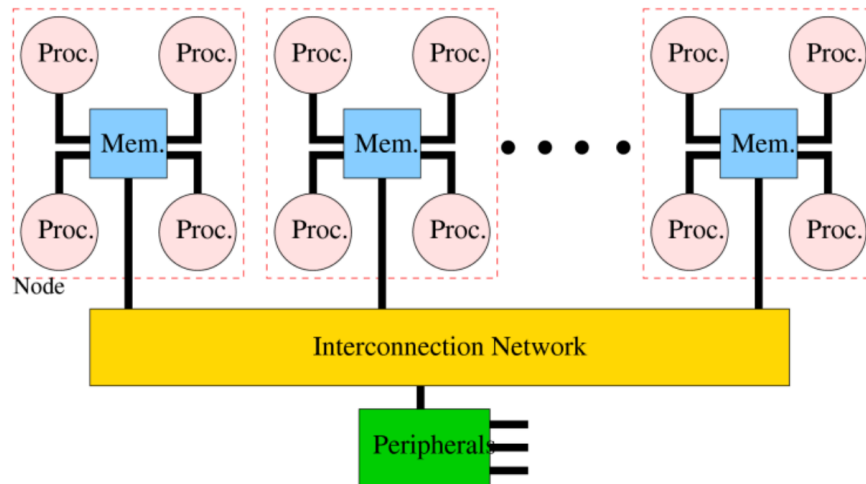
$\#MP \text{ processes} + \#ML \text{ threads} > \#cores$

Parallel programming models

Message Passing

<https://www.scipy.org/topical-software.html#parallel-and-distributed-programming>

Parallel programming models: message passing



- Lower throughput for sharing large amount of data vs. RAM
- Explicit parallelization needed
- MPI is library standard
- Communication incl. exchange of data via interconnect. Infiniband 5 μ s (5000 ns + 2 x RAM access)

Passing a message over Infiniband takes about 5000 ns. This is how many times slower than a memory access?

A	B	C	D	E
1-10	11-50	51-250	251-1000	>1000

Latency

0.5	ns	- CPU L1 dCACHE reference
1	ns	- speed-of-light (a photon) travel a 1 ft (30.5cm) distance
5	ns	- CPU L1 iCACHE Branch mispredict
7	ns	- CPU L2 CACHE reference
71	ns	- CPU cross-QPI/NUMA best case on XEON E5-46*
100	ns	- MUTEX lock/unlock
100	ns	- own DDR MEMORY reference
135	ns	- CPU cross-QPI/NUMA best case on XEON E7-*
202	ns	- CPU cross-QPI/NUMA worst case on XEON E7-*
325	ns	- CPU cross-QPI/NUMA worst case on XEON E5-46*
10,000	ns	- Compress 1K bytes with Zippy PROCESS
20,000	ns	- Send 2K bytes over 1 Gbps NETWORK
250,000	ns	- Read 1 MB sequentially from MEMORY
500,000	ns	- Round trip within a same DataCenter
10,000,000	ns	- DISK seek
10,000,000	ns	- Read 1 MB sequentially from NETWORK
30,000,000	ns	- Read 1 MB sequentially from DISK
150,000,000	ns	- Send a NETWORK packet CA -> Netherlands

			ns
		us	
	ms		

MPI: Message Passing Interface

- See <https://www.open-mpi.org/>
- Runs on large portion of TOP500 super computers
- Mature code base: OpenMPI v1.0 released Nov 2005. MPI 1.0 standard developed 1992-1994
- Python bindings: <http://mpi4py.scipy.org/>



Open MPI v3.0.1 documentation

About
Presentations
Open MPI Team
FAQ
Videos
Performance
Open MPI Software
Download
Documentation
Current
v3.0 (current stable)
Still supported
v2.1 (prior stable)
v2.0 (prior stable)
Older versions
v1.10 (prior stable)
v1.8 (prior stable)
v1.7 (prior feature)
v1.6 (retired)
v1.5 (retired)
v1.4 (ancient)
v1.3 (ancient)
v1.2 (ancient)
v1.1 (ancient)
Source Code Access
Bug Tracking
Regression Testing
Version Information
Sub-Projects
Hardware Locality
Network Locality
MPI Testing Tool
Open MPI User Docs
Open Tool for Parameter
Optimization
Community
Mailing Lists
Getting Help/Support
Contribute
Mirrors
Contact
Licence

Open MPI commands (section 1 man pages)

mpicc++	ompi-dvm	orte-ps	oshmem_info
mpicxx	ompi-ps	orte-server	oshrun
mpicxxc	ompi-server	orte-top	shmemc++
mpicxxc	ompi-top	orted	shmemcc
mpi77	ompi_info	orterun	shmemccx
mpi90	opal_wrapper	oshc++	shmemfort
mpifort	orte-clean	oshcc	shmemnrun
mpinrun	orte-dvm	oshcxc	
ompi-clean	orte-info	oshift	

Open MPI general information (section 7 man pages)

ompi_crcc	orte_file	orte_soapc
opal_crs	orte_hosts	orte_store

MPI API (section 3 man pages)

MPI	MPI_Get_library_version	MPI_T_pvar_start	shmem_float_prod_to
MPIX_Query_cuda_support	MPI_Get_processor_name	MPI_T_pvar_stop	shmem_float_out
MPI_Abort	MPI_Get_version	MPI_T_pvar_stop	shmem_float_out_nbi
MPI_Accumulate	MPI_Graph_create	MPI_Test	shmem_float_set
MPI_Add_error_class	MPI_Graph_get	MPI_Test_cancelled	shmem_float_sum_to
MPI_Add_error_code	MPI_Graph_max	MPI_Testall	shmem_float_swap
MPI_Add_error_string	MPI_Graph_neighbors	MPI_Testany	shmem_free
MPI_Address	MPI_Graph_neighbors_count	MPI_Testsome	shmem_get128
MPI_Aint_add	MPI_Graphdims_get	MPI_Topo_test	shmem_get128_nbi
MPI_Aint_diff	MPI_Group	MPI_Type_c2f	shmem_get16_nbi
MPI_Allgather	MPI_Group_c2f	MPI_Groupquest_start	shmem_get32
MPI_Allgatherv	MPI_Group_compare	MPI_Group_c2f	shmem_get32_nbi
MPI_Allreduce	MPI_Group_difference	MPI_Group_excl	shmem_get64
MPI_Alltoall	MPI_Group_excl	MPI_Group_f2c	shmem_get64_nbi
MPI_Alltoallv	MPI_Group_f2c	MPI_Group_free	shmem_getnm
MPI_Alltoallw	MPI_Group_free	MPI_Group_incl	shmem_getnm_nbi
MPI_Attr_delete	MPI_Group_incl	MPI_Group_intersection	shmem_global_exit
MPI_Attr_get	MPI_Group_intersection	MPI_Group_range_excl	shmem_get128_block
MPI_Attr_put	MPI_Group_range_excl	MPI_Group_range_incl	shmem_get32
MPI_Barrier	MPI_Group_range_incl	MPI_Group_rank	shmem_get64
MPI_Bcast	MPI_Group_rank	MPI_Group_size	shmem_info_get_nam
MPI_Bsend	MPI_Group_size	MPI_Group_translate_ranks	shmem_info_get_vers
MPI_Bsend_init	MPI_Group_union	MPI_Group_delete_attr	shmem_init
MPI_Buffer_attach	MPI_Iallgather	MPI_Type_dup	shmem_int_add
MPI_Buffer_detach	MPI_Iallgatherv	MPI_Type_extnt	shmem_int_and_to_a
MPI_Cancel	MPI_Iallreduce	MPI_Type_f2c	shmem_int_cswap
MPI_Cart_coords	MPI_Ialltoall	MPI_Type_free	shmem_int_fadd
MPI_Cart_create	MPI_Ialltoallv	MPI_Type_get_attr	shmem_int_finc
MPI_Cart_get	MPI_Ibarrier	MPI_Type_get_contents	shmem_int_g
MPI_Cart_map	MPI_Ibcast	MPI_Type_get_envelope	shmem_int_get
MPI_Cart_rank	MPI_Ibcast	MPI_Type_get_extent	shmem_int_get_nbi
MPI_Cart_shift	MPI_Ibcast	MPI_Type_get_extent_x	shmem_int_inc
MPI_Cart_sub	MPI_Ibcast	MPI_Type_get_name	shmem_int_incl
MPI_Cartdim_get	MPI_Ibcast	MPI_Type_get_true_extent	shmem_int_incl_nbi
MPI_Close_port	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_accept	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_c2f	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_call_errhandler	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_compare	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_connect	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_create	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_create_errhandler	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_create_group	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a
MPI_Comm_create_keyval	MPI_Ibcast	MPI_Type_get_true_extent_x	shmem_int_min_to_a

<https://www.open-mpi.org/doc/current/>

MPI_Comm_delete_attr	MPI_Info_create	MPI_Type_size	shmem_int_swap
MPI_Comm_disconnect	MPI_Info_delete	MPI_Type_size_x	shmem_int_wait
MPI_Comm_dup	MPI_Info_dup	MPI_Type_struct	shmem_int_wait_until
MPI_Comm_dup_with_info	MPI_Info_env	MPI_Type_ub	shmem_int_xor_to_all
MPI_Comm_f2c	MPI_Info_f2c	MPI_Type_vector	shmem_iput128
MPI_Comm_free	MPI_Info_free	MPI_Unpack	shmem_iput32
MPI_Comm_free_keyval	MPI_Info_get	MPI_Unpack_external	shmem_iput64
MPI_Comm_get_attr	MPI_Info_get_nkeys	MPI_Unpublish_name	shmem_long_add
MPI_Comm_get_errhandler	MPI_Info_get_nthkey	MPI_Unpublish_name	shmem_long_and_to_all
MPI_Comm_get_info	MPI_Info_get_valuelen	MPI_Wait	shmem_long_cswap
MPI_Comm_get_name	MPI_Info_set	MPI_Waitall	shmem_long_fadd
MPI_Comm_get_parent	MPI_Init	MPI_Waitany	shmem_long_fetch
MPI_Comm_group	MPI_Init_thread	MPI_Win_allocate	shmem_long_finc
MPI_Comm_idup	MPI_Initialized	MPI_Win_allocate_shared	shmem_long_g
MPI_Comm_join	MPI_Intercomm_create	MPI_Win_attach	shmem_long_get
MPI_Comm_rank	MPI_Intercomm_merge	MPI_Win_c2f	shmem_long_get_nbi
MPI_Comm_remote_group	MPI_Iprobe	MPI_Win_call_errhandler	shmem_long_inc
MPI_Comm_remote_size	MPI_Irecv	MPI_Win_complete	shmem_long_inc
MPI_Comm_set_attr	MPI_Ireduce	MPI_Win_create	shmem_long_inc
MPI_Comm_set_errhandler	MPI_Ireduce_scatter	MPI_Win_create_dynamic	shmem_long_max_to_all
MPI_Comm_set_info	MPI_Ireduce_scatter_block	MPI_Win_create_errhandler	shmem_long_min_to_all
MPI_Comm_set_name	MPI_Irsend	MPI_Win_create_keyval	shmem_long_or_to_all
MPI_Comm_size	MPI_Is_thread_main	MPI_Win_delete_attr	shmem_long_p
MPI_Comm_spawn	MPI_Iscan	MPI_Win_fence	shmem_long_prod_to_all
MPI_Comm_spawn_multiple	MPI_Iscatter	MPI_Win_f2c	shmem_long_put
MPI_Comm_split	MPI_Iscatterv	MPI_Win_flush	shmem_long_put_nbi
MPI_Comm_split_type	MPI_Isend	MPI_Win_flush	shmem_long_set
MPI_Comm_test_inter	MPI_Issend	MPI_Win_flush_all	shmem_long_sum_to_all
MPI_Comm_test_swap	MPI_Keyval_create	MPI_Win_flush_local	shmem_long_swap
MPI_Dims_create	MPI_Keyval_free	MPI_Win_flush_local_all	shmem_long_wait
MPI_Dist_graph_create	MPI_Lookup_name	MPI_Win_free	shmem_long_wait_until
MPI_Dist_graph_create_adjacent	MPI_Message_c2f	MPI_Win_free_keyval	shmem_long_xor_to_all
MPI_Dist_graph_neighbors	MPI_Message_f2c	MPI_Win_get_attr	shmem_longdouble_g
MPI_Dist_graph_neighbors_count	MPI_Mprobe	MPI_Win_get_errhandler	shmem_longdouble_get
MPI_Errhandler_create	MPI_Mrecv	MPI_Win_get_group	shmem_longdouble_get_nbi
MPI_Errhandler_free	MPI_Neighbor_allgather	MPI_Win_get_info	shmem_longdouble_inc
MPI_Errhandler_set	MPI_Neighbor_allgatherv	MPI_Win_get_name	shmem_longdouble_iput
MPI_Error_class	MPI_Neighbor_alltoall	MPI_Win_lock	shmem_longdouble_max_to_all
MPI_Error_string	MPI_Neighbor_alltoallv	MPI_Win_lock_all	shmem_longdouble_min_to_all
MPI_Exscan	MPI_Neighbor_alltoallv	MPI_Win_post	shmem_longdouble_p
MPI_File_close	MPI_Op_c2f	MPI_Win_set_attr	shmem_longdouble_prod_to_all
MPI_File_call_errhandler	MPI_Op_commutative	MPI_Win_set_errhandler	shmem_longdouble_put
MPI_File_close	MPI_Op_create	MPI_Win_set_info	shmem_longdouble_put_nbi
MPI_File_delete	MPI_Op_f2c	MPI_Win_set_name	shmem_longdouble_set
MPI_File_f2c	MPI_Op_free	MPI_Win_shared_query	shmem_longdouble_swap
MPI_File_get_amode	MPI_Open_start	MPI_Win_start	shmem_longdouble_wait
MPI_File_get_atomicsity	MPI_Pack	MPI_Win_sync	shmem_longdouble_wait
MPI_File_get_byte_offset	MPI_Pack_external	MPI_Win_test	shmem_longdouble_wait
MPI_File_get_errhandler	MPI_Pack_external_size	MPI_Win_unlock	shmem_longdouble_wait
MPI_File_get_group	MPI_Pack_size	MPI_Win_unlock_all	shmem_longdouble_wait
MPI_File_get_info	MPI_Pcontrol	MPI_Win_wait	shmem_longdouble_wait
MPI_File_get_position	MPI_Probe	MPI_Wtick	shmem_longdouble_wait
MPI_File_get_position_shared	MPI_Publish_name	MPI_Whence	shmem_longdouble_wait
MPI_File_get_size	MPI_Put	QMPI_Affinity_str	shmem_longdouble_wait
MPI_File_get_type_extent	MPI_Query_thread	OpenMPI	shmem_longdouble_wait
MPI_File_get_view	MPI_Raccumulate	OpenSHMEM	shmem_longdouble_wait
MPI_File_inread	MPI_Recv	_my_pe	shmem_longdouble_wait
MPI_File_inread_all	MPI_Recv_init	_num_pes	shmem_longdouble_wait
MPI_File_inread_at	MPI_Reduce	intro_shmem	shmem_longdouble_wait
MPI_File_inread_at_all	MPI_Reduce_local	shfns	shmem_longdouble_wait
MPI_File_inread_shared	MPI_Reduce_scatter	shmem_addr_accessible	shmem_longdouble_wait
MPI_File_iread	MPI_Register_datarg	shmem_align	shmem_longdouble_wait
MPI_File_iread_at	MPI_Request_c2f	shmem_alltoall32	shmem_longdouble_wait
MPI_File_iread_at_all	MPI_Request_f2c	shmem_alltoall64	shmem_longdouble_wait
MPI_File_iread_shared	MPI_Request_free	shmem_alltoall32	shmem_longdouble_wait
MPI_File_iread_at	MPI_Request_get_status	shmem_alltoall64	shmem_longdouble_wait
MPI_File_iread_at_all	MPI_Rget	shmem_barrier	shmem_longdouble_wait
MPI_File_iread_shared	MPI_Rget_accumulate	shmem_barrier_all	shmem_longdouble_wait
MPI_File_open	MPI_Rput	shmem_broadcast32	shmem_longdouble_wait
MPI_File_preallocate	MPI_Rsend	shmem_broadcast64	shmem_longdouble_wait
MPI_File_read	MPI_Rsend_init	shmem_char_g	shmem_longdouble_wait
MPI_File_read_all	MPI_Scan	shmem_char_get	shmem_longdouble_wait
MPI_File_read_all_begin	MPI_Scatter	shmem_char_get_nbi	shmem_longdouble_wait
MPI_File_read_all_end	MPI_Scatterv	shmem_char_p	shmem_longdouble_wait

<https://www.open-mpi.org/doc/current/>

MPI_File_read_at	MPI_Send	shmem_char_put	shmem_put16_nbi
MPI_File_read_at_all	MPI_Send_init	shmem_char_put_nbi	shmem_put32
MPI_File_read_at_all_begin	MPI_Sendrecv	shmem_clear_cache_inv	shmem_put32_nbi
MPI_File_read_at_all_end	MPI_Sendrecv_replace	shmem_clear_cache_line_inv	shmem_put64
MPI_File_read_ordered	MPI_Sizeof	shmem_clear_lock	shmem_put64_nbi
MPI_File_read_ordered_begin	MPI_Ssend	shmem_collect32	shmem_putmem
MPI_File_read_ordered_end	MPI_Ssend_init	shmem_collect64	shmem_putmem_nbi
MPI_File_read_shared	MPI_Start	shmem_complex_prod_to_all	shmem_putmem_nbi
MPI_File_seek	MPI_Startall	shmem_complex_sum_to_all	shmemQuiet
MPI_File_seek_shared	MPI_Status_c2f	shmem_complexx_prod_to_all	shmemQuiet
MPI_File_set_atomicsity	MPI_Status_f2c	shmem_complexx_sum_to_all	shmemQuiet
MPI_File_set_errhandler	MPI_Status_set_cancelled	shmem_double_fetch	shmemQuiet
MPI_File_set_info	MPI_Status_set_elements	shmem_double_get	shmemQuiet
MPI_File_set_size	MPI_Status_set_elements_x	shmem_double_q	shmemQuiet
MPI_File_set_view	MPI_T_category_changed	shmem_double_get_nbi	shmemQuiet
MPI_File_sync	MPI_T_category_get_categories	shmem_double_iput	shmemQuiet
MPI_File_write	MPI_T_category_get_cvars	shmem_double_max_to_all	shmemQuiet
MPI_File_write_all	MPI_T_category_get_info	shmem_double_min_to_all	shmemQuiet
MPI_File_write_all_begin	MPI_T_category_get_num	shmem_double_p	shmemQuiet
MPI_File_write_all_end	MPI_T_category_get_pvars	shmem_double_p	shmemQuiet
MPI_File_write_at	MPI_T_cvar_get_info	shmem_double_prod_to_all	shmemQuiet
MPI_File_write_at_all	MPI_T_cvar_get_info	shmem_double_put	shmemQuiet
MPI_File_write_at_all_begin	MPI_T_cvar_handle_alloc	shmem_double_put_nbi	shmemQuiet
MPI_File_write_at_all_end	MPI_T_cvar_handle_free	shmem_double_set	shmemQuiet
MPI_File_write_ordered	MPI_T_cvar_write	shmem_double_swap	shmemQuiet
MPI_File_write_ordered_begin	MPI_T_enum_get_info	shmem_collect32	shmemQuiet
MPI_File_write_ordered_end	MPI_T_enum_get_item	shmem_collect64	shmemQuiet
MPI_Finalize	MPI_Finalize	shmem_finalize	shmemQuiet
MPI_Finalized	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Free_mem	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Gather	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Gatherv	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Get	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Get_accumulate	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Get_address	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Get_count	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Get_elements	MPI_Finalized	shmem_finalize	shmemQuiet
MPI_Get_elements_x	MPI_Finalized	shmem_finalize	shmemQuiet

Open MPI is an Associated Project of the Software in the Public Interest
non-profit organization.



Contact Open MPI webmaster

Open MPI server)
I'm:

Page last modified: 1
Mar-20
©2004-2018 The Open M
Proj

2018-

<https://www.open-mpi.org>

<https://computing.llnl.gov/tutorials/openMP/>

```

# Run with mpiexec -n 4 python ./mpi-montecarlo.py -s 1000000
from mpi4py import MPI
import argparse
import random
from math import pi

comm = MPI.COMM_WORLD
# This script will be executed once per process
size = comm.Get_size() # Number of processes available
rank = comm.Get_rank() # Rank of the process executing this process

def compute_pi(args):
    n = int(args.steps / size)

    if rank == 0:
        steps = [n] * size
    else:
        steps = None

    steps = comm.scatter(steps, root=0)
    print(f"Rank={rank}, steps={steps}")

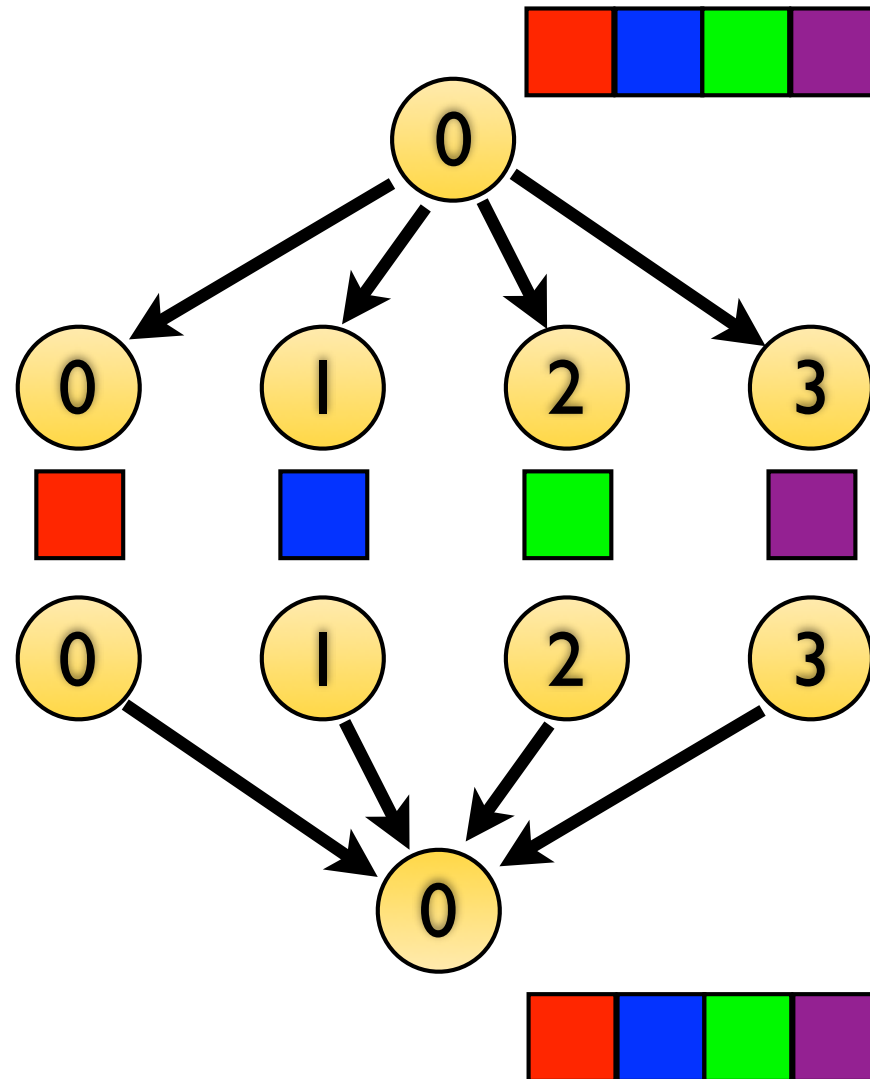
    s = sample_pi(steps)

    result = comm.gather(s, root=0)
    if rank == 0:
        n_total = n*size
        s_total = sum(result)
        pi_est = (4.0*s_total)/n_total
        print(" Steps\tSuccess\tPi est.\tError")
        print("%6d\t%7d\t%1.5f\t%1.5f" % (n_total, s_total, pi_est, pi-pi_est))

```

MPI: Scatter/Gather

Convention: rank 0
process is root

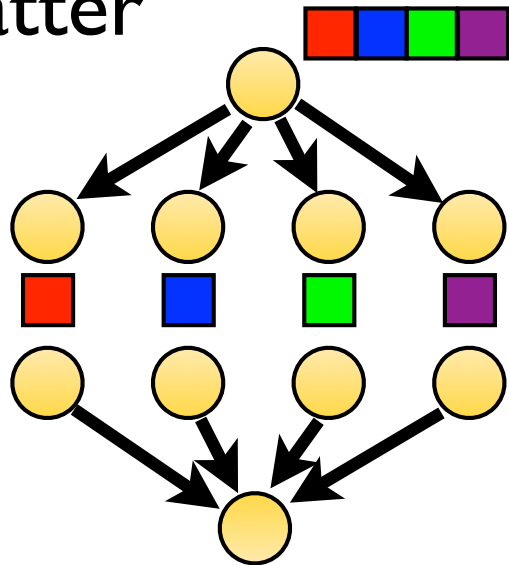


Scatter

Gather

MPI scatter/gather vs. mp.map

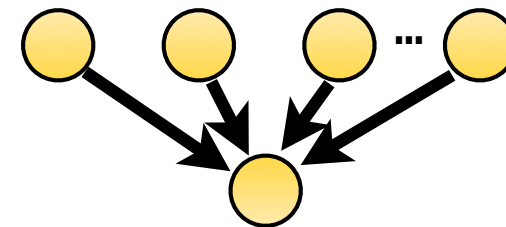
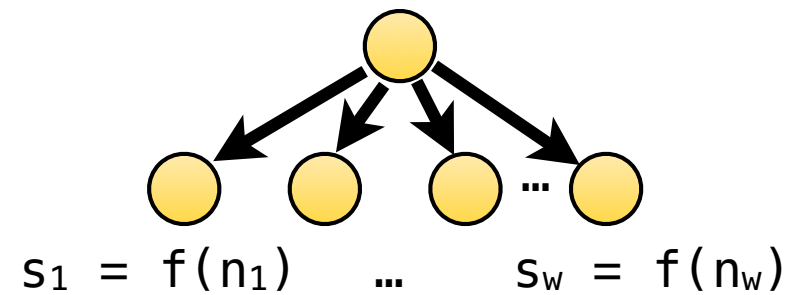
Scatter



Gather

$s = p.\text{map}(f, [n] * w)$

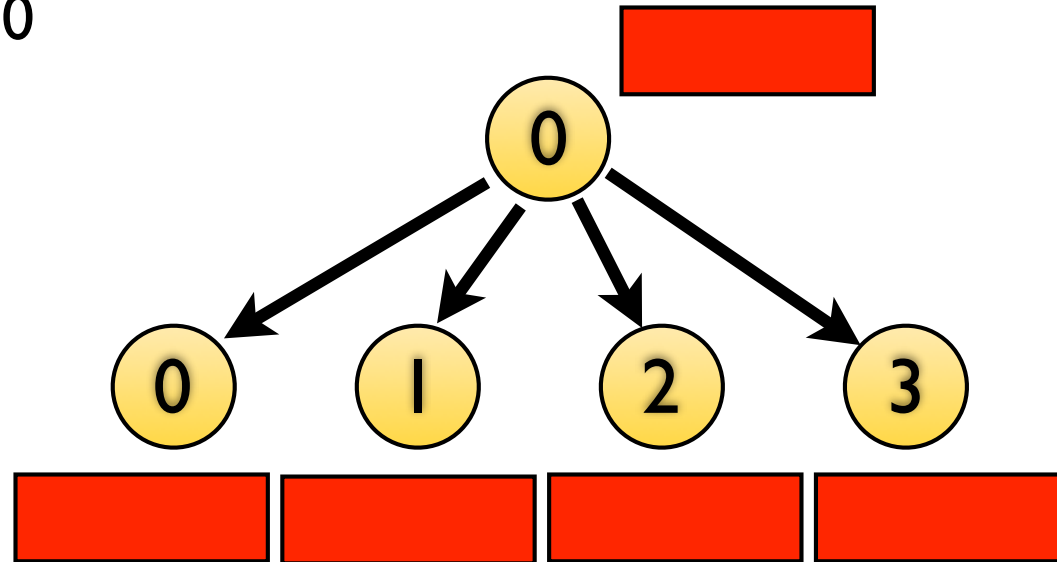
$[n_1, n_2, n_3, \dots, n_w]$



$[s_1, s_2, s_3, \dots, s_w]$

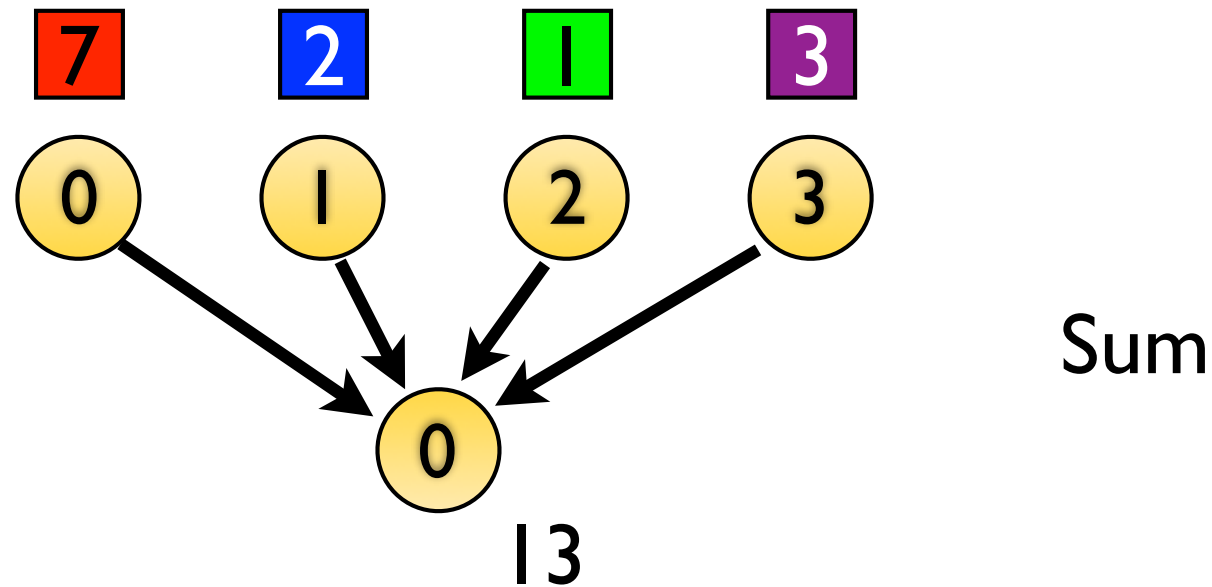
MPI Broadcast

Convention: rank 0
process is root



Broadcast

MPI Reduce



Reduce supports binary operations w/
associative law: Sum, Prod, Max, Min,
Argmax, Argmin, Logic operations,...

```

def compute_pi(args):
    n = int(args.steps / size)

    if rank == 0:
        steps = [n] * size
    else:
        steps = None

    steps = comm.scatter(steps, root=0)
    print(f"Rank={rank}, steps={steps}")

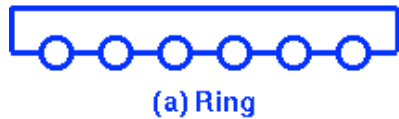
    s = sample_pi(steps)

    s_total = 0
    s = np.asarray(s)
    s_total = np.asarray(s_total)

    comm.Reduce(s, s_total, op=MPI.SUM, root=0)

    if rank == 0:
        n_total = n*size
        pi_est = (4.0*s_total)/n_total
        print(" Steps\tSuccess\tPi est.\tError")
        print("%6d\t%7d\t%1.5f\t%1.5f" % (n_total, s_total, pi_est, pi-pi_est))

```



During a **gather** operation, when one message is passed from every one of the n nodes to the root, the total number of message passing operations is ...

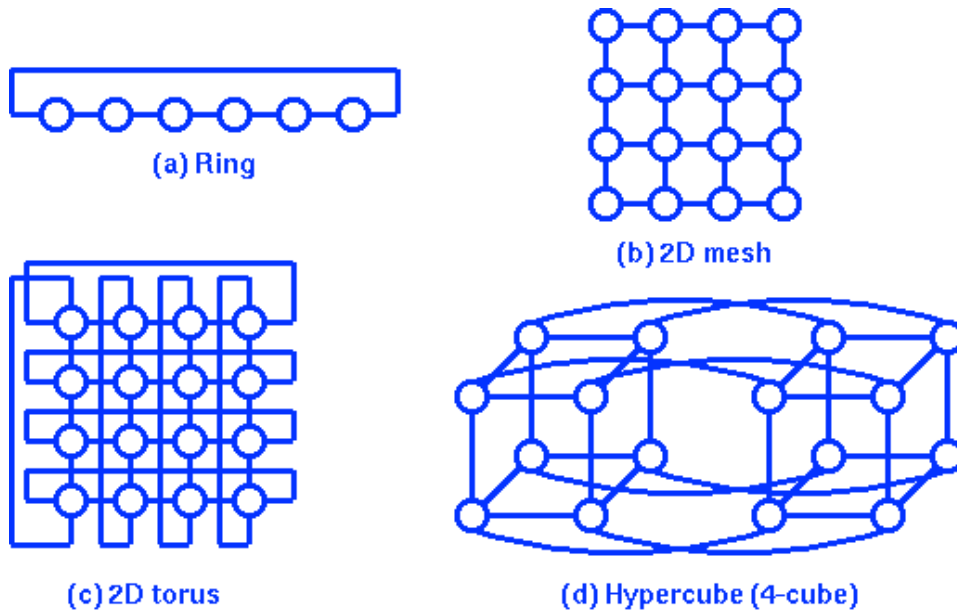
n nodes, n odd

For node farthest from root 1
 second-farthest 2, 3rd-farthest 3
 $+2+3+4+5+\dots+(n-1)/2 = 1/8 (n^2-1)$
 using Gauss

Here a message passing operation is passing the node's own data, or a message received from another node. Assume that no data is aggregated.

A	B	C	D	E
$O(\log n)$	$O(n)$	$O(n \log n)$	$O(n^2)$	$O(n^3)$

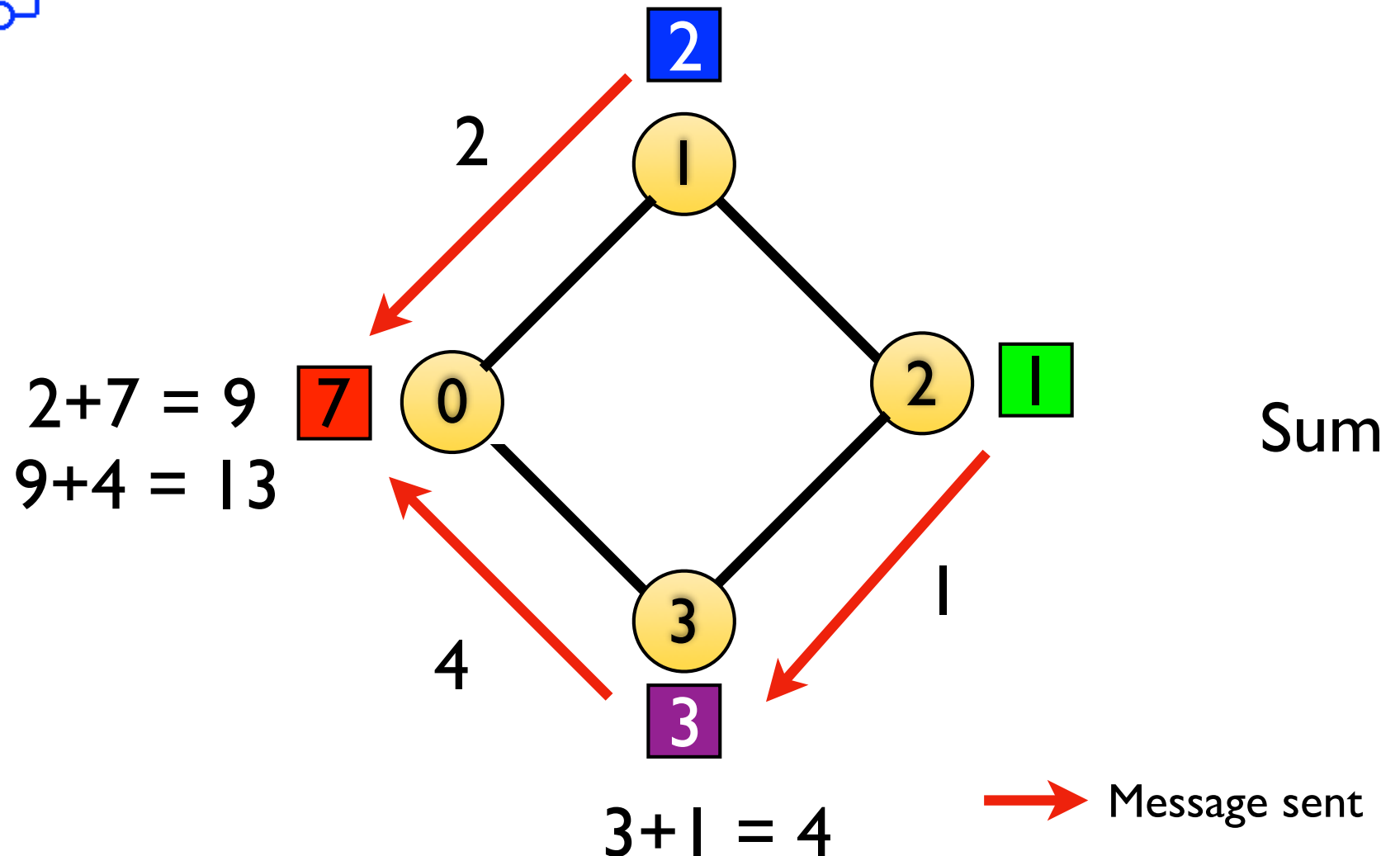
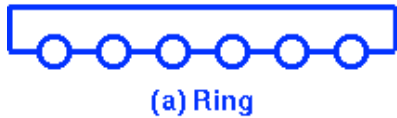
Reduce is *not just* Gather + Function call



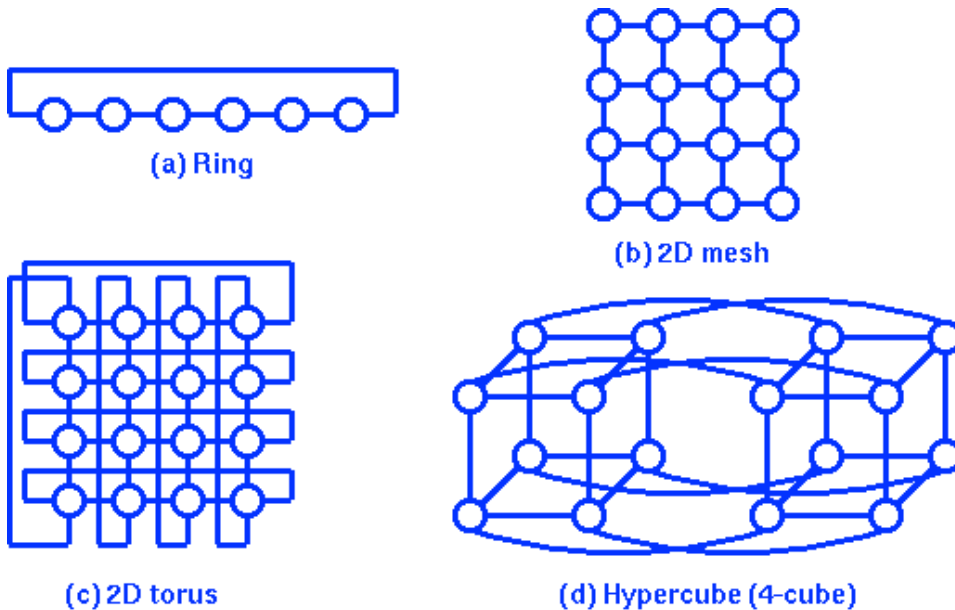
Example Interconnect topologies

- Communication Bottlenecks close to MPI root
- Substantial Latencies for thousands of cores
- Broadcast, Gather, Scatter similarly adapt to interconnect topology

Reduce is *not* Gather + Function call



Importance of Reduce



Example Interconnect topologies

- Communication Bottlenecks close to root
- Substantial Latencies for thousands of cores
- Reduce operations yield $O(\text{nodes})$ message passing operations instead of $O(\text{nodes}^2)$
- Broadcast, Gather, Scatter similarly adapt to interconnect topology

Reflection

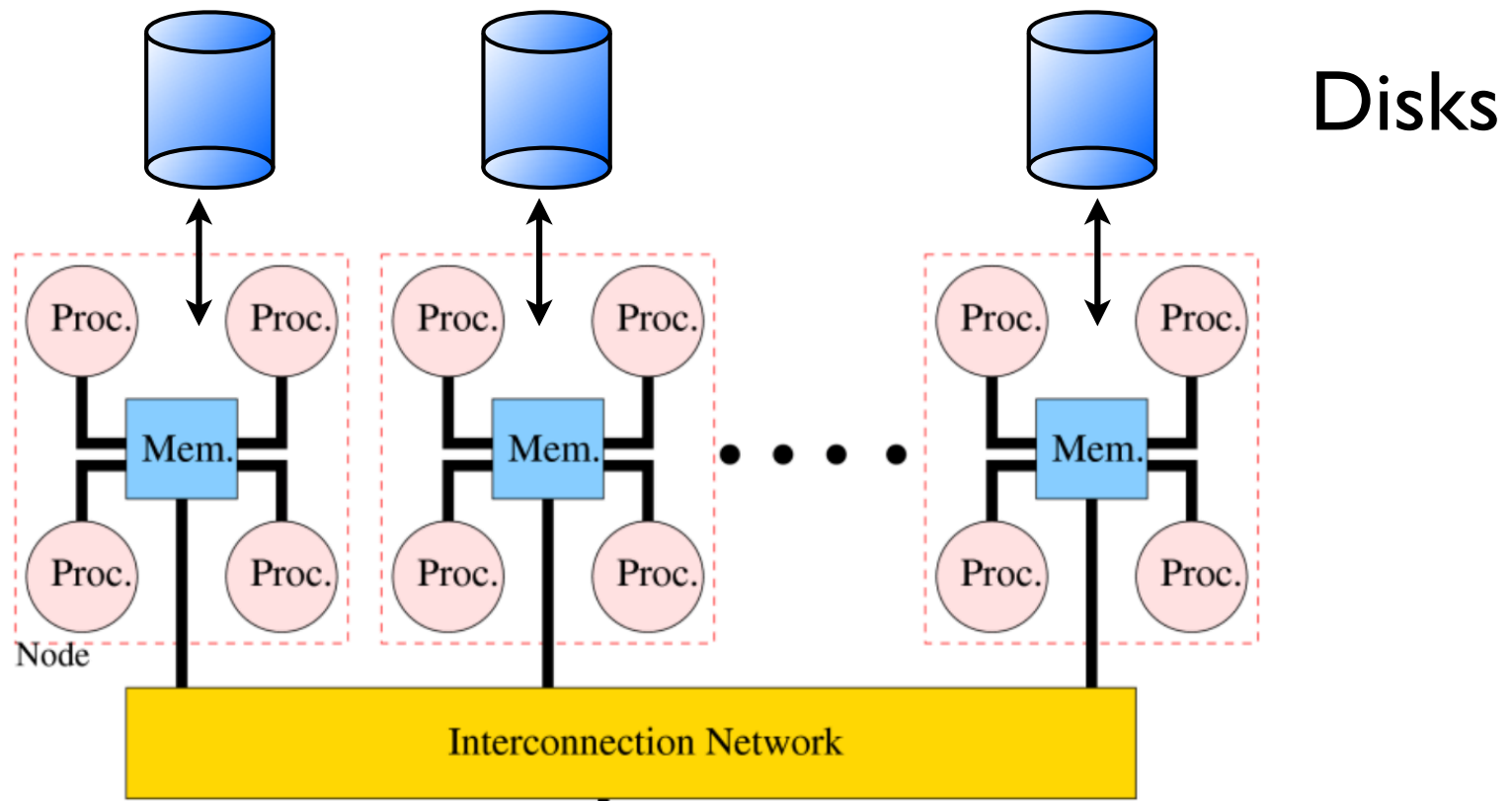
Architectures & Models

- SMP & Multithreading:
 - Shared memory, fast communication between all threads
- Distributed Memory and Message Passing
 - High latency, interconnection topology matters: limited #connections

Workload Comparison

	HPC	Data Science
CPU	Lots pf cycles Floating Point (FP) Little memory	Medium Mixed integer/FP Large memory
Data/File IO	Often very little IO	Massive IO Computations IO-bound
Tasks	Standard Tools (e.g. matrix solvers)	Custom tools and frequently new, changing problems

System architecture for data-centric computation



Aspects of data-centric computation

- Scale of Data: Parallel IO, Data-compute co-location
- Cost of computation: commodity hardware
- Reliability of systems: Design for failures
- Software is run once: support rapid development:
 - automated parallelization
 - restricted model of computation

Parallel programming models

Map Reduce

See Skiena, Data Science Design Manual, Chap 12.6

<http://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>

Map/reduce assumptions

- Input files are distributed over nodes (implicit)
- All data are (key, value) tuples
- Main parallel operations are:
 - Map:
 $(\text{key}, \text{value}) \rightarrow (\text{key}_1, \text{value}_1)[, (\text{key}_2, \text{value}_2), \dots]$
 - Reduce:
 $(\text{key}_1, [\text{value}_1, \text{value}_2, \text{value}_3, \dots]) \rightarrow$
 $(\text{key}_{\text{new}}, \text{value}_{\text{new}})[, (\text{key}_{\text{new}2}, \text{value}_{\text{new}2}), \dots]$

Example: Count number of links to webpages

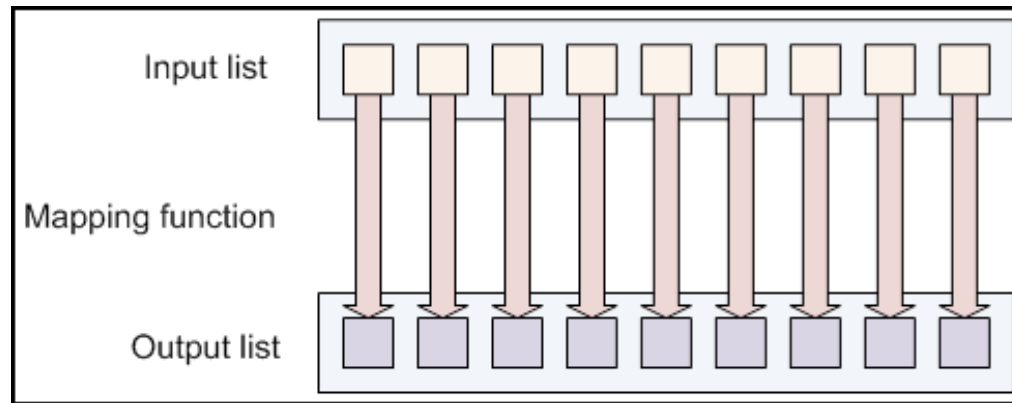
Map(DocumentURL, DocumentHTML):

For each url in DocumentHTML return (url, 1)

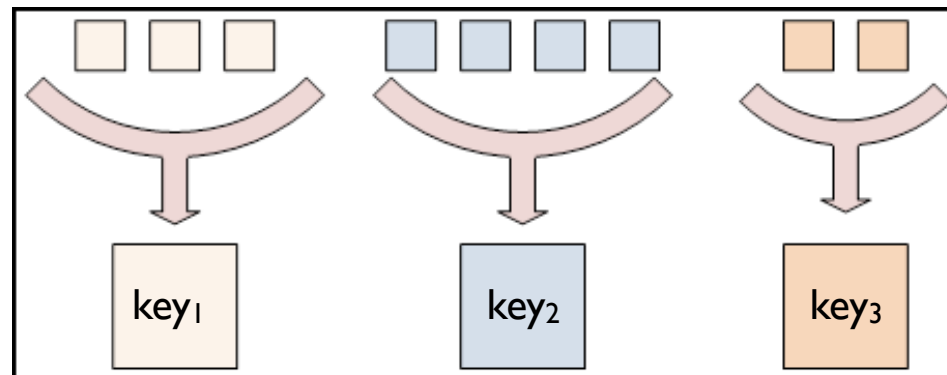
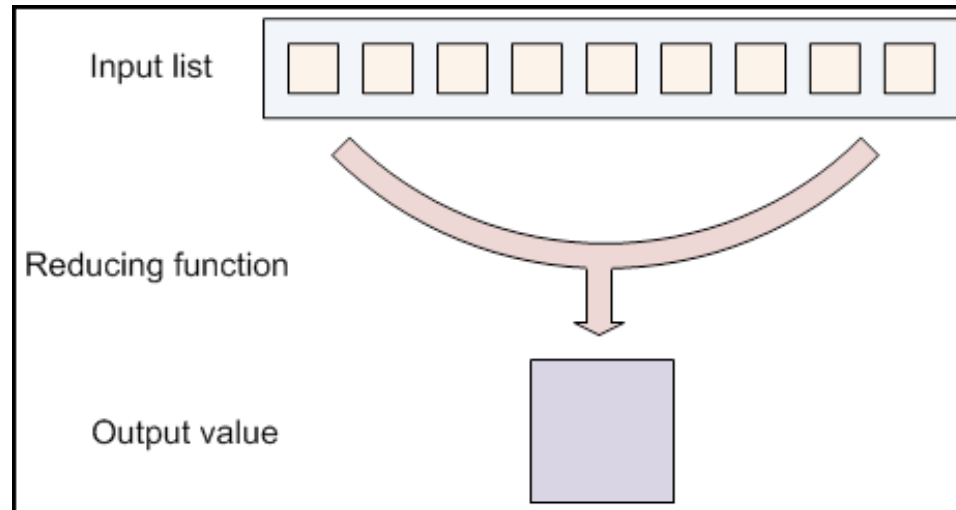
Reduce(url, counts) # counts = [count1, count2,...]

return (url, sum(counts))

Map

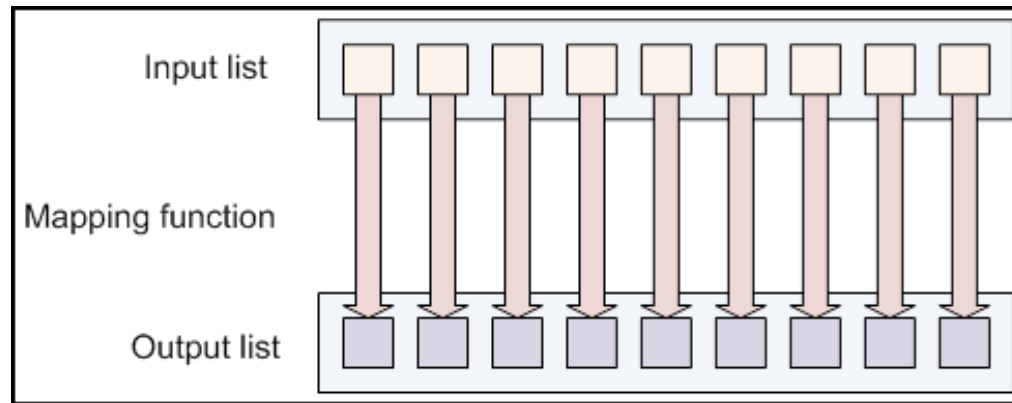


Reduce



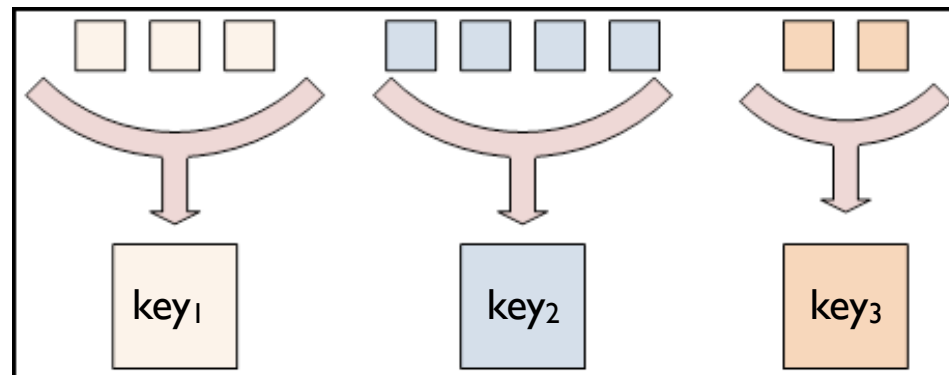
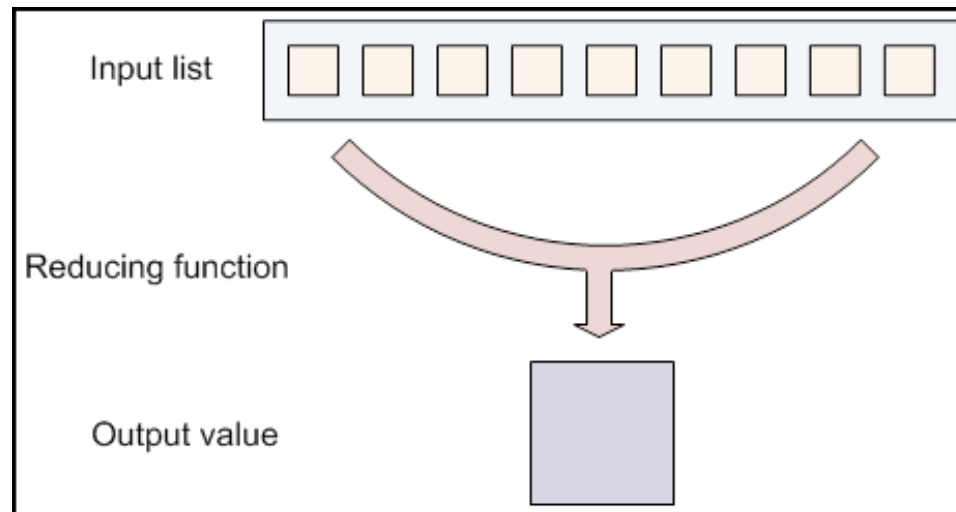
Reduce receives
input aggregated by
key

Map



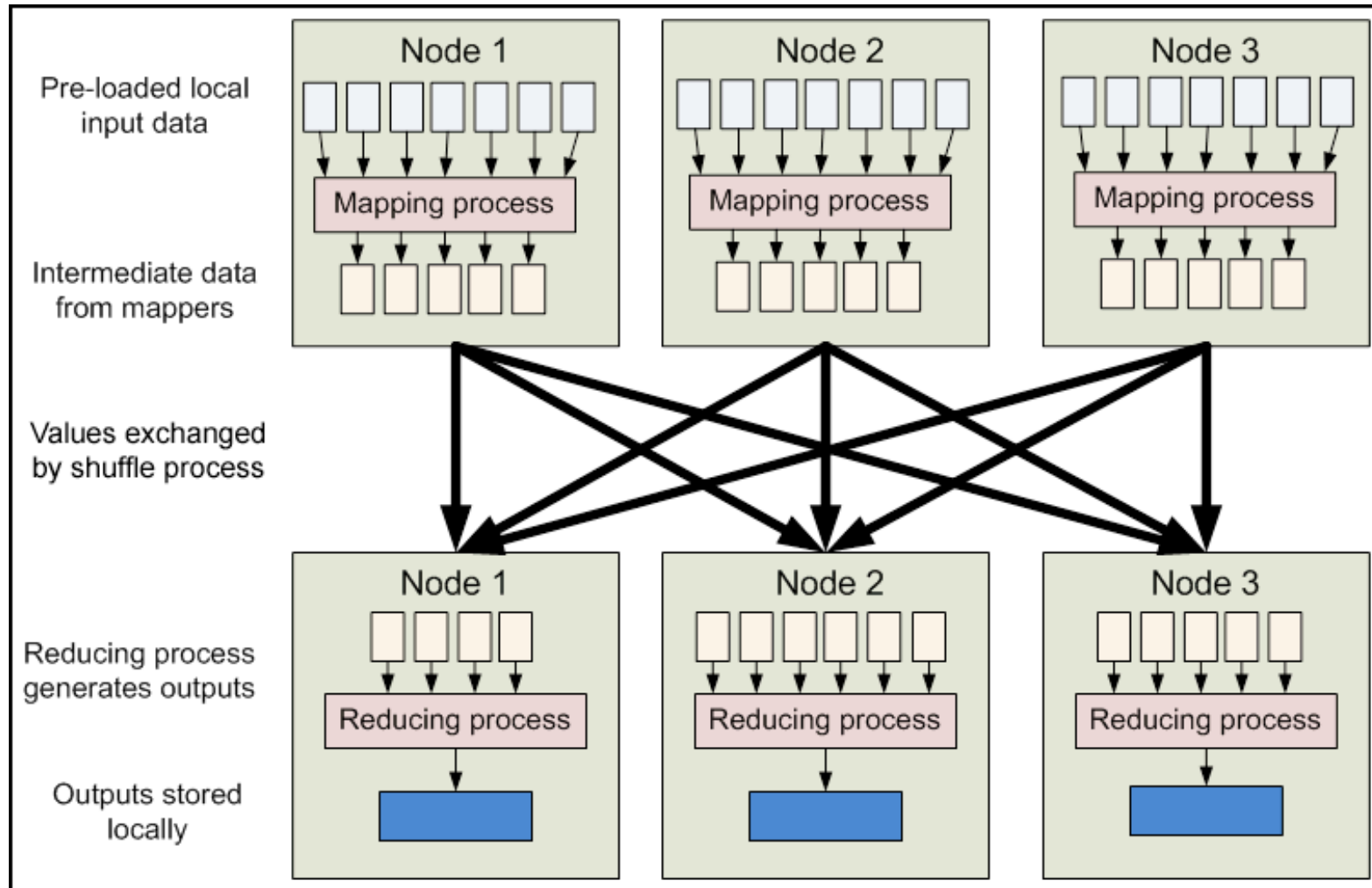
Each map call only needs its input; no assertions about other inputs, execution order, assignment to nodes

Reduce



Each reduce call receives all values for its key; no assertions about other keys, execution order, assignment to nodes

Map/Reduce Parallelism



MapReduce vs multi-threading and MPI

- No concept of process, thread nor hardware
- No control over data flow
- Reduce obtains all values with the same key
- *Substantial effort behind the scenes*

MrJob: Map/Reduce in Python

- <https://github.com/Yelp/mrjob>
- <https://mrjob.readthedocs.io/en/stable/>
- Pure Python
- Runs locally and in the cloud


```
""" Find duplicate keys in a file containing lines consisting of
    key,value
"""
from mrjob.job import MRJob

class FindDuplicates(MRJob):

    def mapper(self, _, line):
        key, value = line.split(',')
        yield (key, 1)

    def reducer(self, key, counts):
        s = sum(counts)
        if s > 1:
            yield (key, s)

if __name__ == '__main__':
    FindDuplicates.run()
```

Implementing wordcount with Map-reduce

```
schliep@MacBookPro13:> wc mpi-summary-statistics.py  
63      197    1918 mpi-summary-statistics.py
```

```
from mrjob.job import MRJob

class WC(MRJob):

    def mapper(self, _, line):
        nr_words = len(line.split())
        nr_chars = len(line)
        yield ("words", nr_words)
        yield ("chars", nr_chars)

    def reducer(self, key, counts):
        if key == "words":
            yield ("words", sum(counts))
        elif key == "chars":
            yield ("chars", sum(counts))

if __name__ == '__main__':
    WC.run()
```

```
def mapper(self, _, line):
    nr_words = len(line.split())
    nr_chars = len(line)
    yield ("words", nr_words)
    yield ("chars", nr_chars)
```

n = # lines

```
def reducer(self, key, counts):
    if key == "words":
        yield ("words", sum(counts))
    elif key == "chars":
        yield ("chars", sum(counts))
```

.

The maximal speedup achievable by the mapper is ...

Theoretical n
real-world k = #nodes

A	B	C	D	E
1	2	4	log(n)	n

n = # lines

```
from mrjob.job import MRJob
```

```
class WC(MRJob):
```

```
    def mapper(self, _, line):  
        nr_words = len(line.split())  
        nr_chars = len(line)  
        yield ("words", nr_words)  
        yield ("chars", nr_chars)
```

Creates 2n tuples

Parallel: up to one
node per line

```
    def reducer(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

```
if __name__ == '__main__':  
    WC.run()
```

```
def mapper(self, _, line):  
    nr_words = len(line.split())  
    nr_chars = len(line)  
    yield ("words", nr_words)  
    yield ("chars", nr_chars)
```

n = # lines

```
def reducer(self, key, counts):  
    if key == "words":  
        yield ("words", sum(counts))  
    elif key == "chars":  
        yield ("chars", sum(counts))
```

The maximal speedup achievable by the reducer is ...

A	B	C	D	E
1	2	4	$\log(n)$	n

```
from mrjob.job import MRJob

class WC(MRJob):

    def mapper(self, _, line):
        nr_words = len(line.split())
        nr_chars = len(line)
        yield ("words", nr_words)
        yield ("chars", nr_chars)

    def reducer(self, key, counts):
        if key == "words":
            yield ("words", sum(counts))
        elif key == "chars":
            yield ("chars", sum(counts))

if __name__ == '__main__':
    WC.run()
```

Is called twice.
 $O(n)$ input

Parallel: up to
two nodes

$O(n)$ communication

```
class WC(MRJob):
```

```
    def mapper(self, _, line):  
        nr_words = len(line.split())  
        nr_chars = len(line)  
        yield ("words", nr_words)  
        yield ("chars", nr_chars)
```

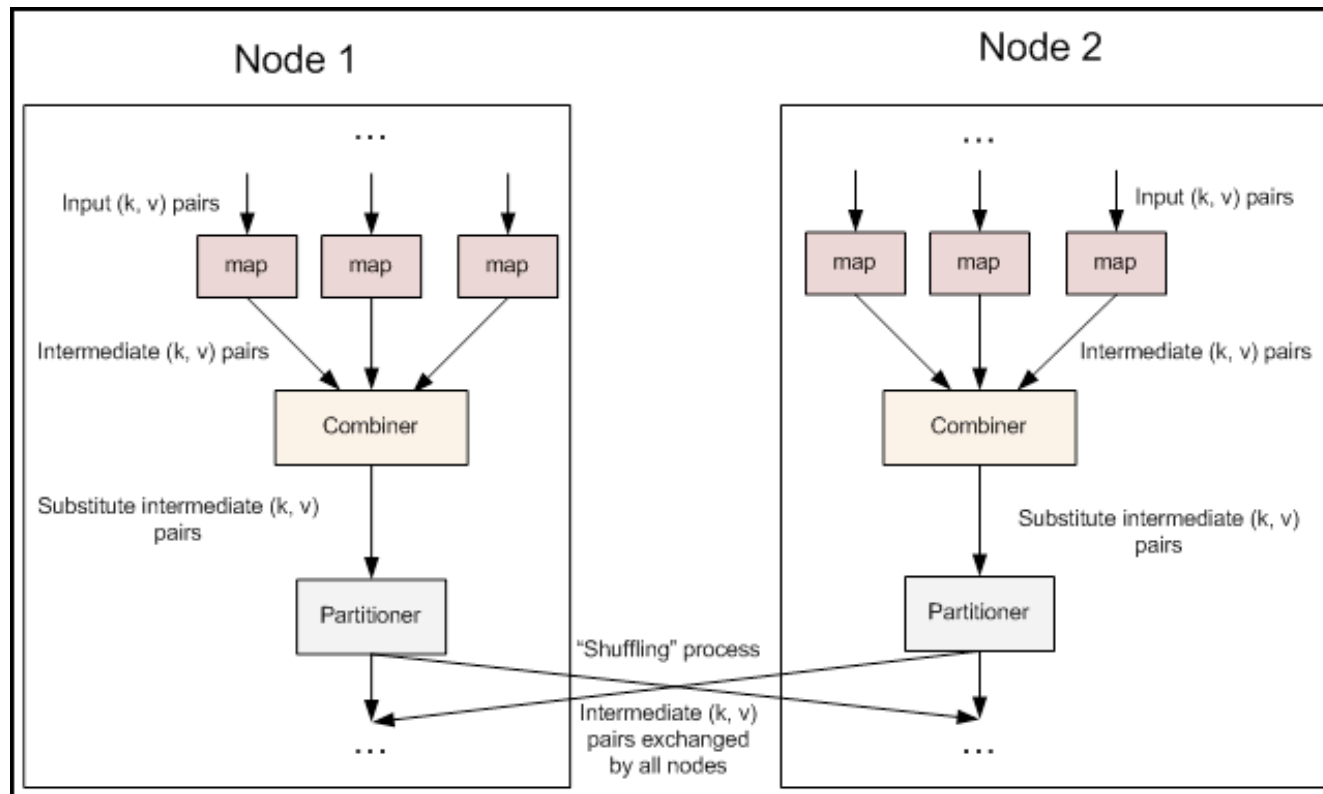
```
    def combiner(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

```
    def reducer(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

```
if __name__ == '__main__':  
    WC.run()
```

n = # lines
k = #nodes for
mapper

Combiners



```
class WC(MRJob):
```

```
    def mapper(self, _, line):  
        nr_words = len(line.split())  
        nr_chars = len(line)  
        yield ("words", nr_words)  
        yield ("chars", nr_chars)
```

```
    def combiner(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

```
    def reducer(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

```
if __name__ == '__main__':  
    WC.run()
```

n = # lines
k = #nodes for
mapper

Creates 2k tuples

As parallel as
mapper

*Same function as
reducer on a single
node*

```
class WC(MRJob):
```

```
    def mapper(self, _, line):  
        nr_words = len(line.split())  
        nr_chars = len(line)  
        yield ("words", nr_words)  
        yield ("chars", nr_chars)
```

```
    def combiner(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

```
    def reducer(self, key, counts):  
        if key == "words":  
            yield ("words", sum(counts))  
        elif key == "chars":  
            yield ("chars", sum(counts))
```

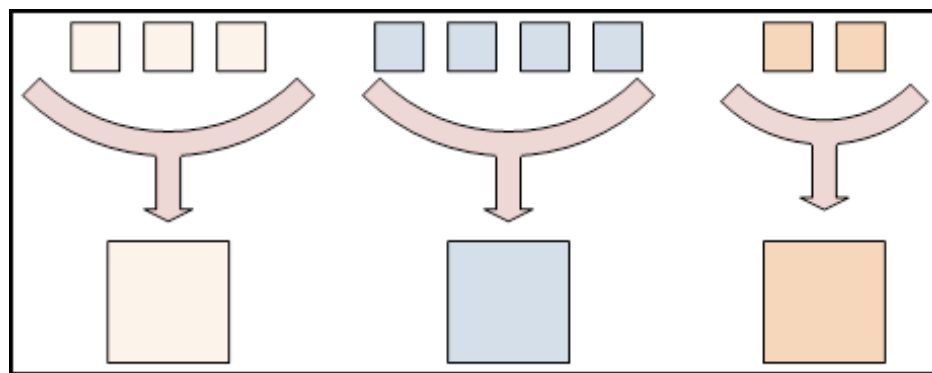
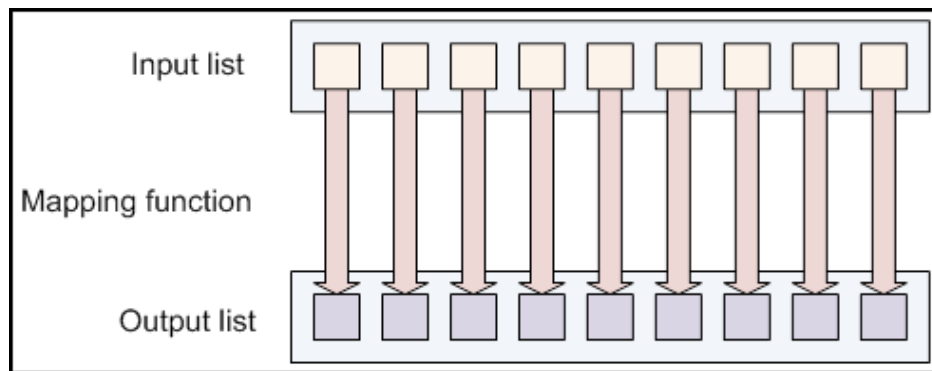
```
if __name__ == '__main__':  
    WC.run()
```

n = # lines
k = #nodes for
mapper

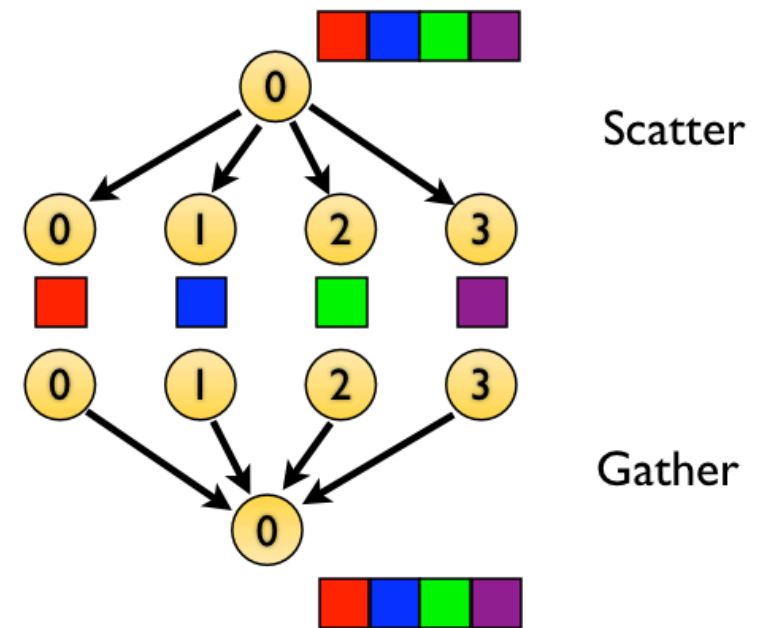
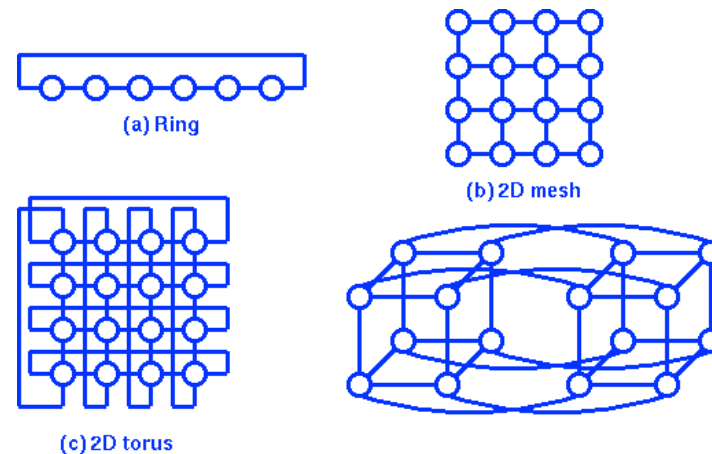
Called twice with
 $O(k)$ input

Parallel: two nodes

$O(k)$ communication



keys define
communication structure
and effort implicitly



Explicit parallelization in
MPI and threads

Questions & Answers from Chat

Q: Do mapper and combiner always use the same number of nodes? A: A combiner is run on all (key, value) pairs produced by a mapper on one node. So conceptually, yes.

Q: What is the speedup of mapper? A: The theoretical maximal speedup is $n = \text{\#lines}$ if executed on n nodes. Real-world: the speedup is roughly equal to the number of nodes the mapper is running on

Q: What does reduce do? A: Mapper produces ("words", 2), ("words", 4), ("words", 5) [on some machines, maybe in parallel] then reduce gets called with "words", [2,4,5] and thus can count the total number of words.

Q: Does that mean that the reducer has to sum n values together? A: If there is no combiner, yes. This is clearly bad, so that is why the combiner has been introduced. So on each of the nodes the combiner has to add up n/k values (this is parallel on k nodes; on each node multiple cores could be used). The reducer only get k values to add up.

Agenda for today:

- Using Multi-threaded programming
- Message passing:
 - Parallel primitives
 - Reduce
- Map/Reduce

Upcoming Deadlines



Python Programming 2

Not available until Apr 1 | Due Apr 15 at 10am | -/5 pts



Preparation for Lecture 4/15

Not available until Apr 1 | Due Apr 15 at 10am



Assignment 2

Available until Apr 20 | Due Apr 17 at 10am | -/14 pts