# Parallel & Concurrent Programming: Collective Communication

Emery Berger CMPSCI 691W Spring 2006



# Outline

Last time:

- Distributed parallel programming via message-passing
- MPI library approach
- Today:
  - Moving beyond point-to-point:
    - collective communication



some slides from van de Geijn et al.

# **Collective Communication**

- Instead of single target, collective communication operates on entire communicator
  - e.g., MPI\_COMM\_WORLD
  - Advantages over point-to-point:
    - Higher-level
    - Typically far more efficient
      - e.g., O(n)  $\Rightarrow$  O(log n)
      - Pipelining



# Simple Analytical Model

- To send n bytes:
  - α = startup latency
  - β = per-byte cost
    - **α + n**β
- Naive broadcast (p processors):
  - α(p-1) + (p-1)nβ
  - Minimum spanning tree:
    - α(log p + p 1) + 2(p-1)/p \* nβ
    - Note: does not model contention
      - In practice, can avoid it



# MPI Collective Ops

 MPI has rich set of collective communication operations (& duals):

- Synchronization
  - Barrier
- Communication
  - Broadcast
  - Scatter, Gather
  - Reduce, Scan

#### Communicator management



### **Barrier Svnchronization**

#### MPI\_Barrier:

progress continues past barrier only after all processes have arrived

// perform computation

--

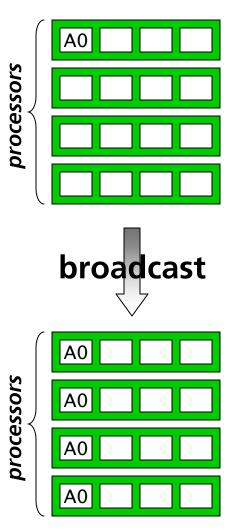
MPI\_Barrier (MPI\_COMM\_WORLD); // all done here // perform computation (faster)

MPI\_Barrier (MPI\_COMM\_WORLD); // all done here

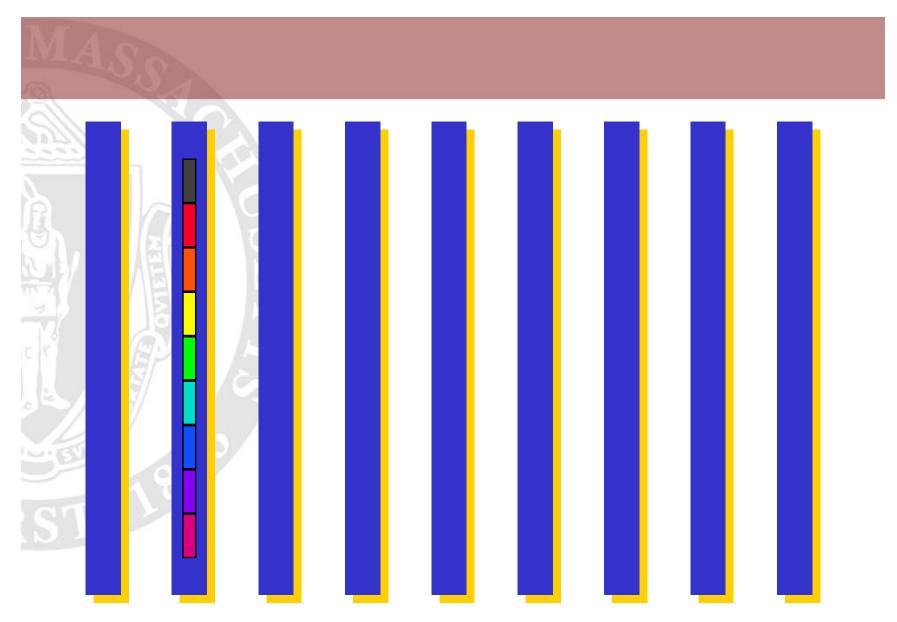


### From root: Broadcast

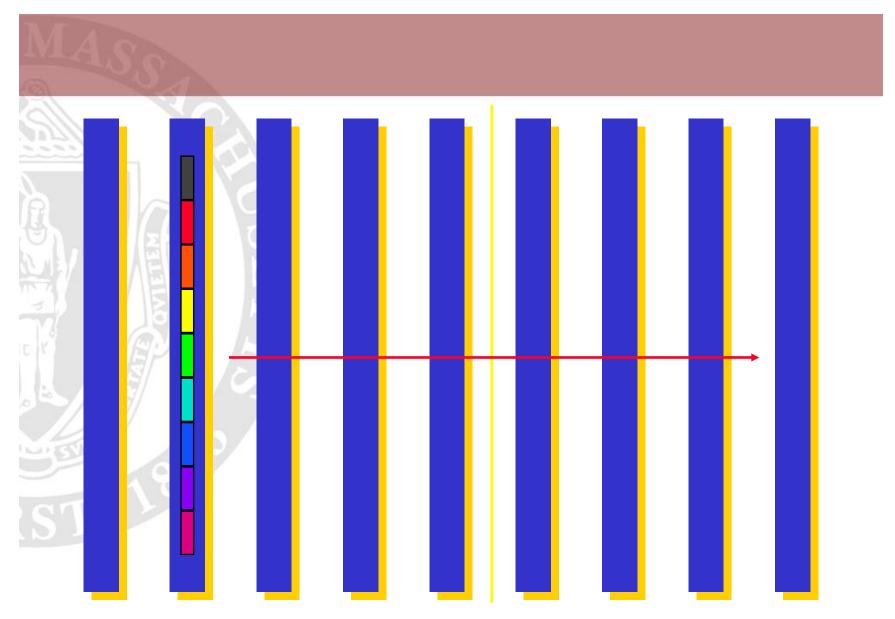
 MPI\_Bcast: broadcast single datum across all processors



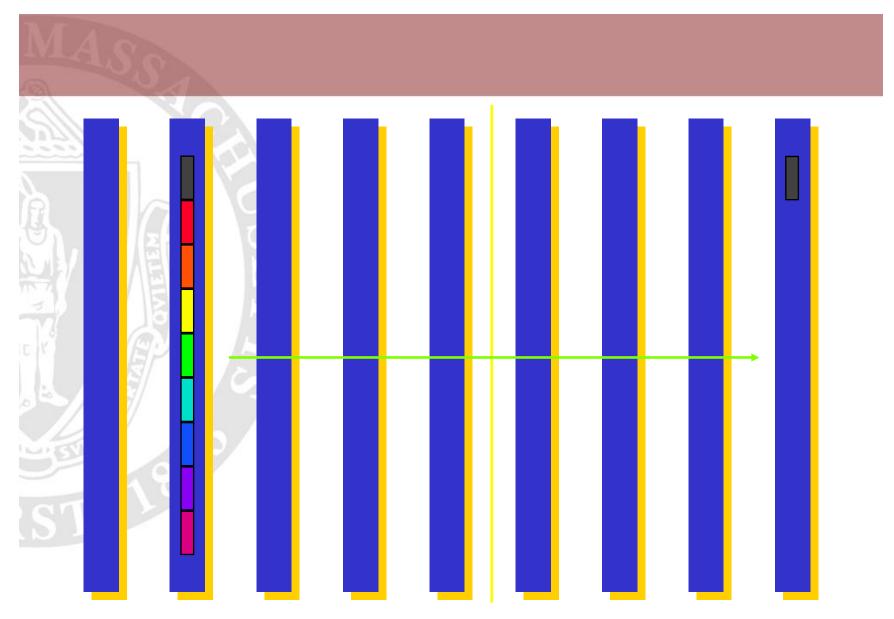




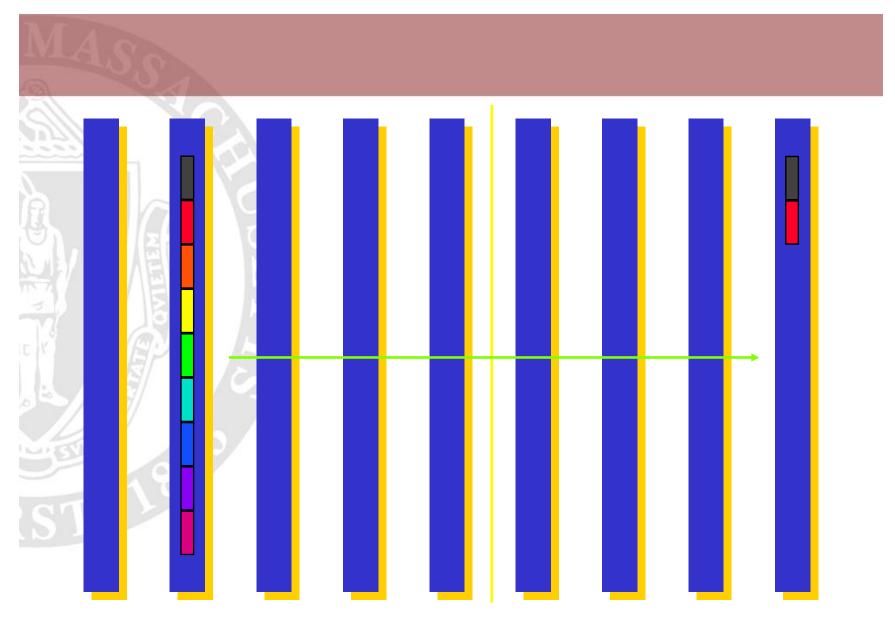




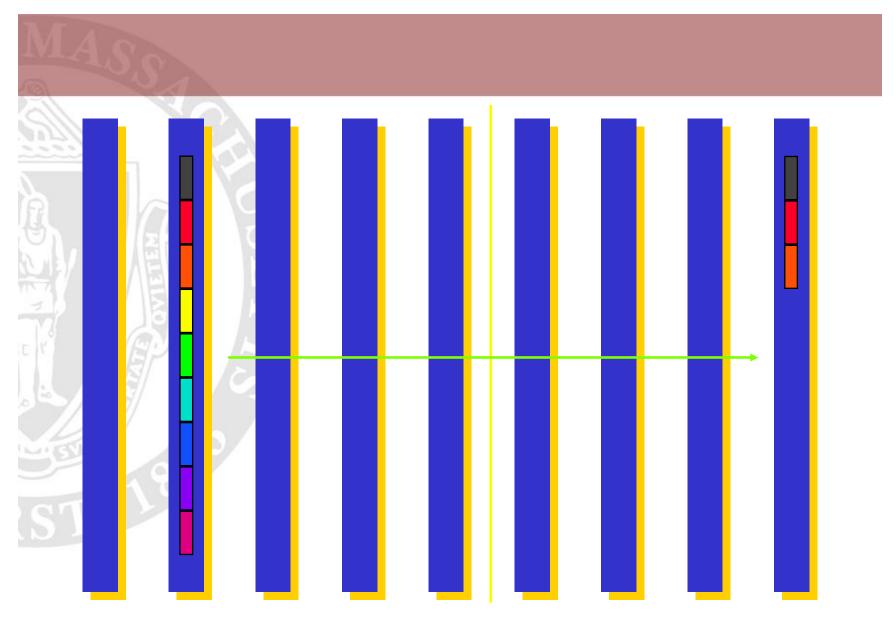




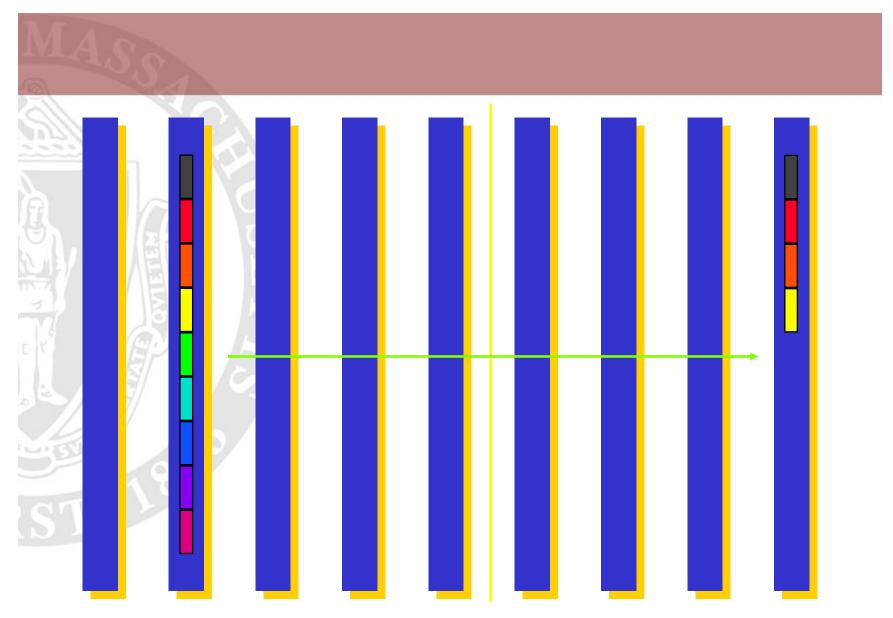




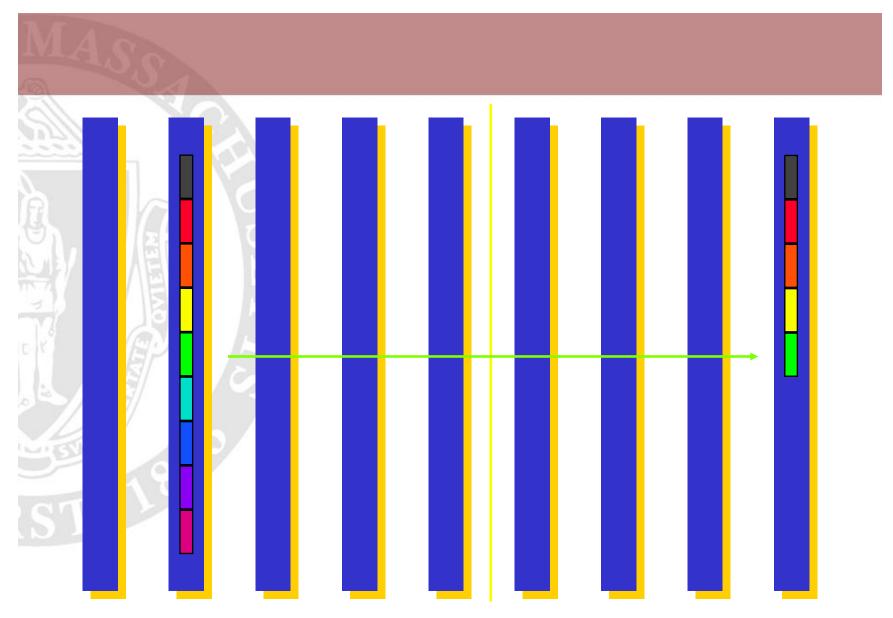




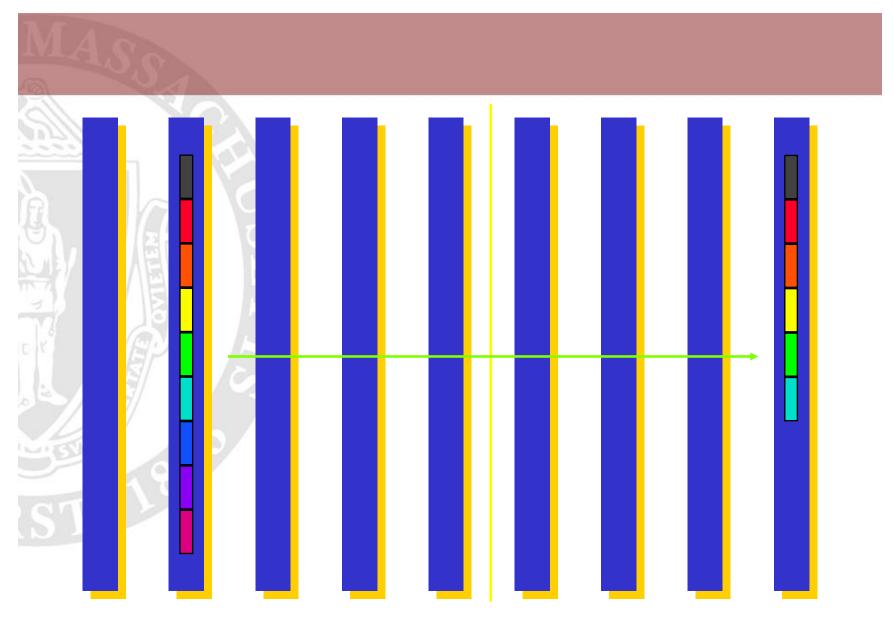




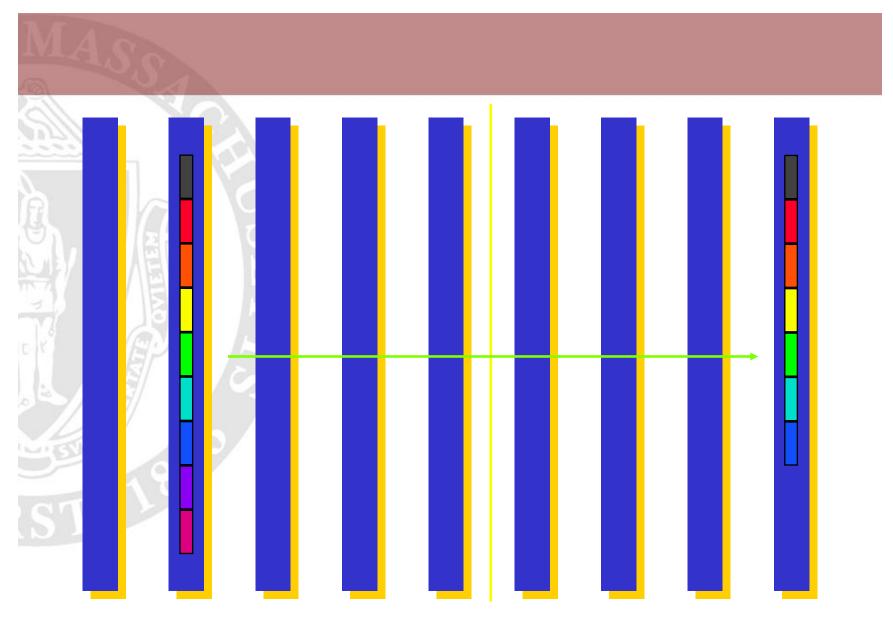




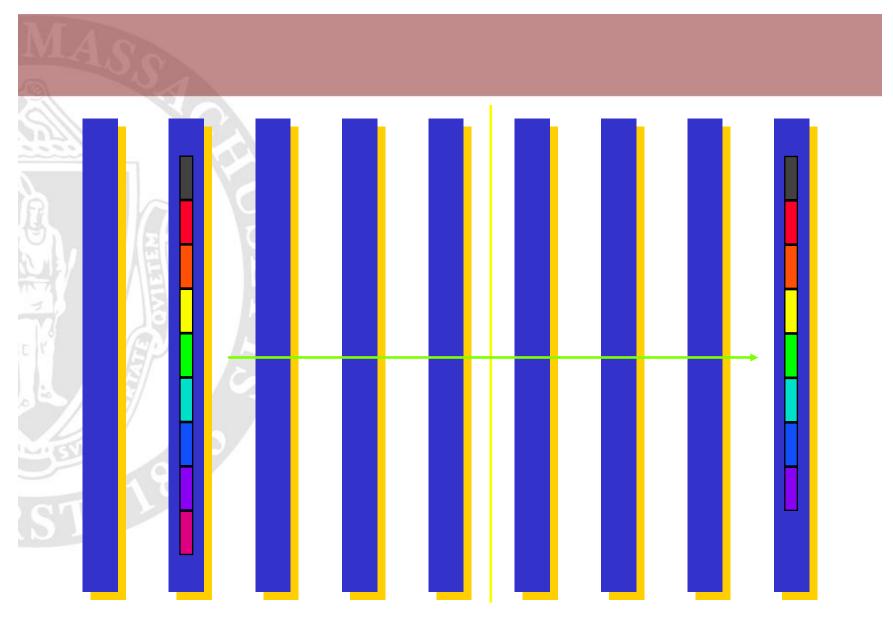




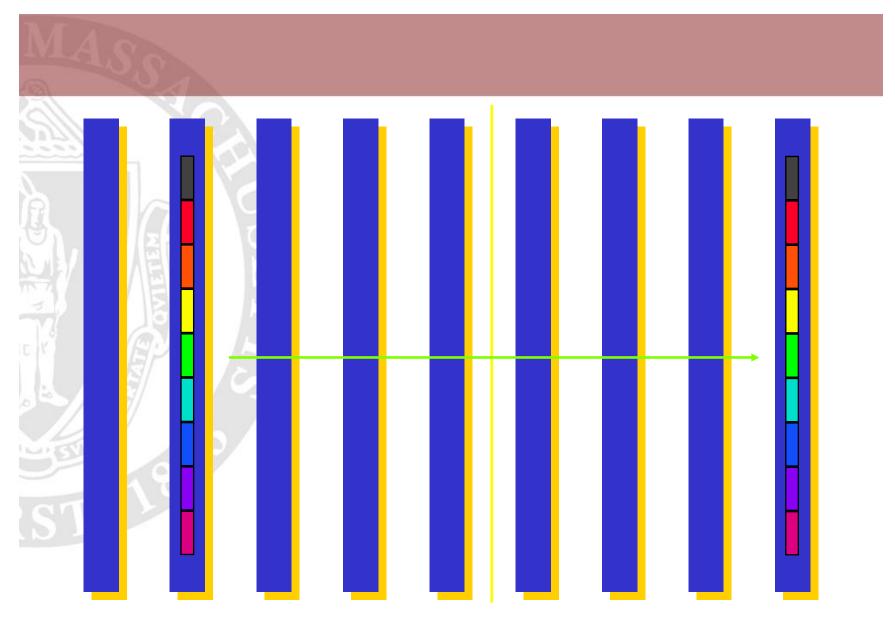




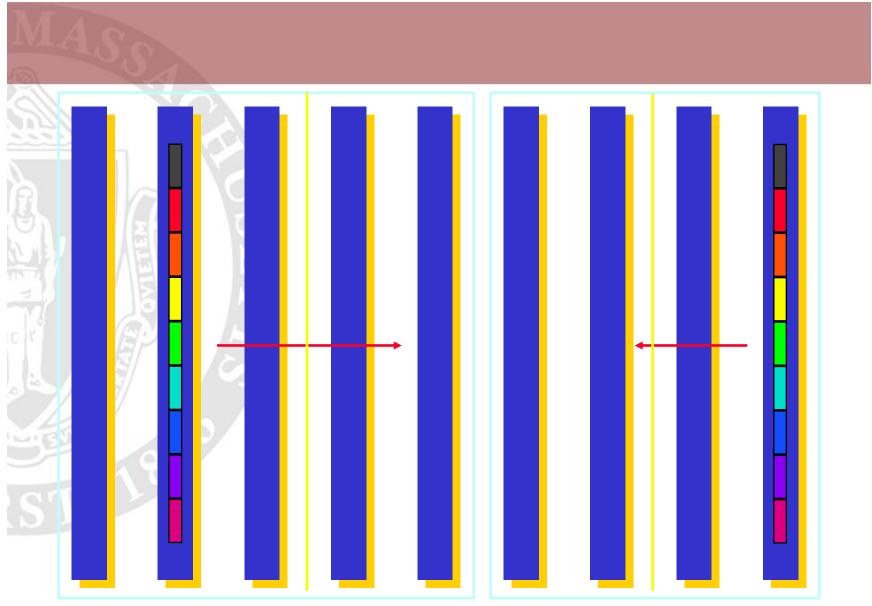




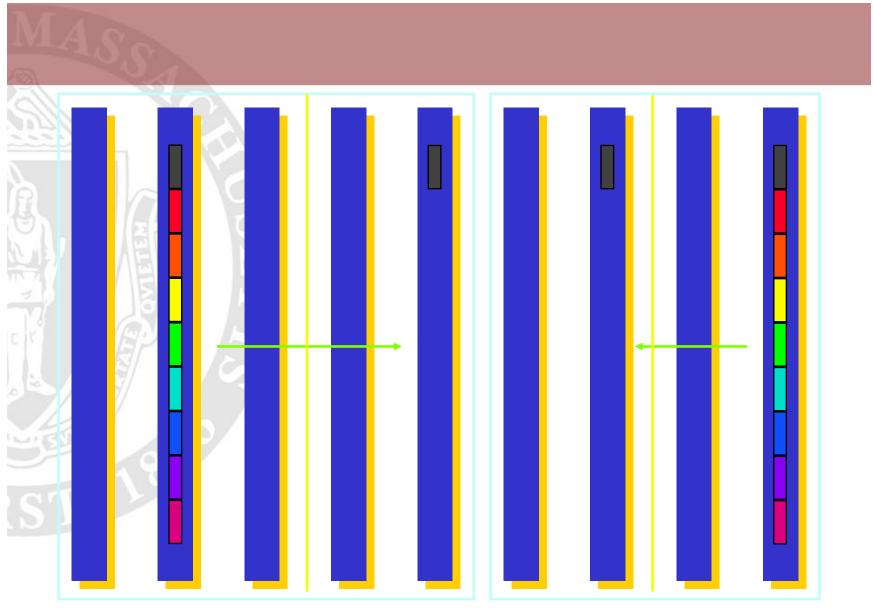








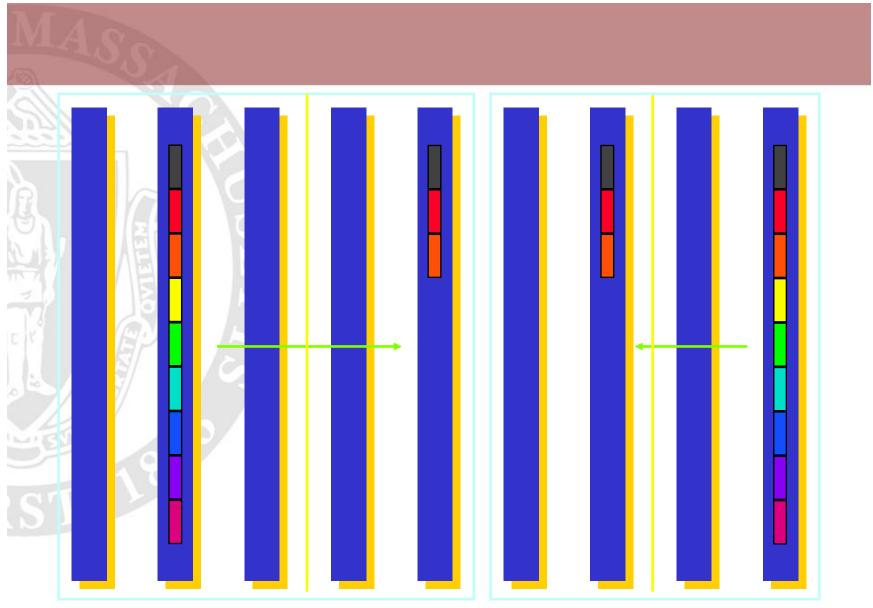




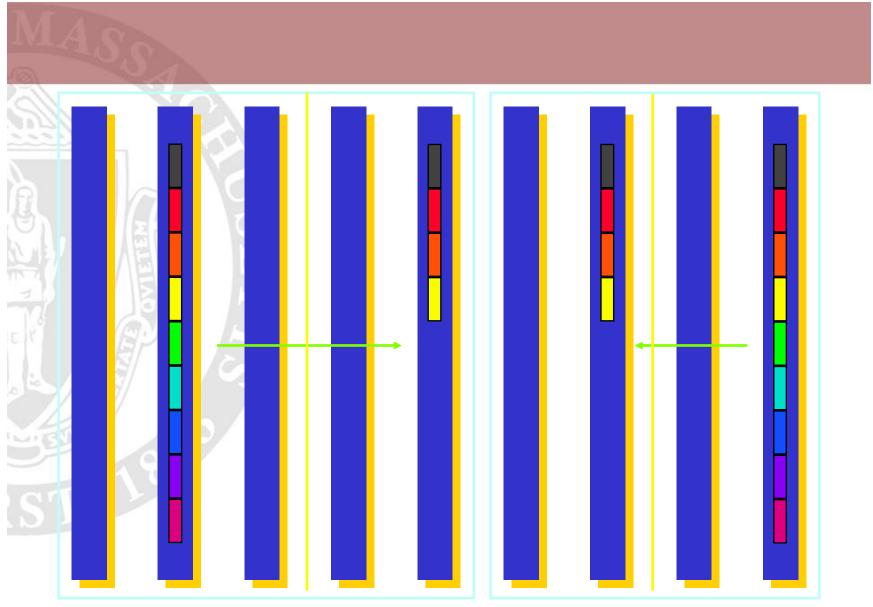




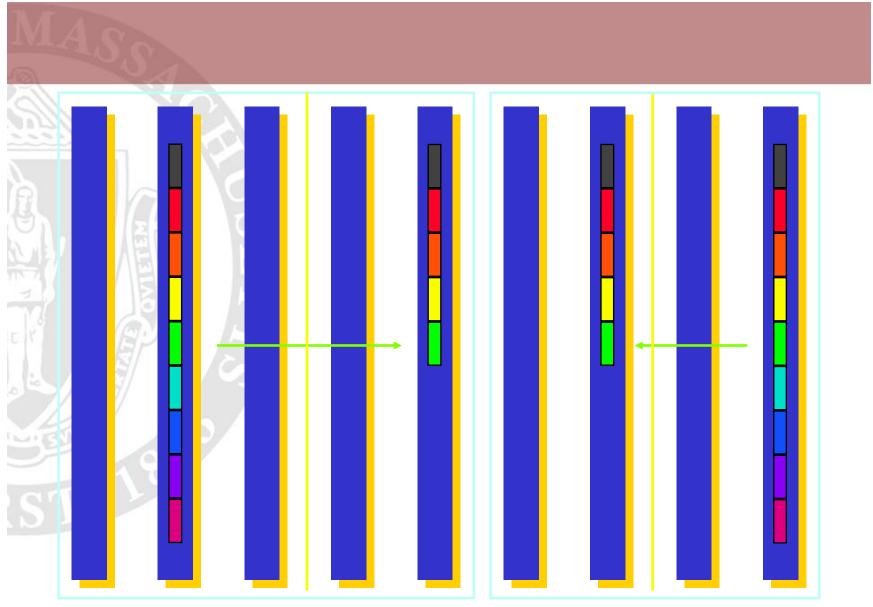




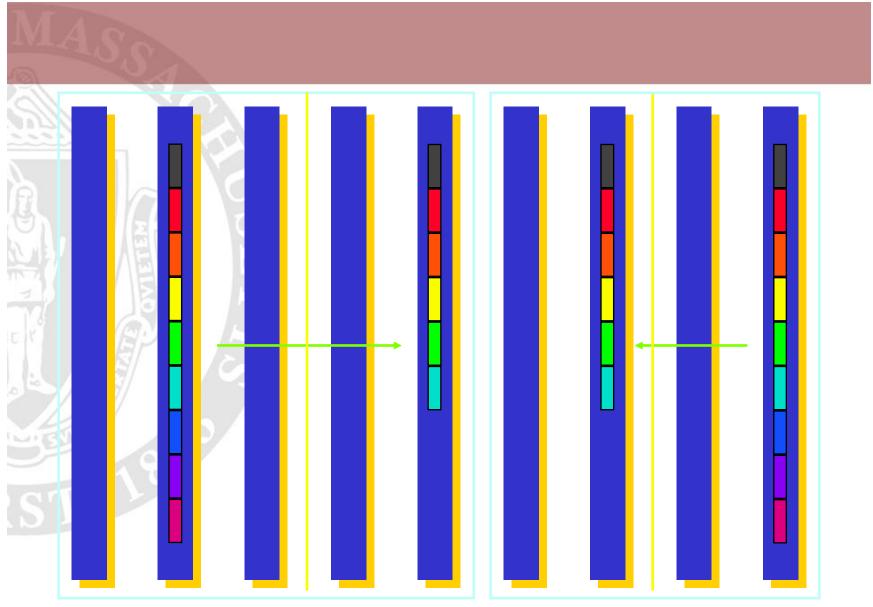




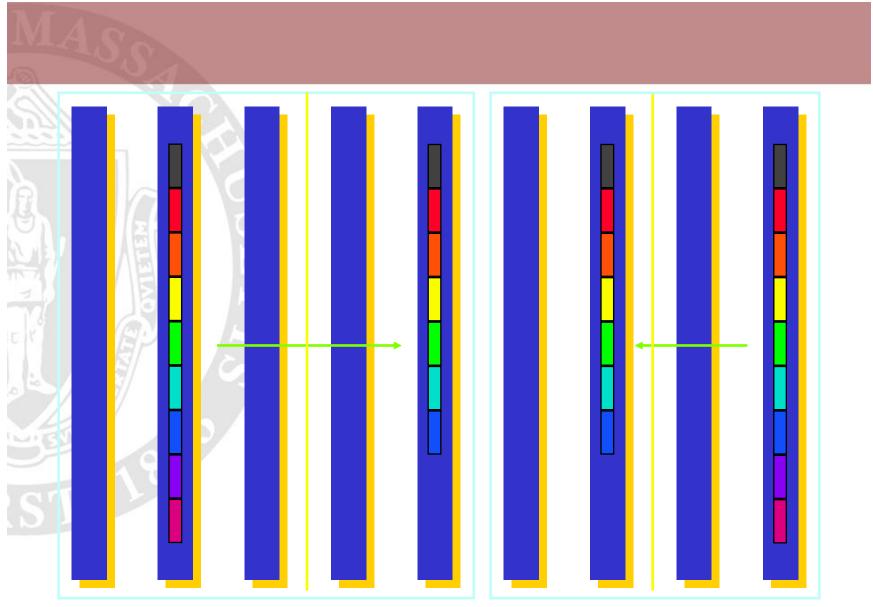




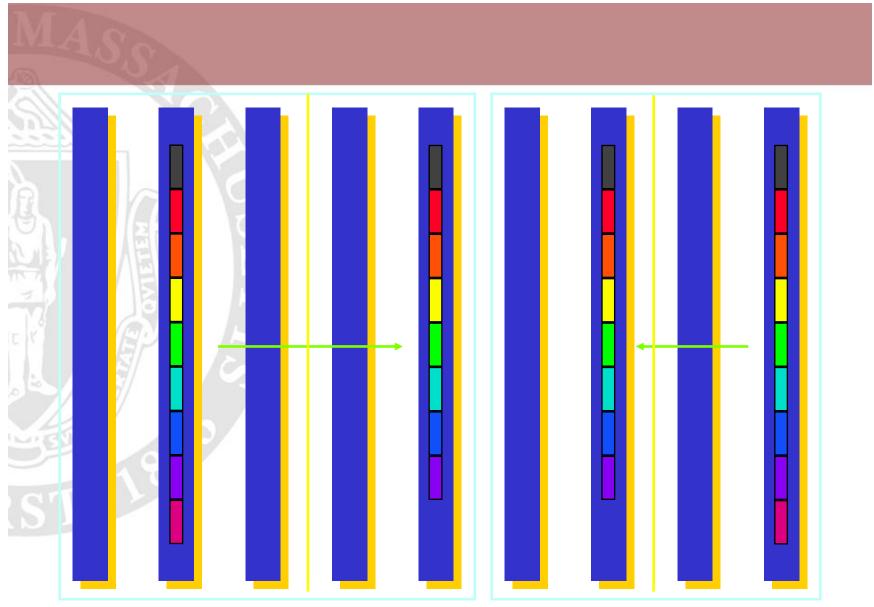




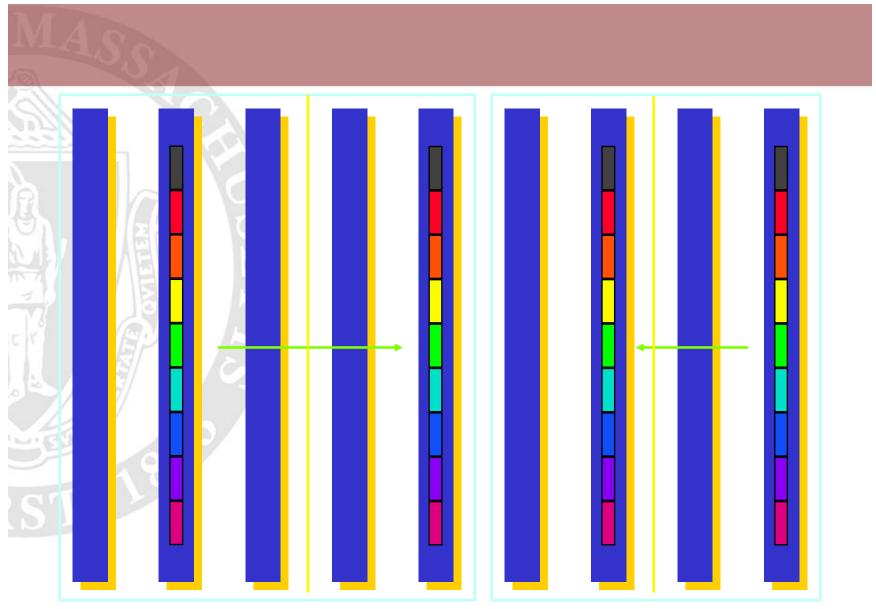




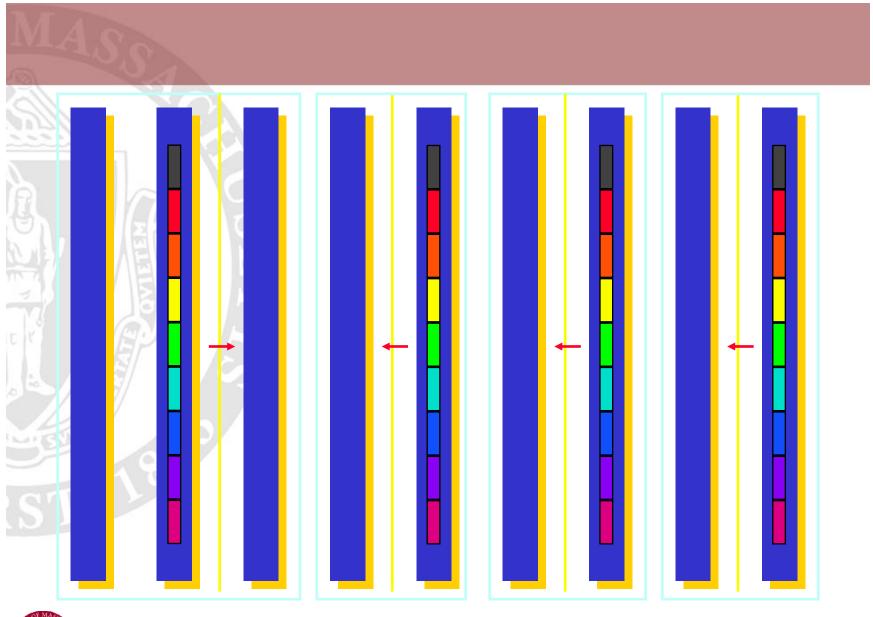




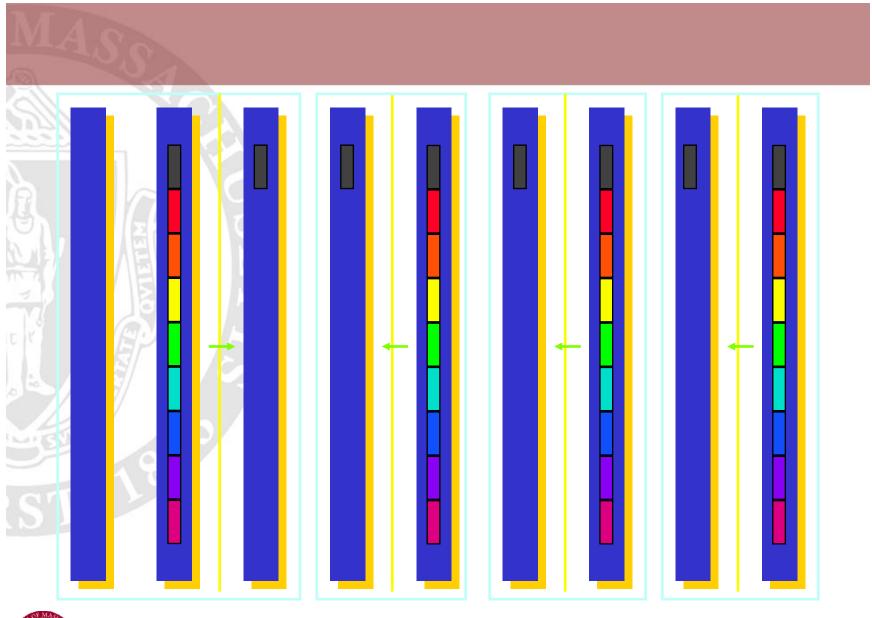




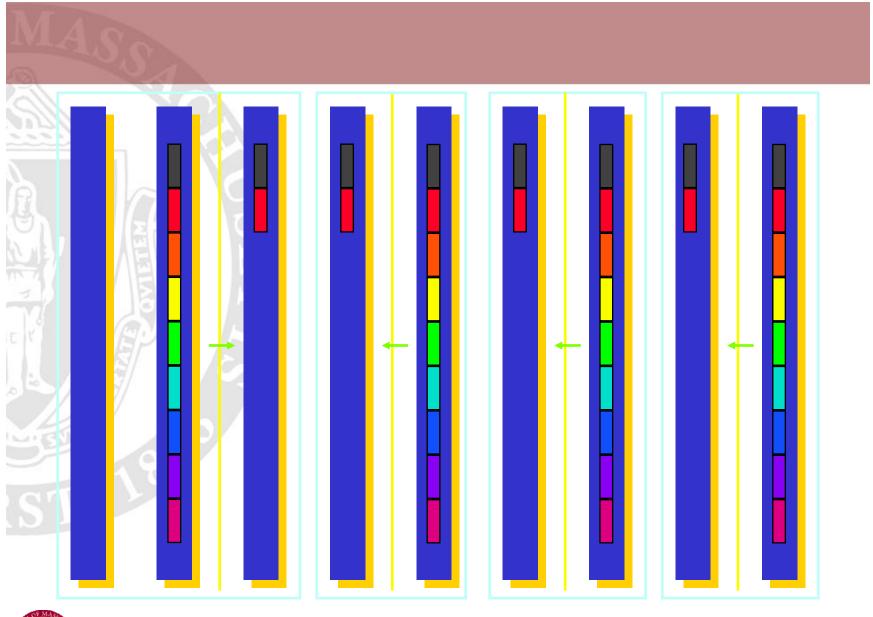




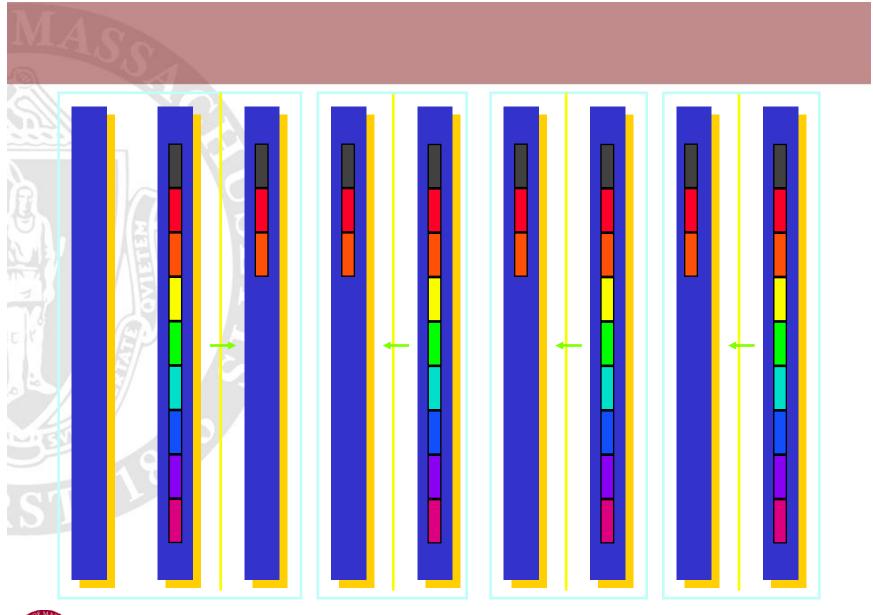




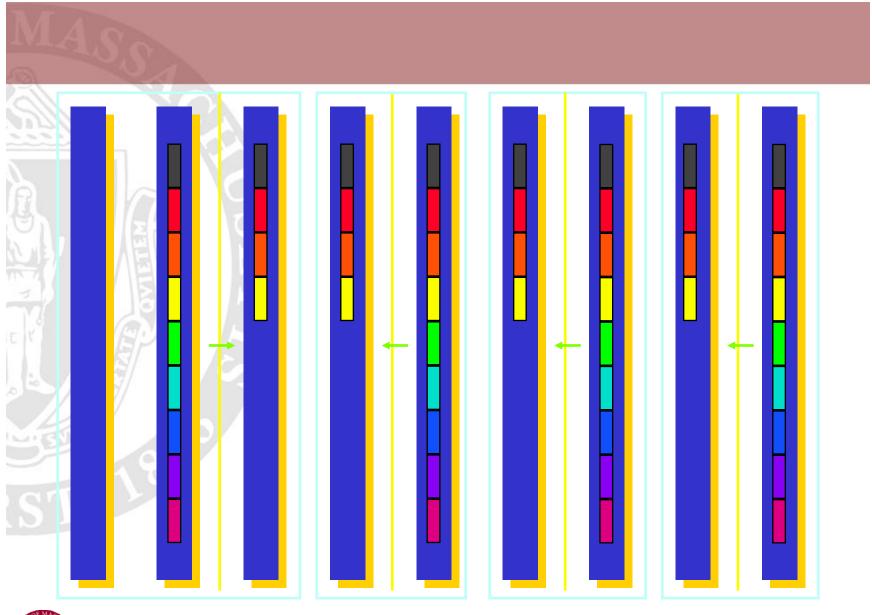




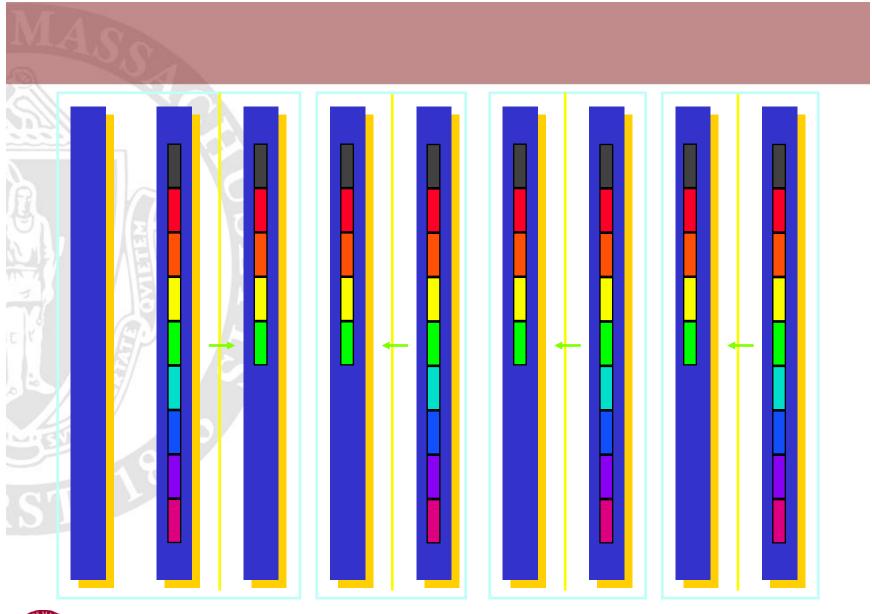




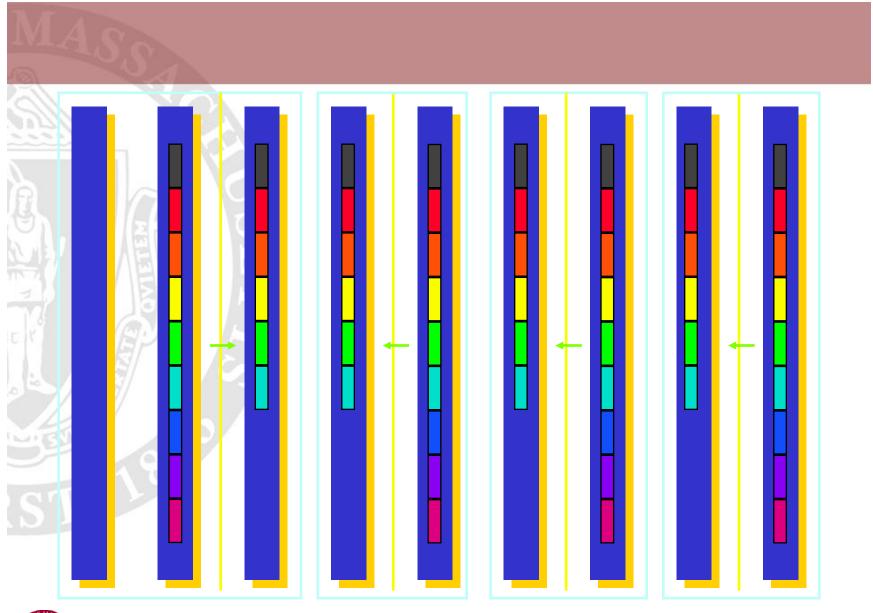




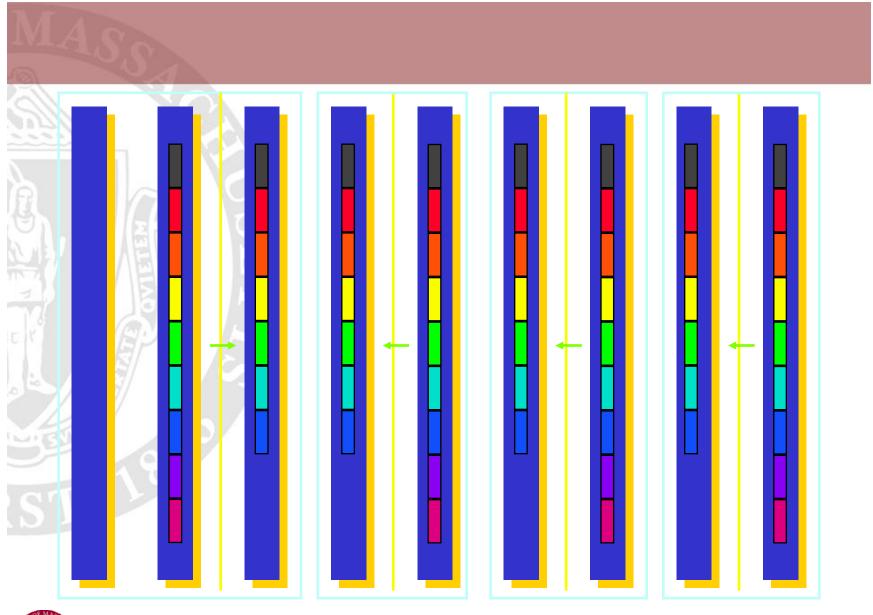




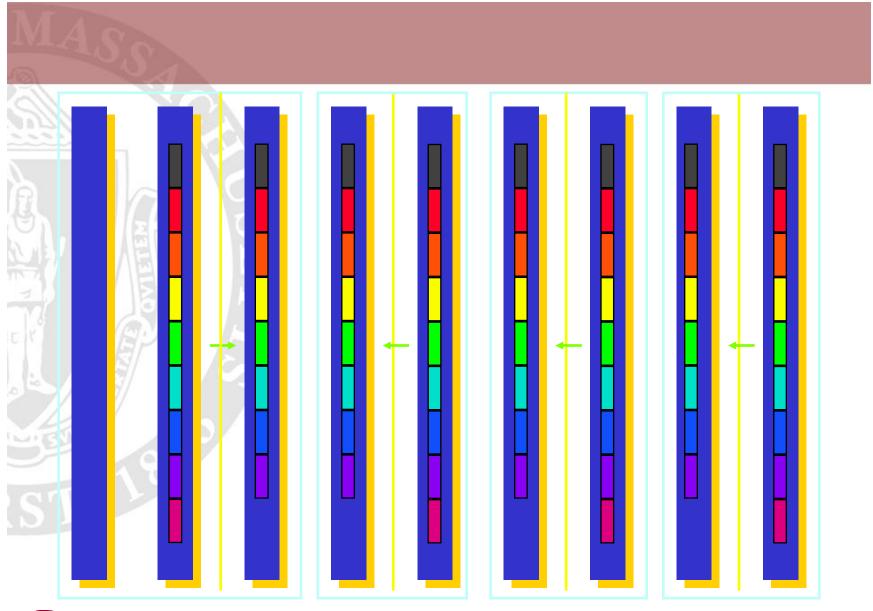




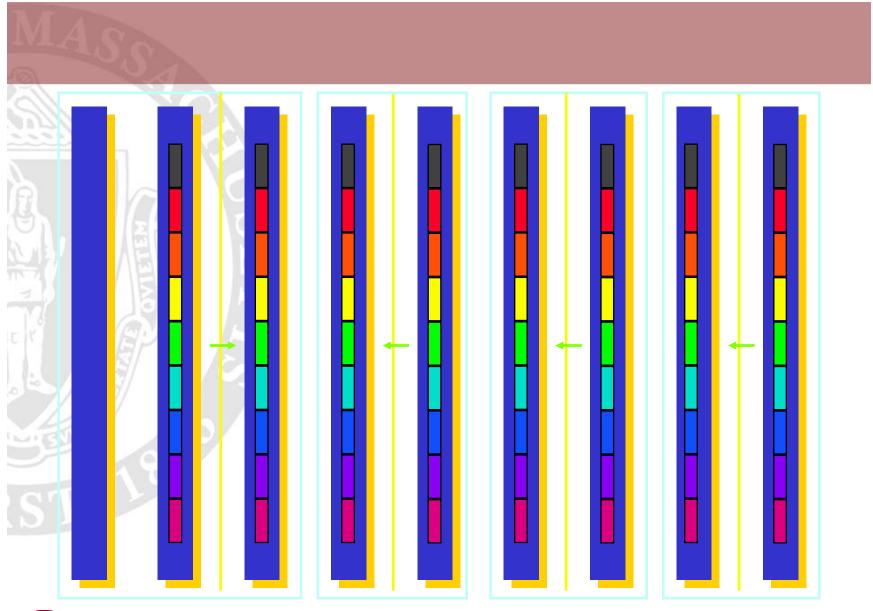




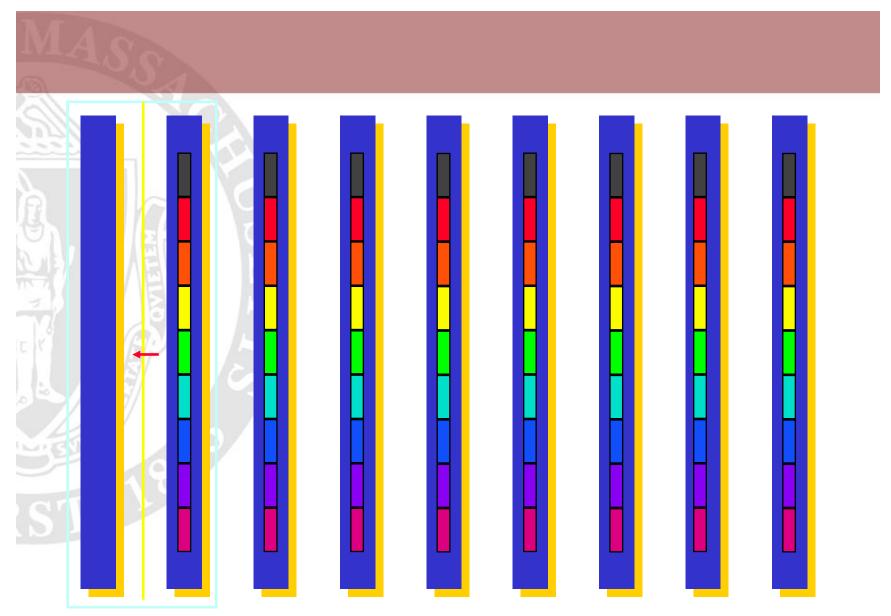




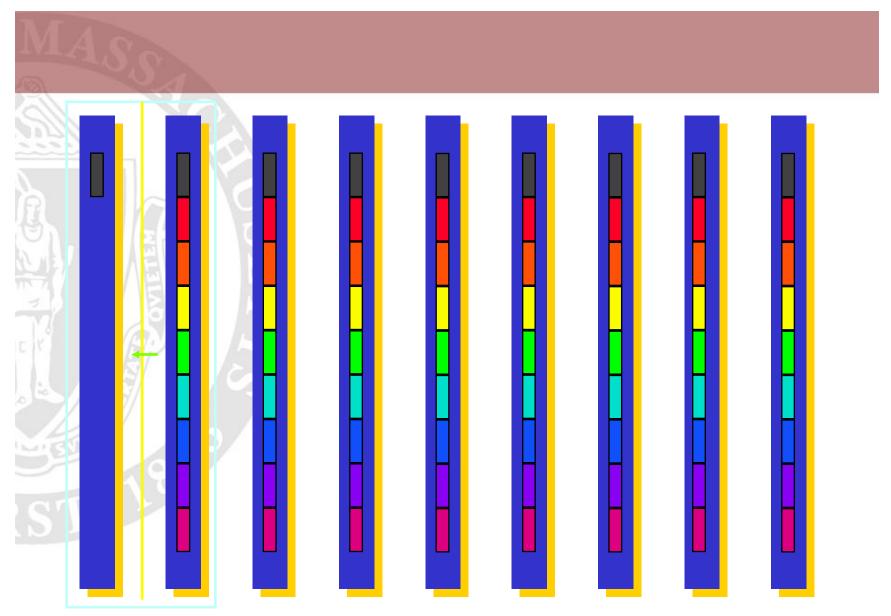




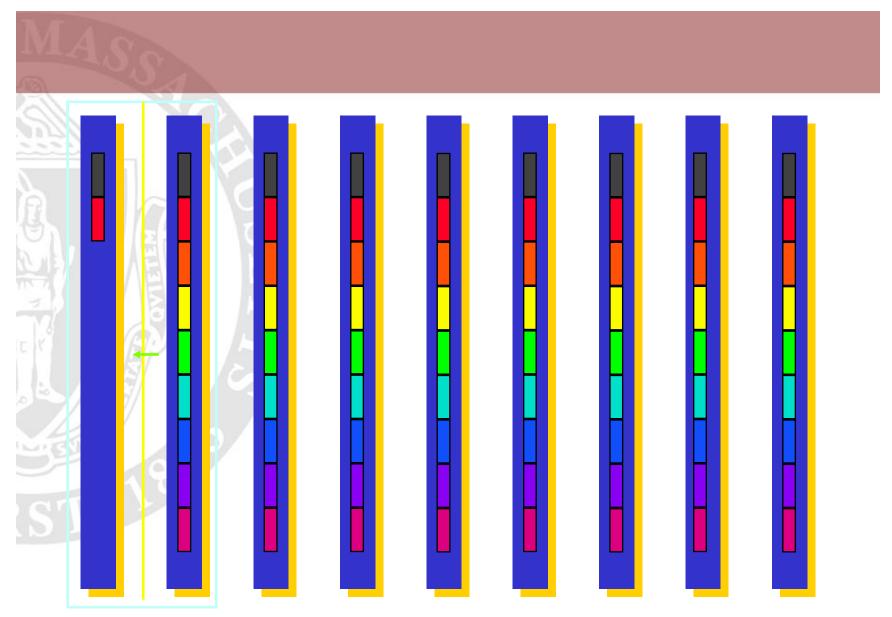




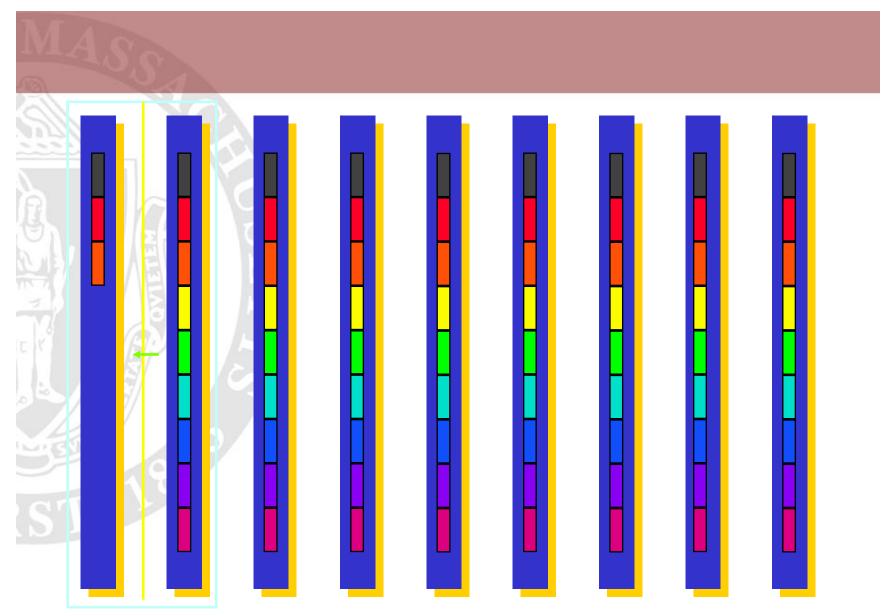




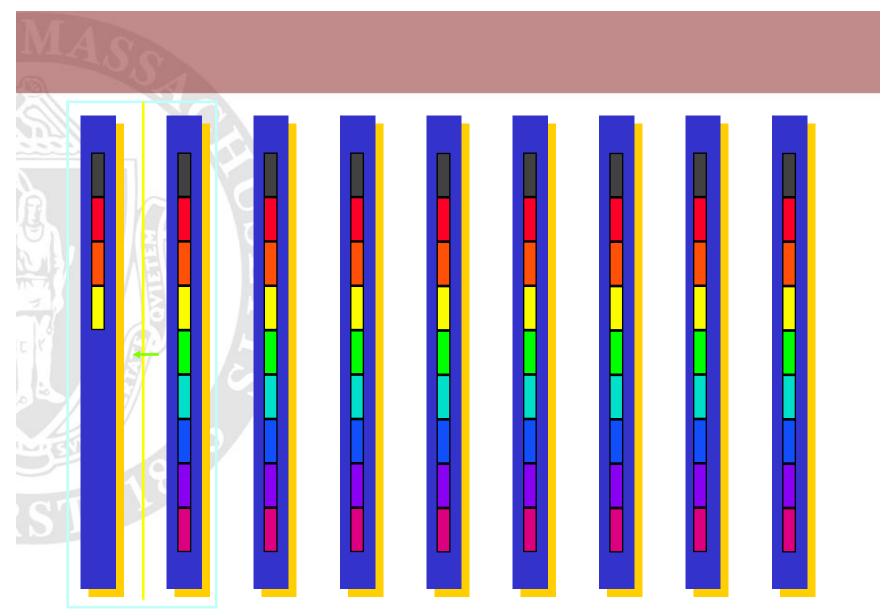




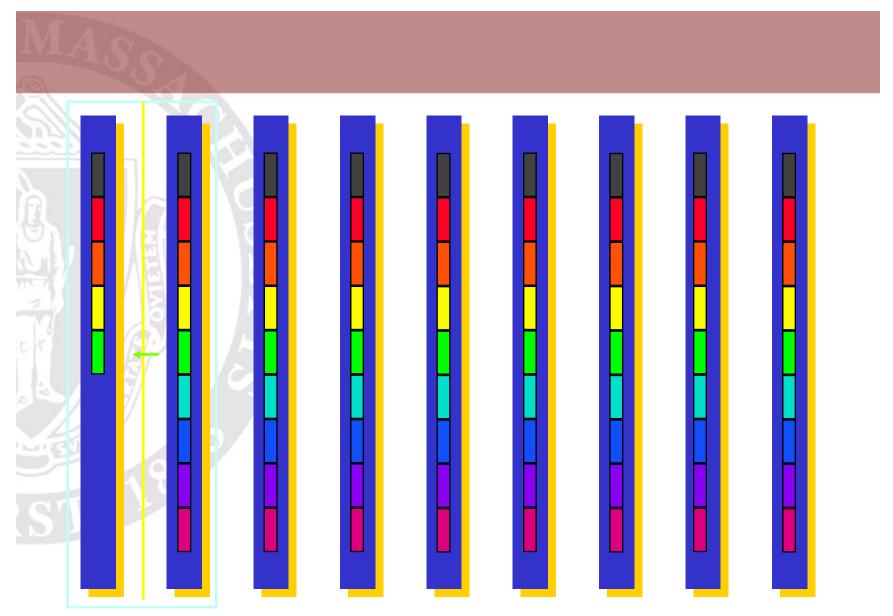




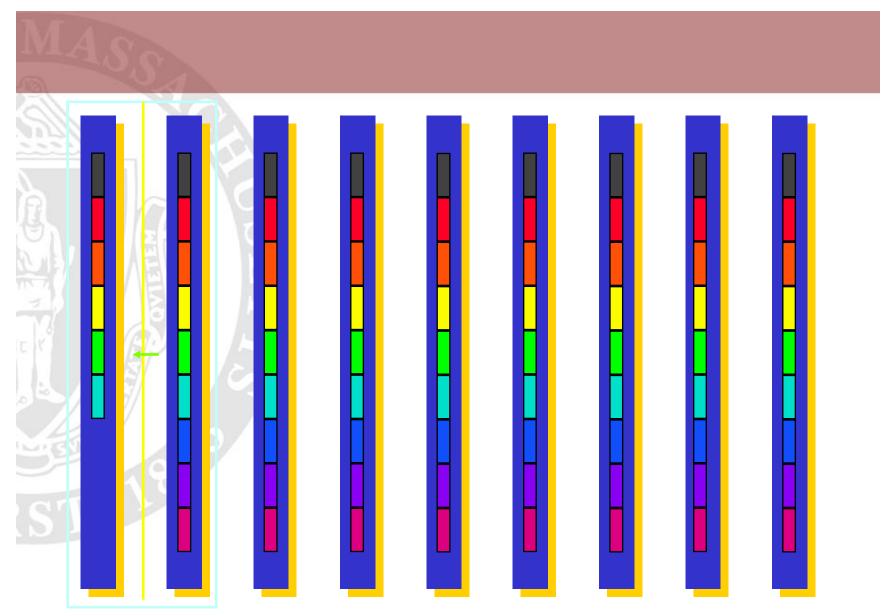




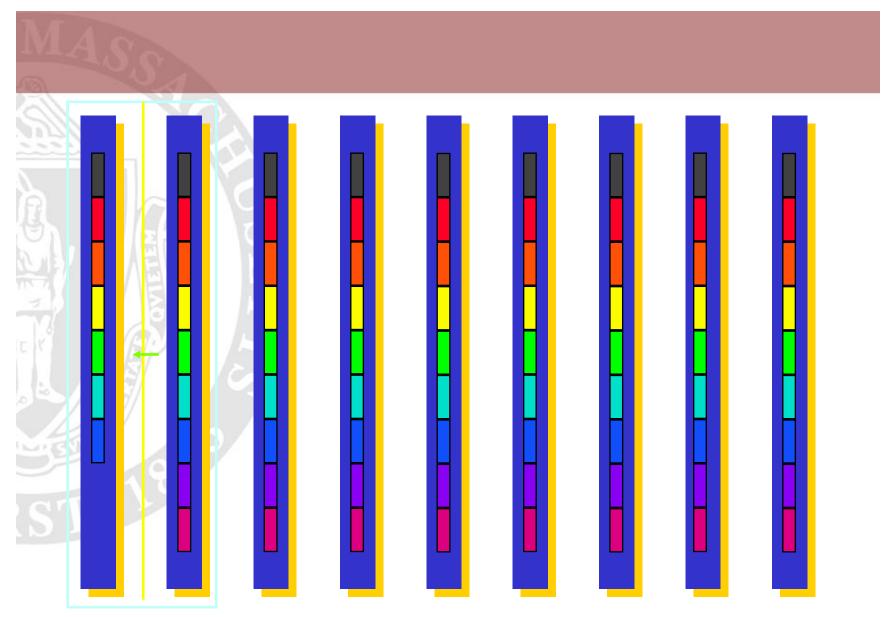




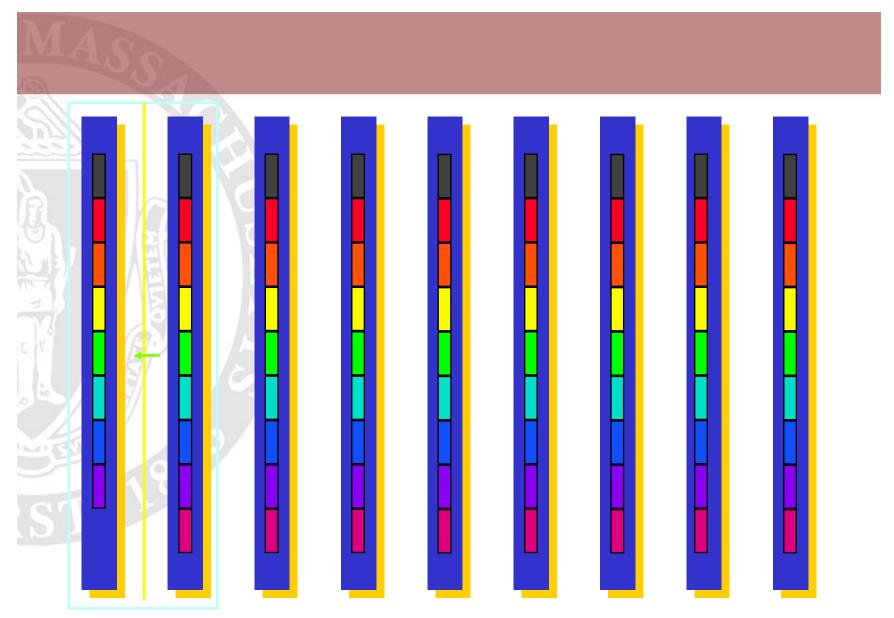




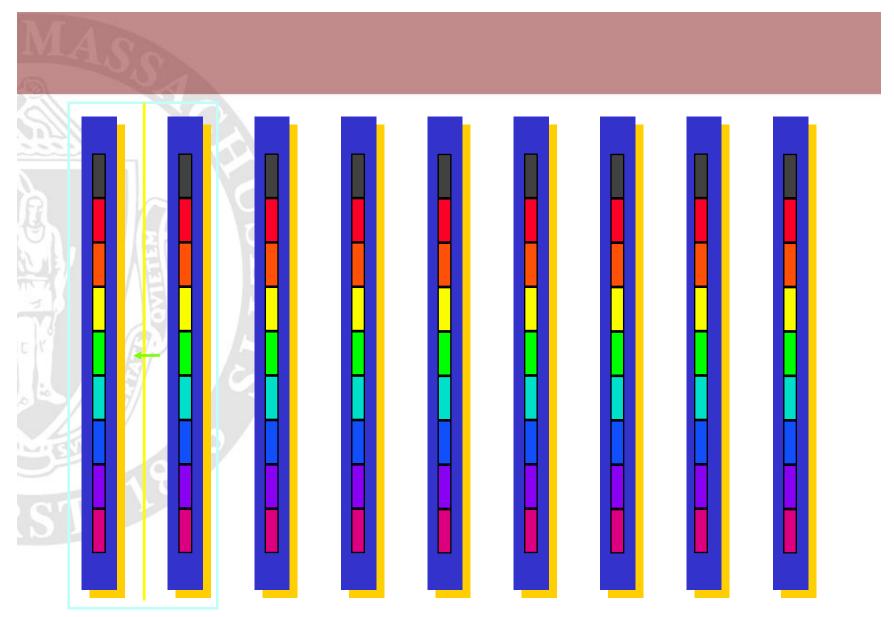




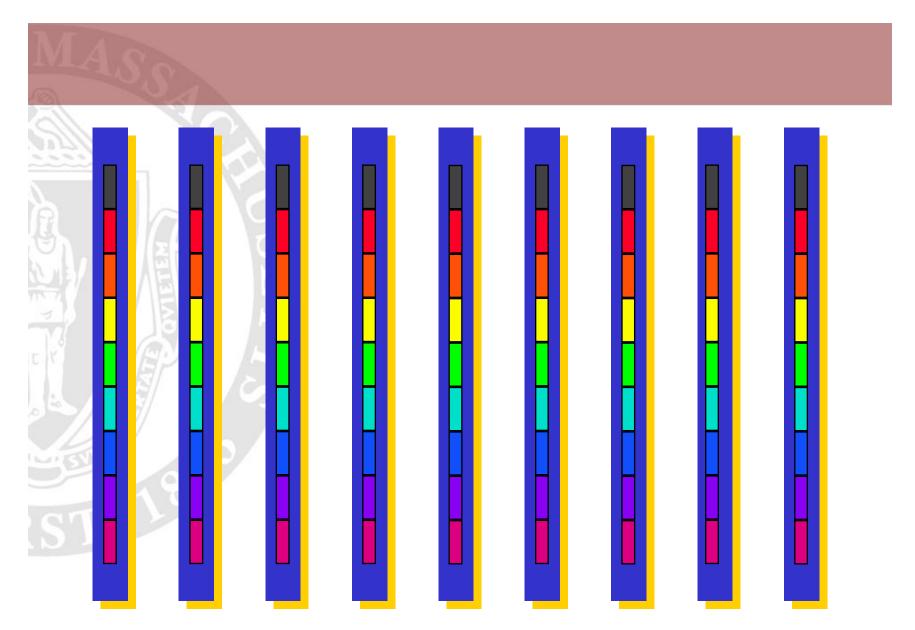










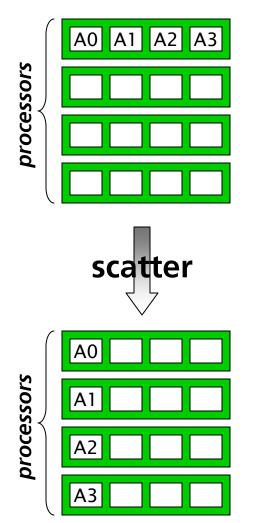




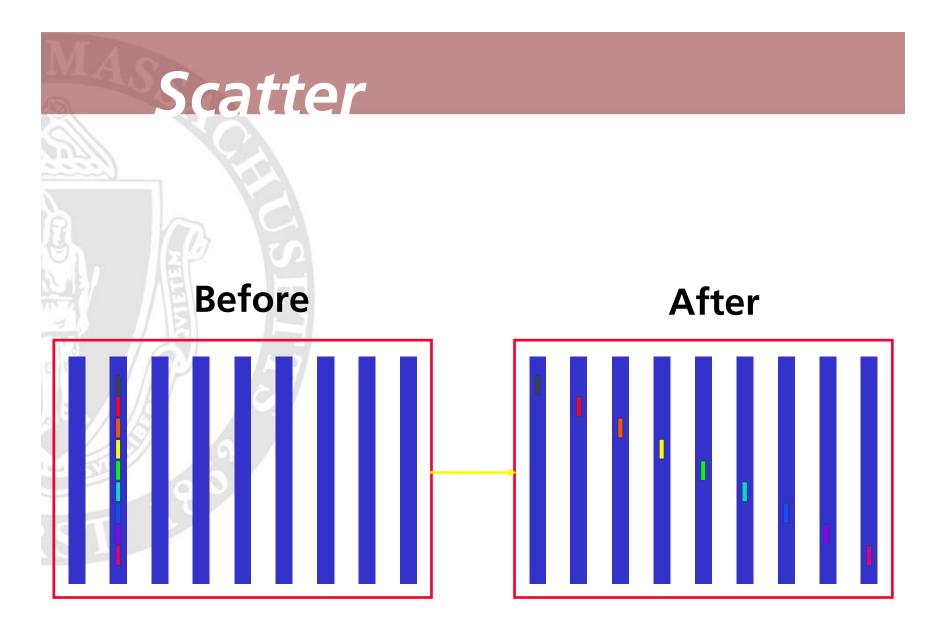
## From root: Scatter

 MPI\_Scatter: spread data across all processors

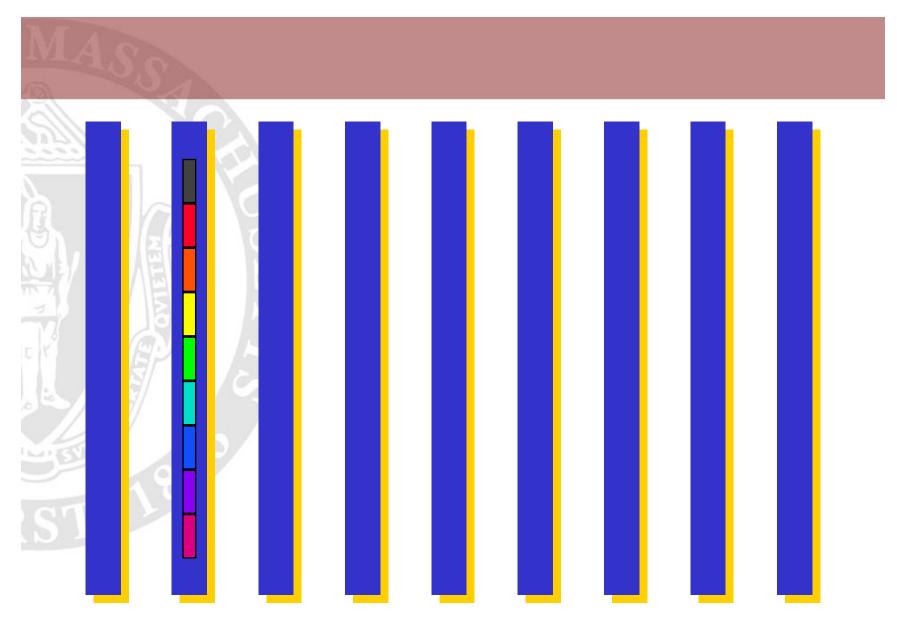
 variant –
 MPI\_Scatterv: scatters buffer in parts (specified in array)



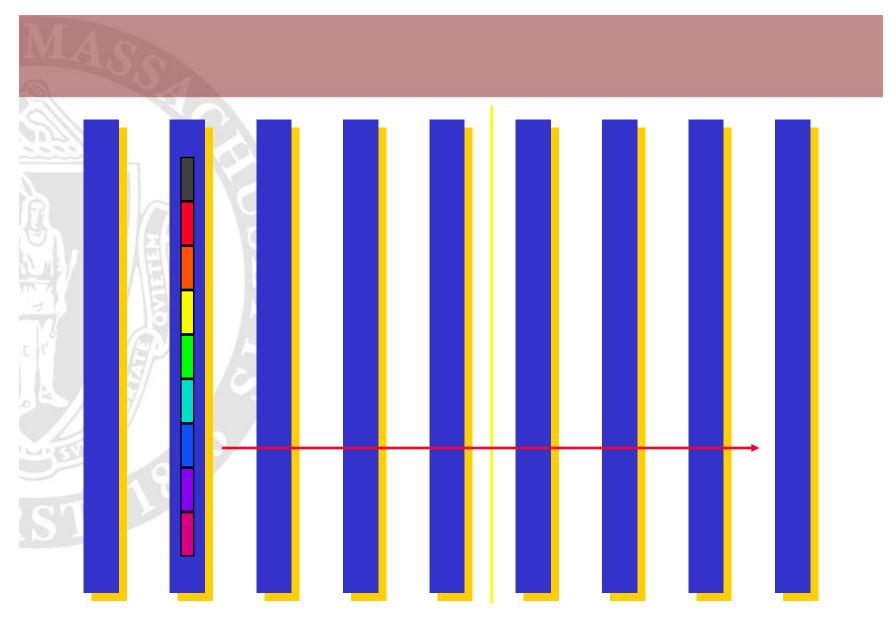




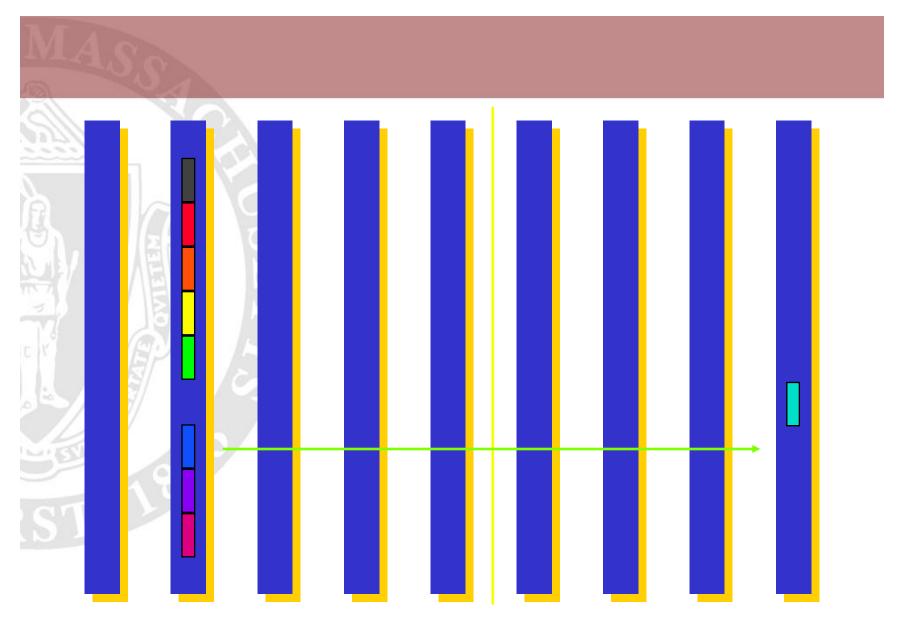




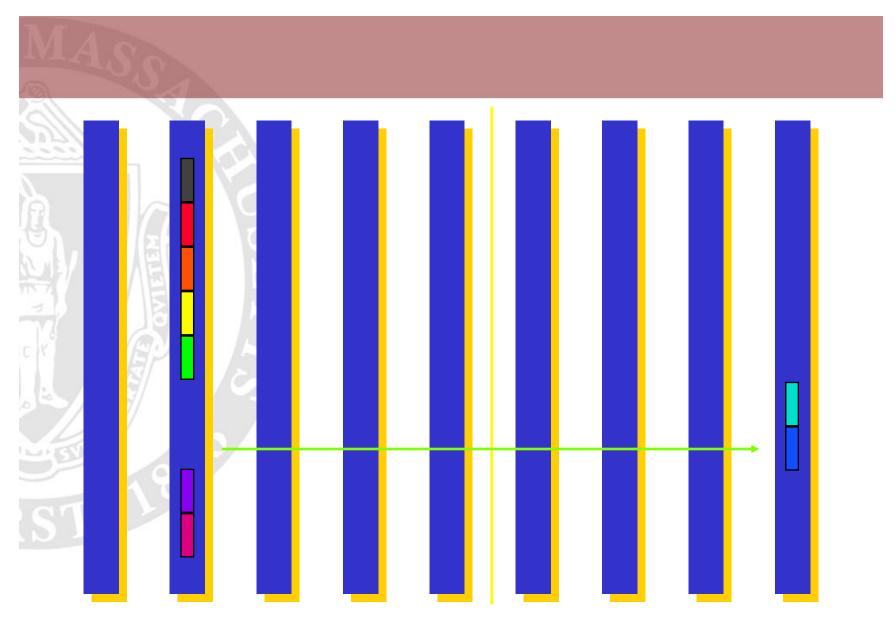




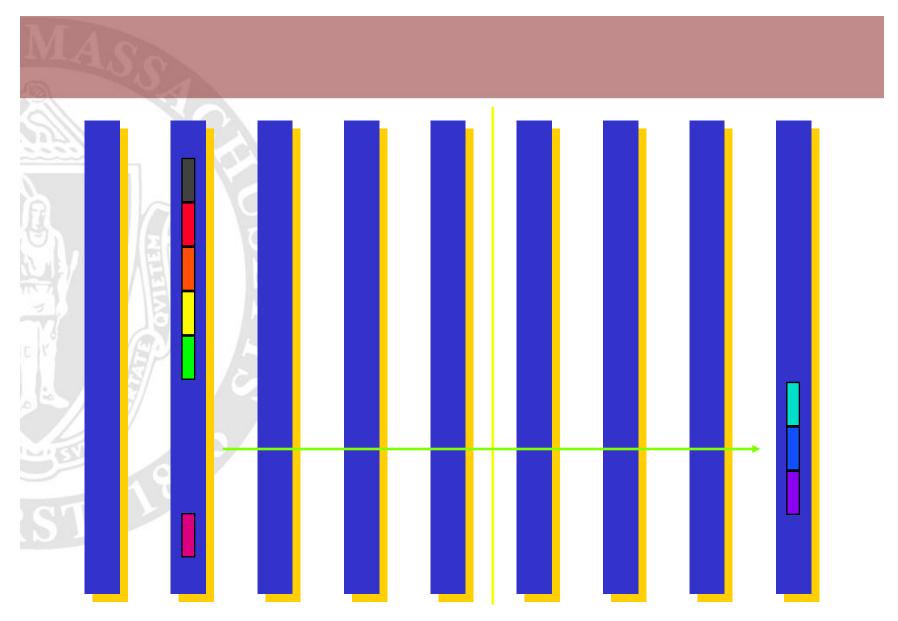




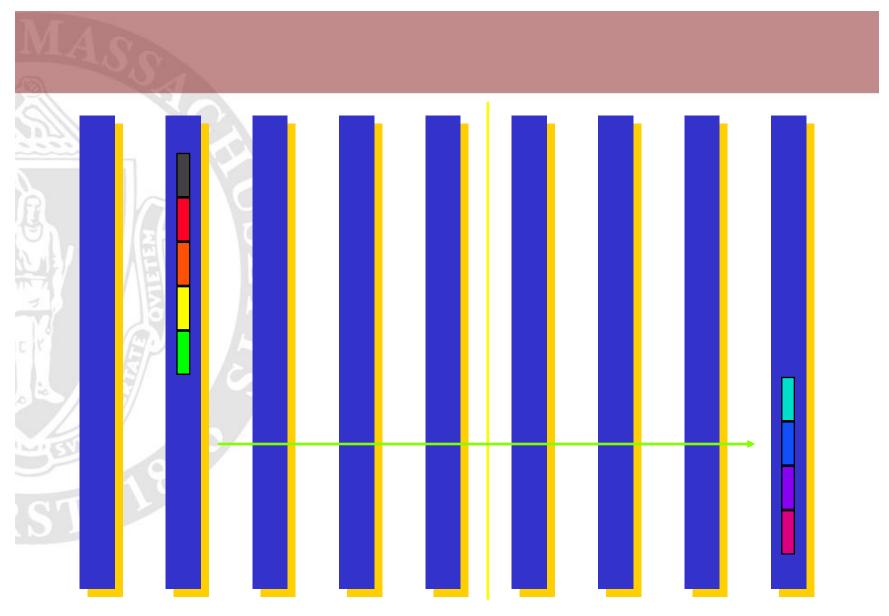




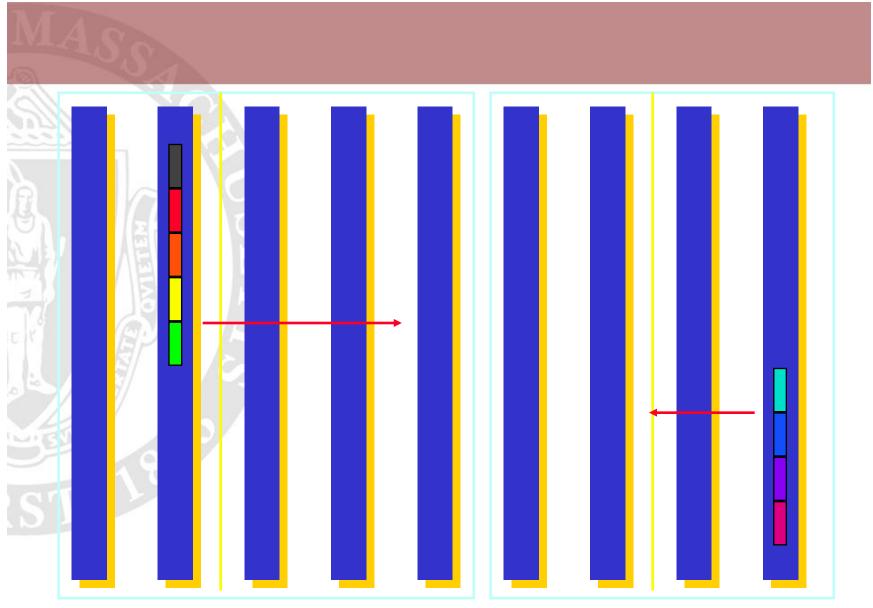




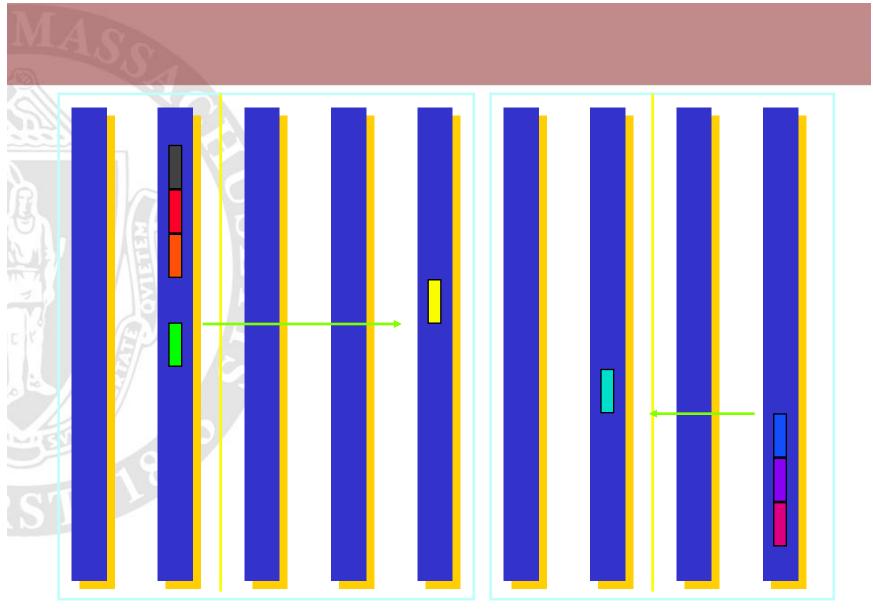




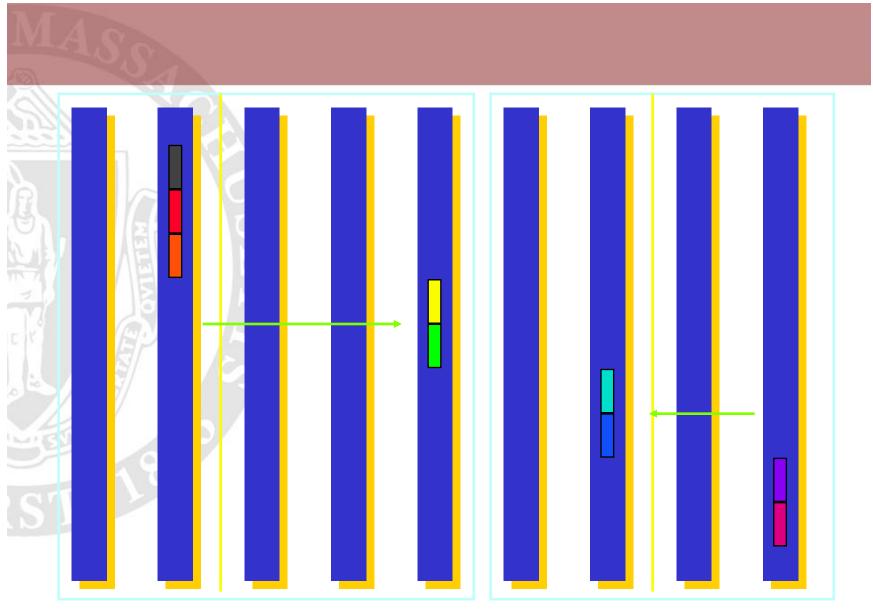




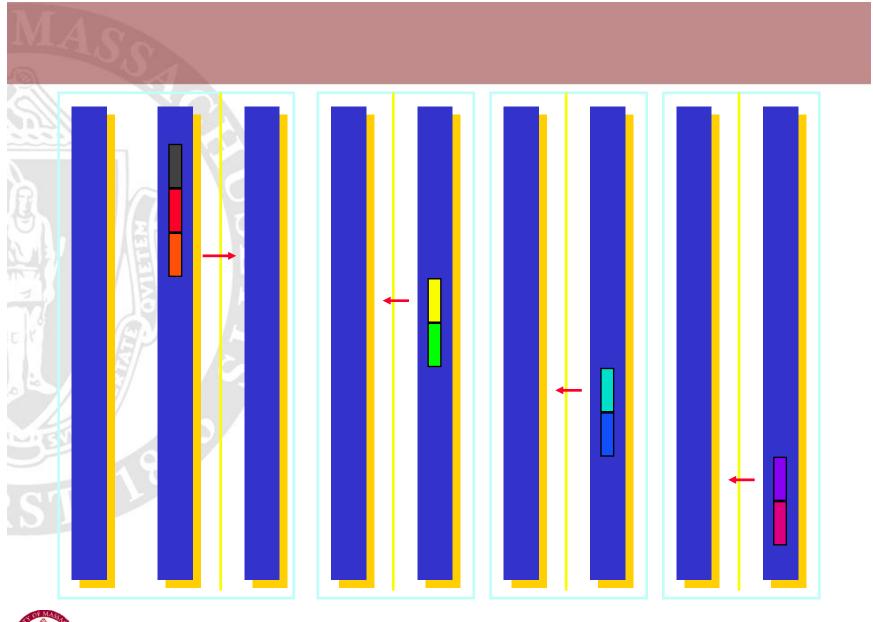




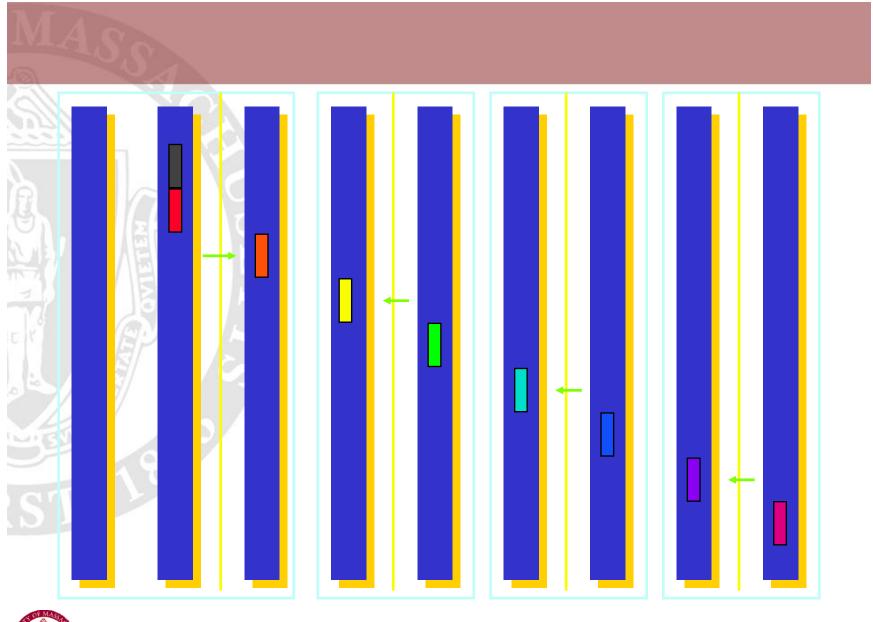




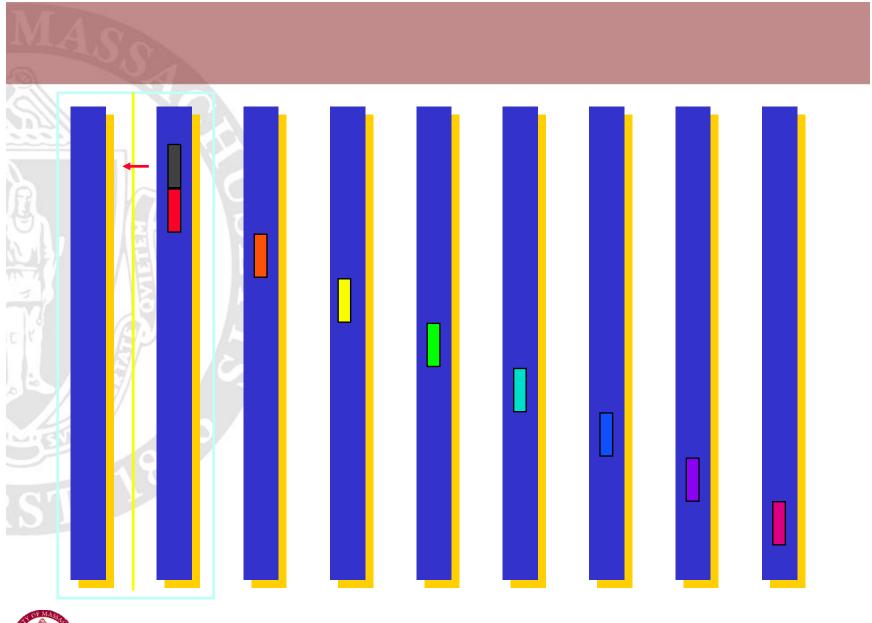




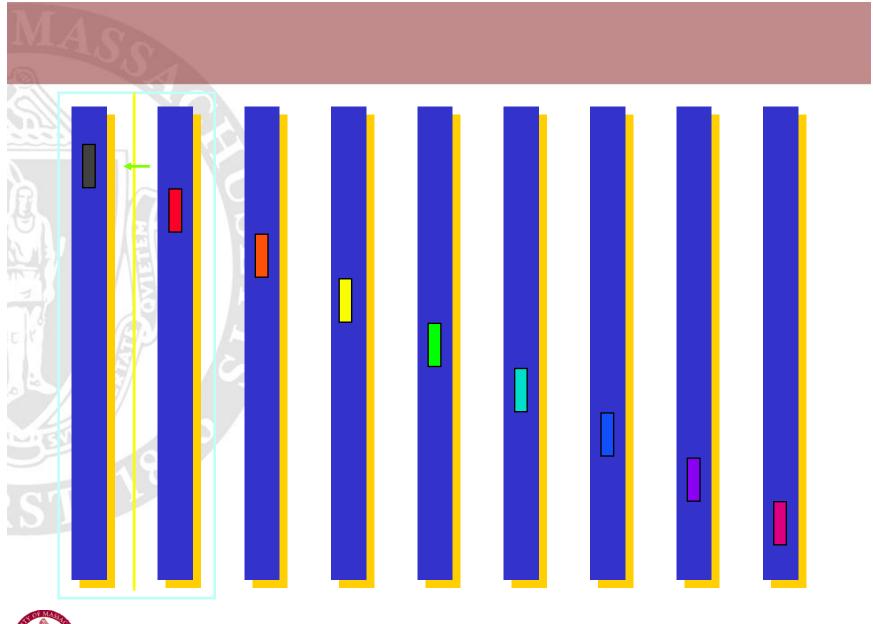




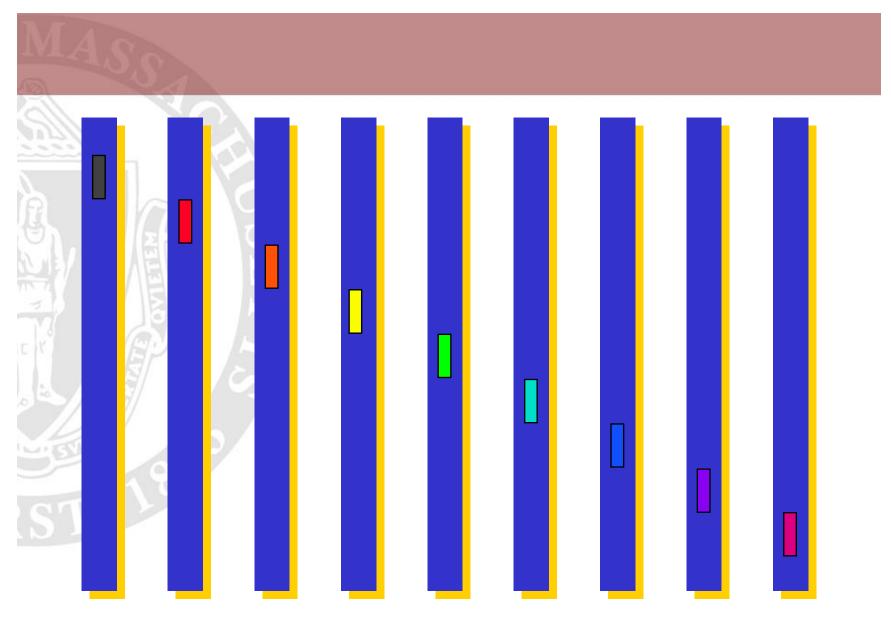










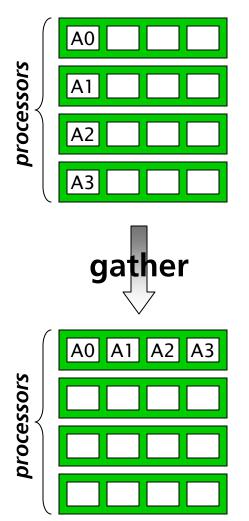




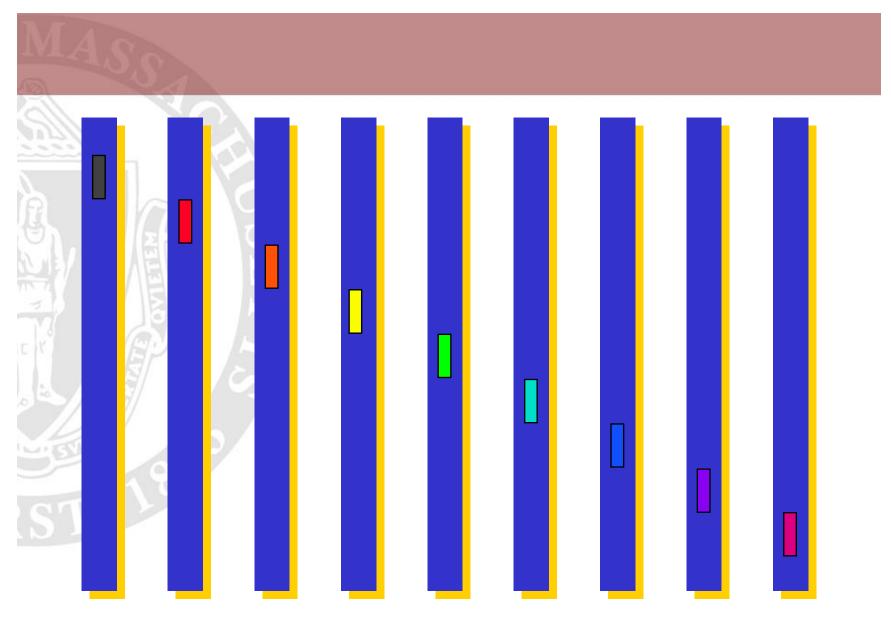
## To root: Gather

 MPI\_Gather: gather data from all processors to one processor

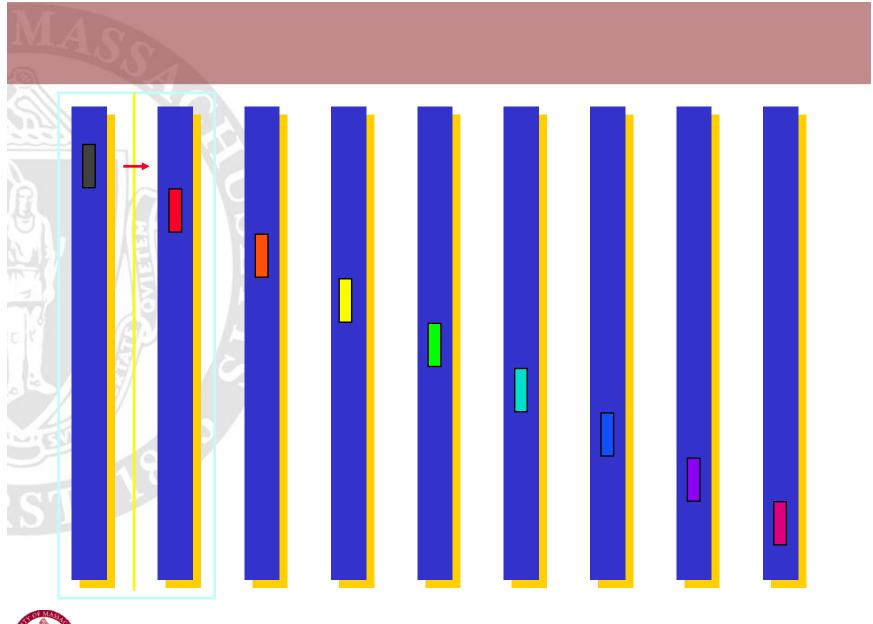
variant –
 MPI\_Gatherv



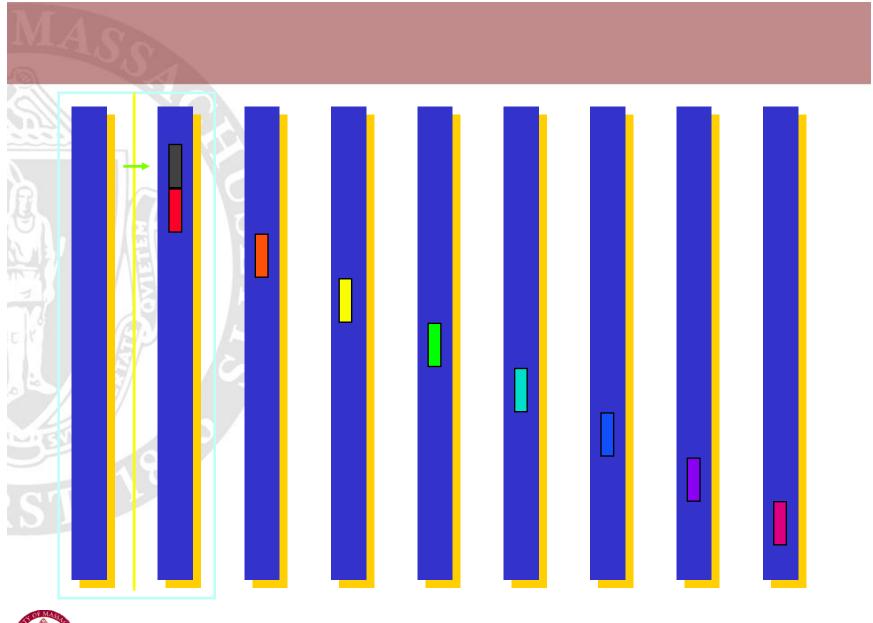




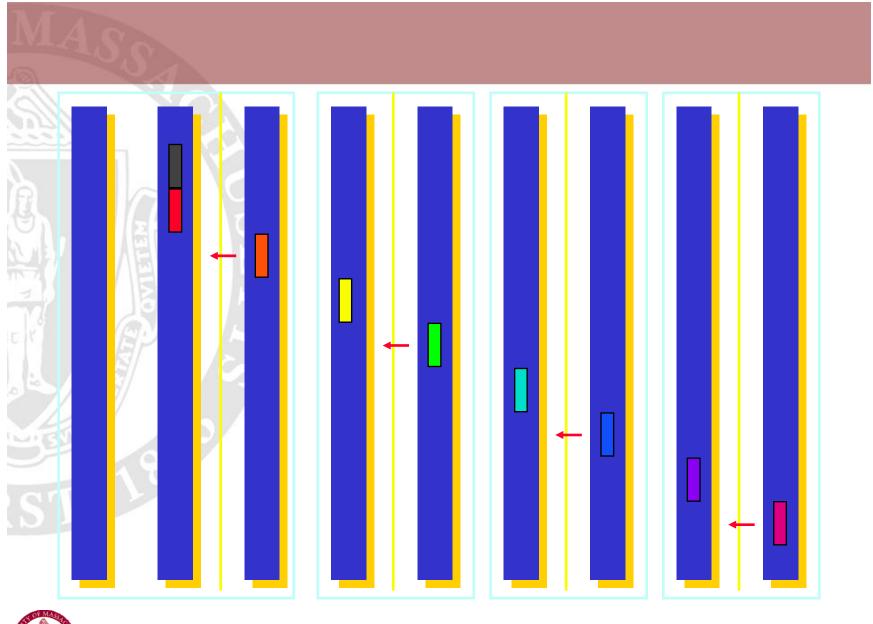




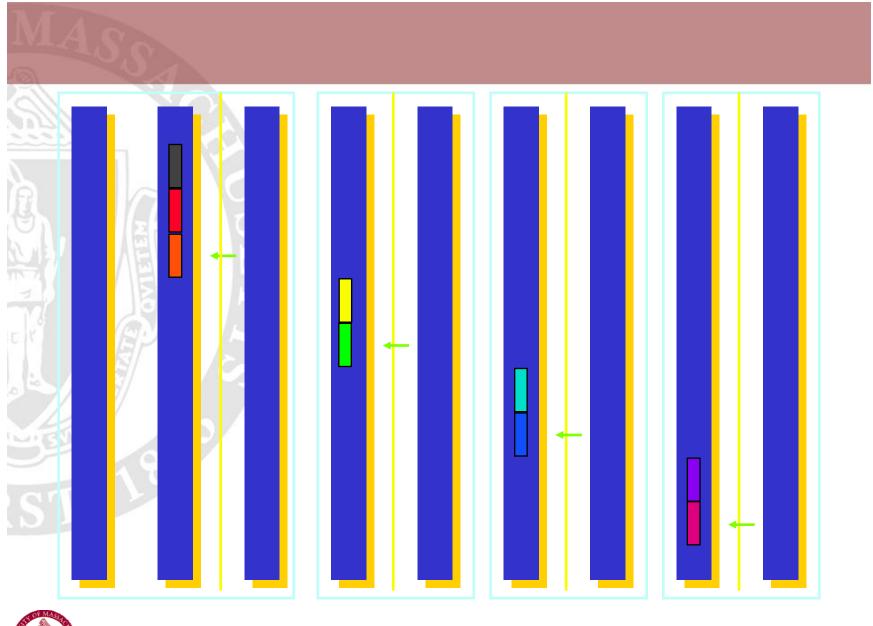




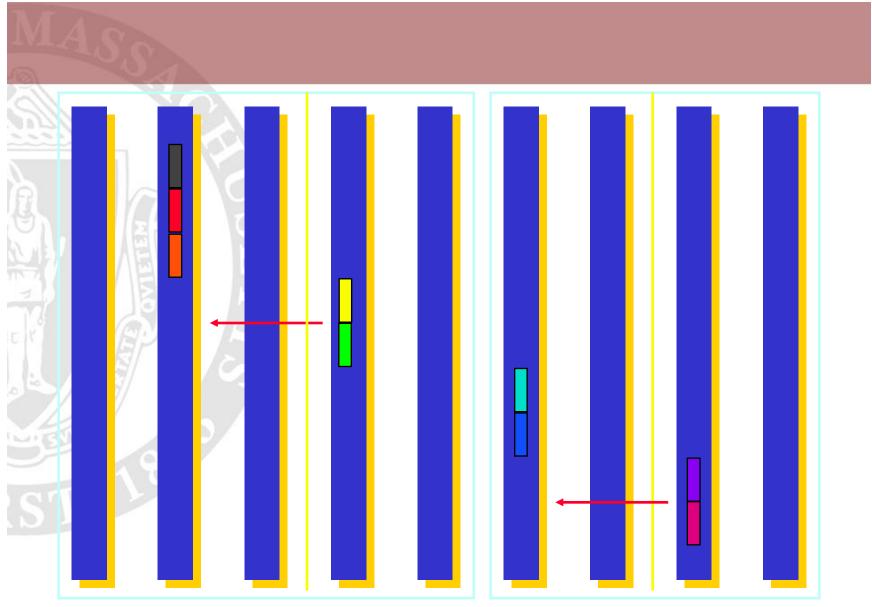




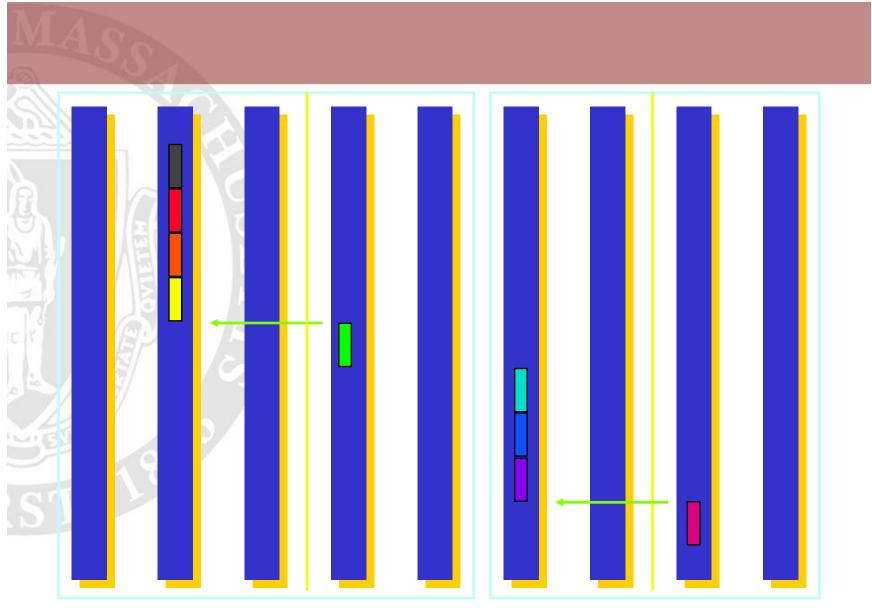




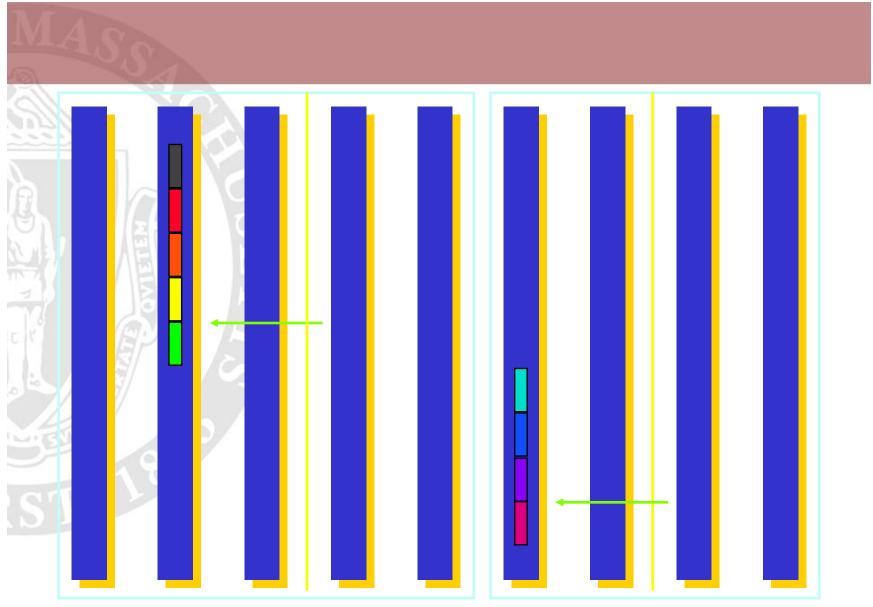




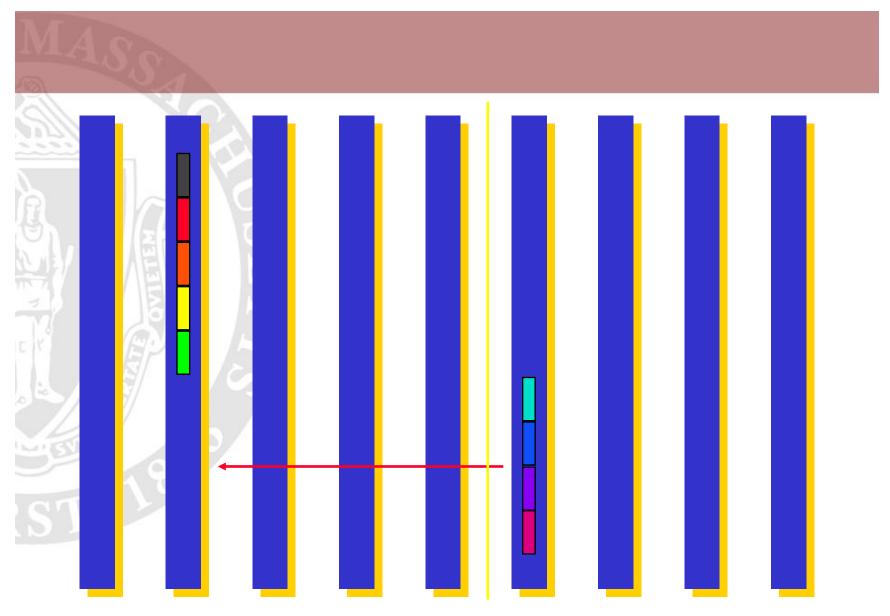




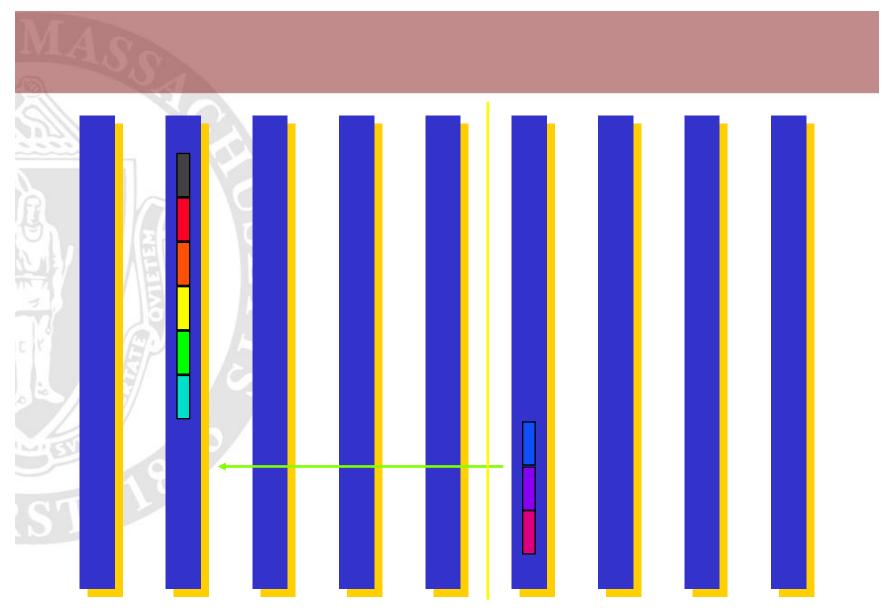




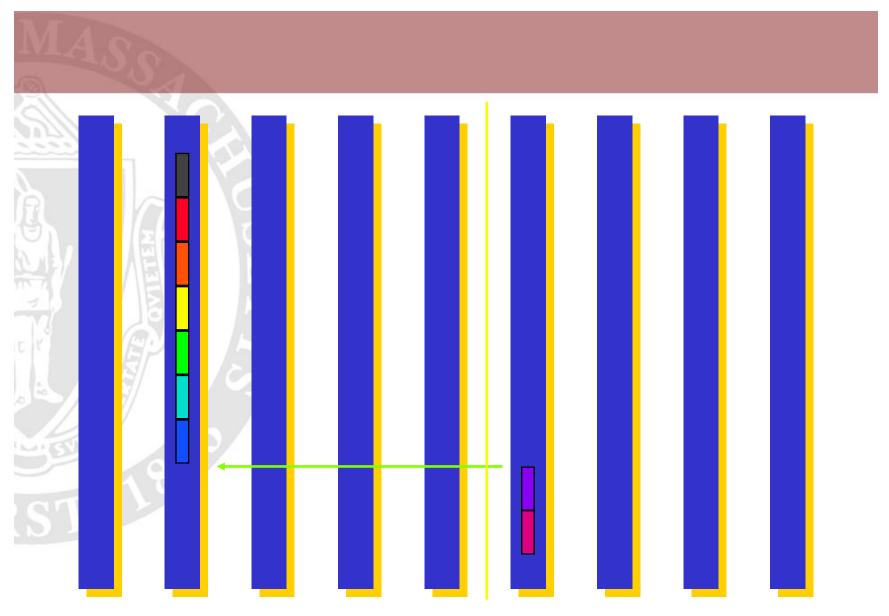




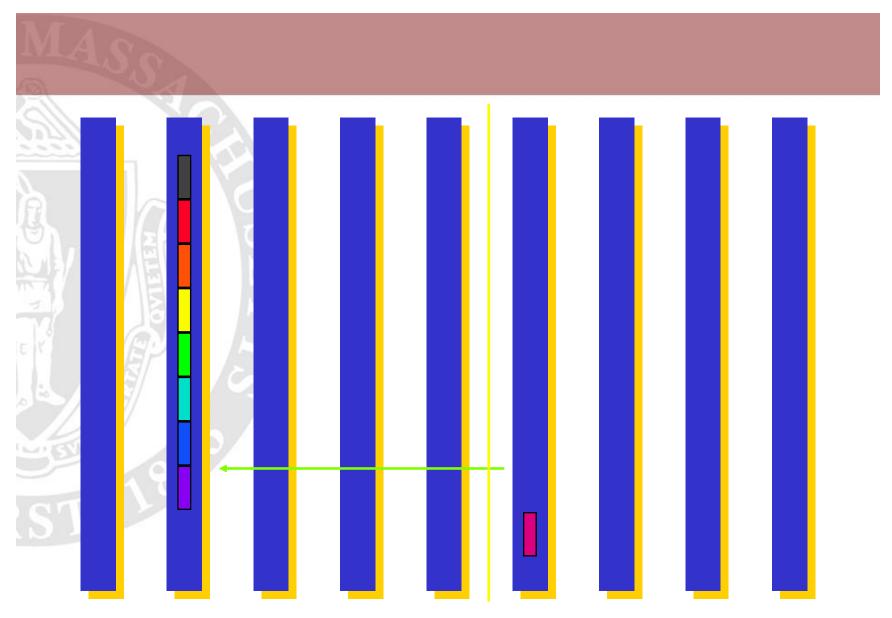




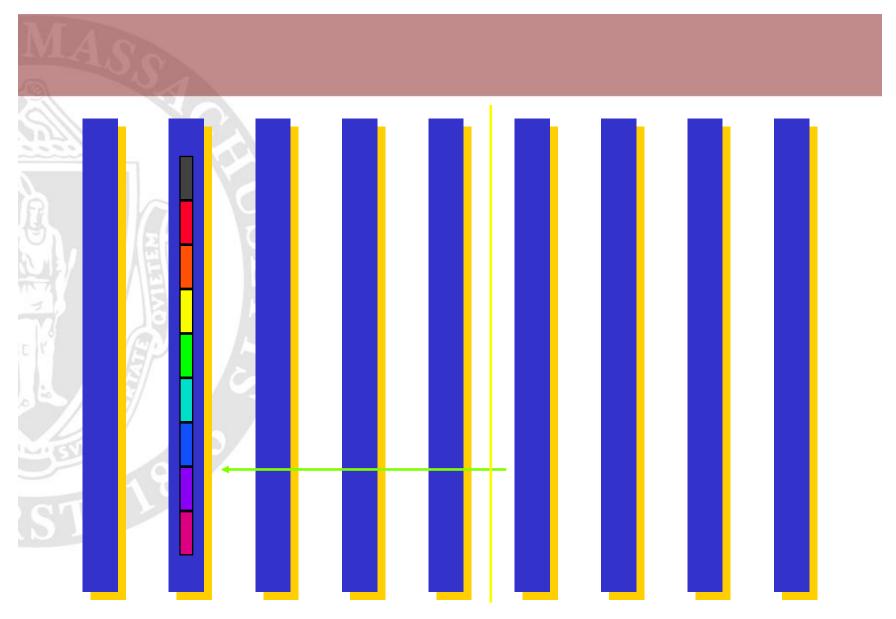




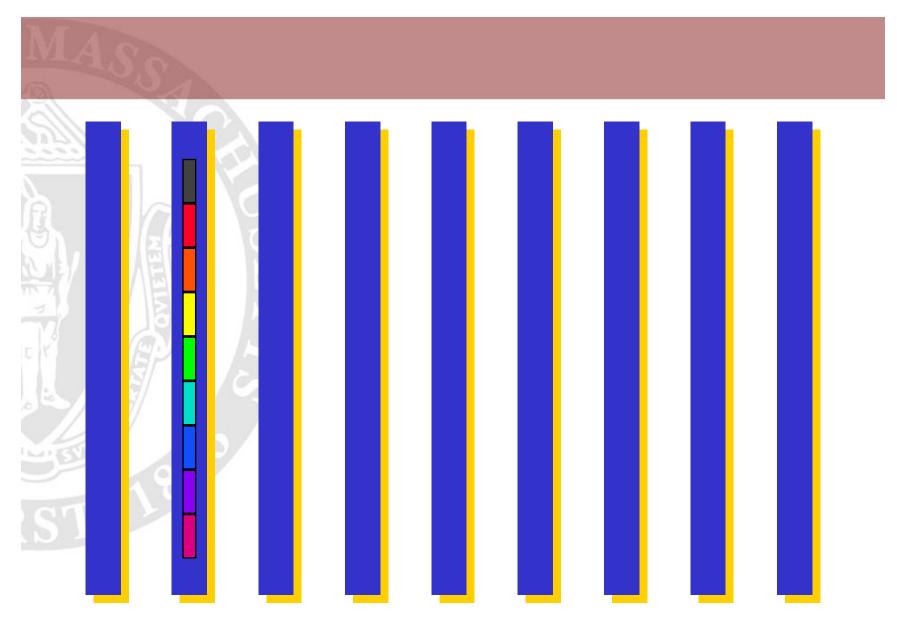








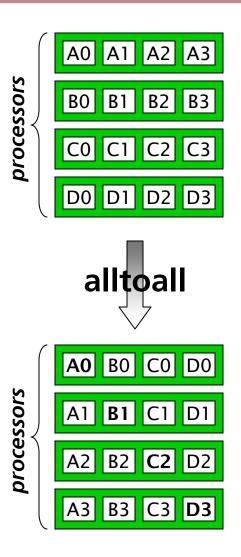






# All: Alltoall

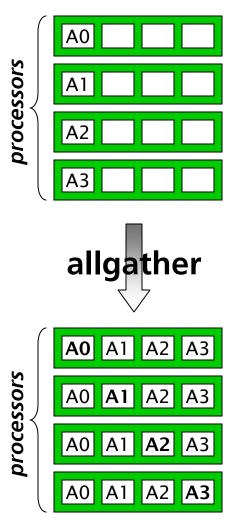
 MPI\_Alltoall: send data from all processors to all processors





## All: Allgather

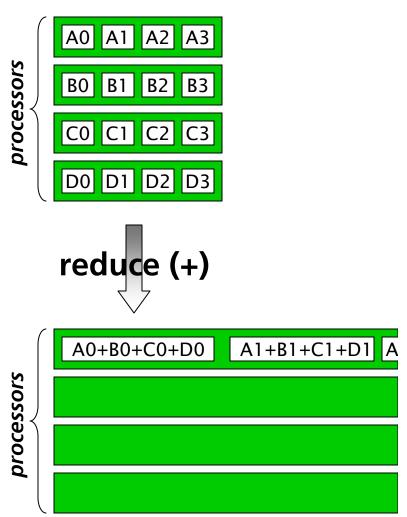
 MPI\_Allgather: send data from all processors to one vector across all processors





### **Reductions: Reduce**

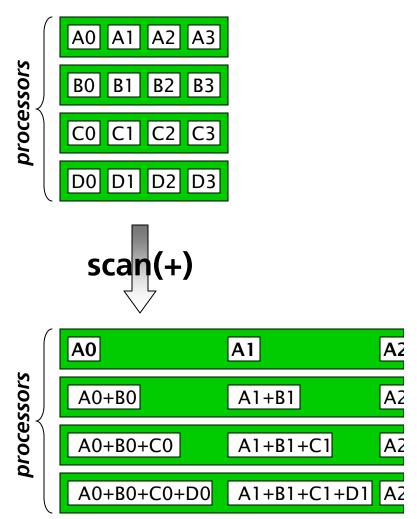
 MPI\_Reduce: collect data from all processors, apply operator to each, put on one processor



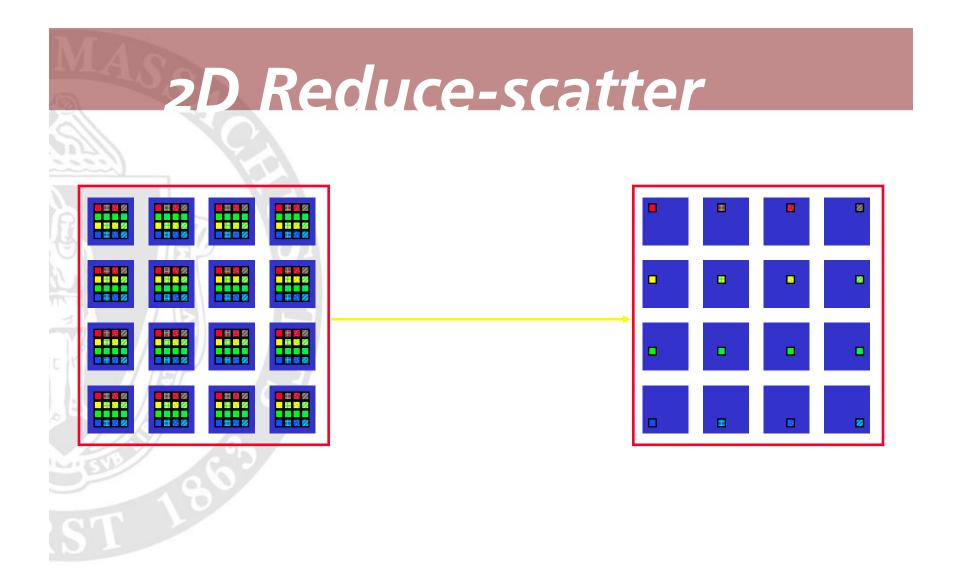


### **Reductions: Scan**

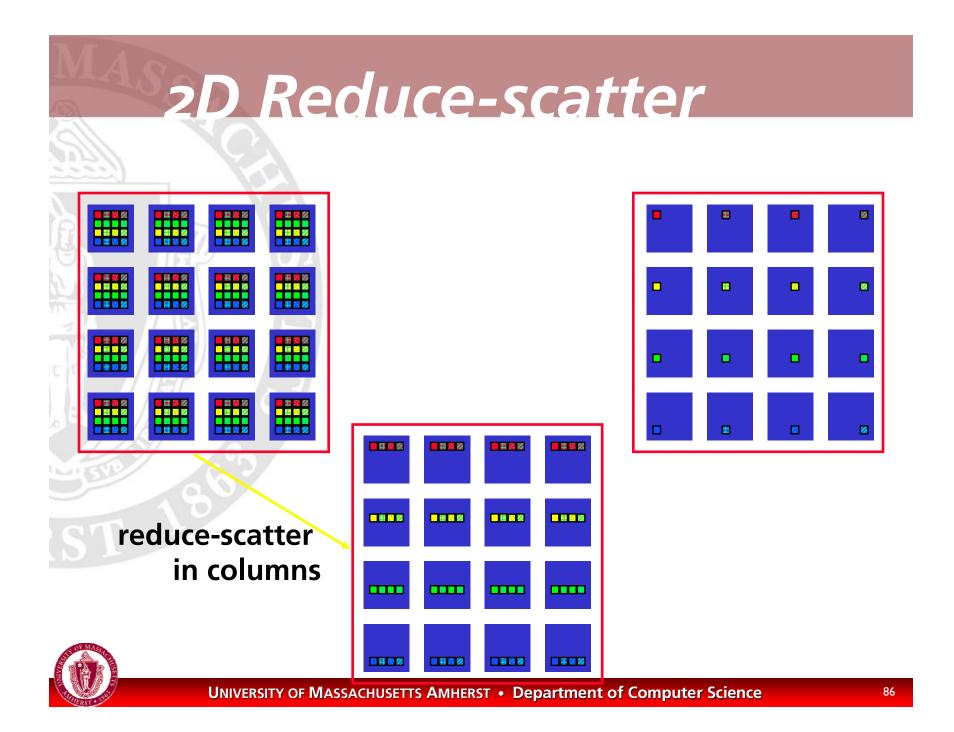
 MPI\_Scan: apply partial reductions
 Operation applied to increasinglylong prefixes

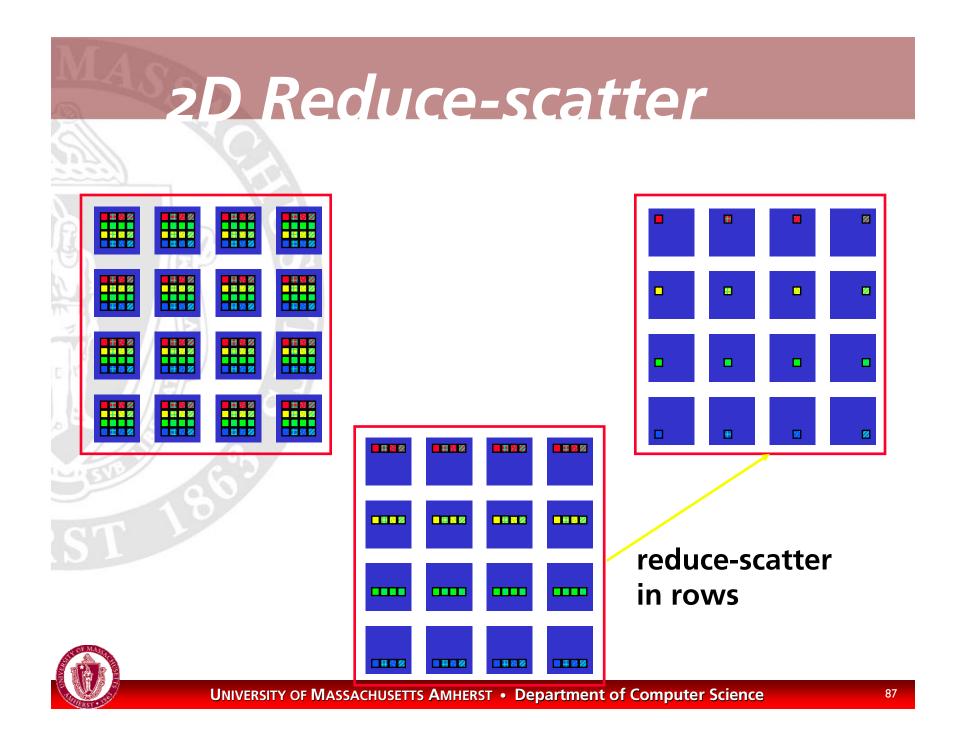




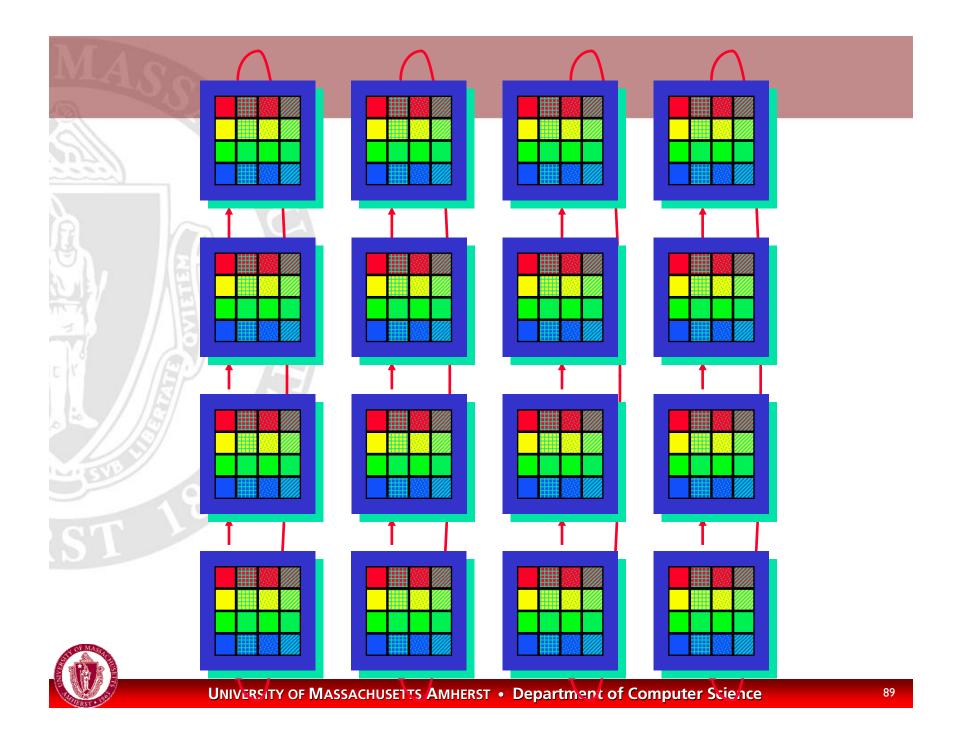


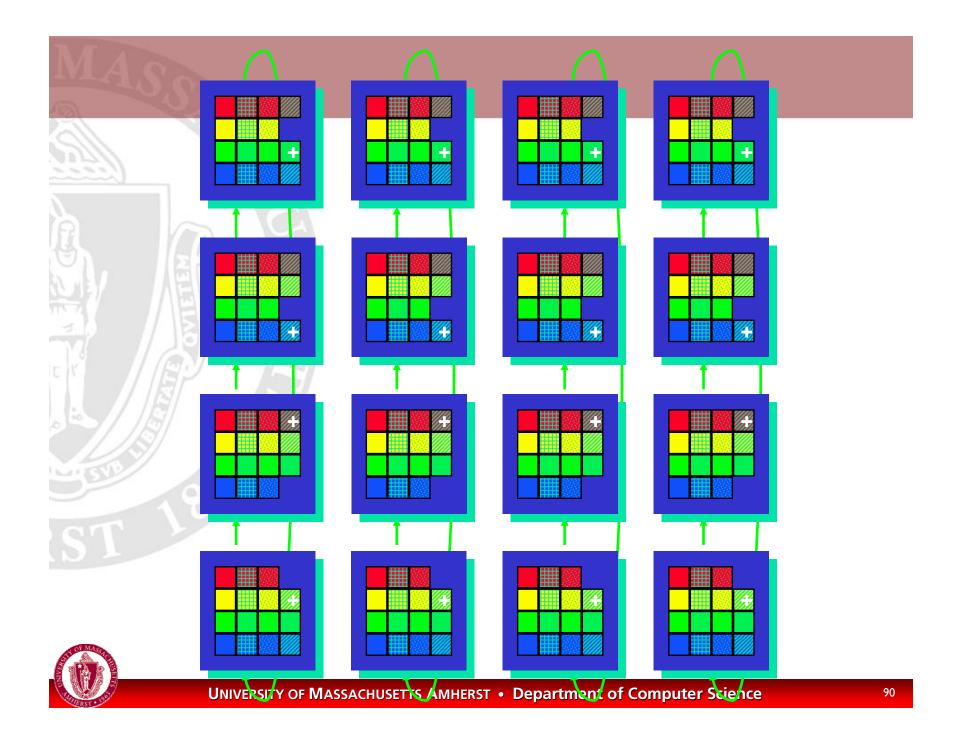


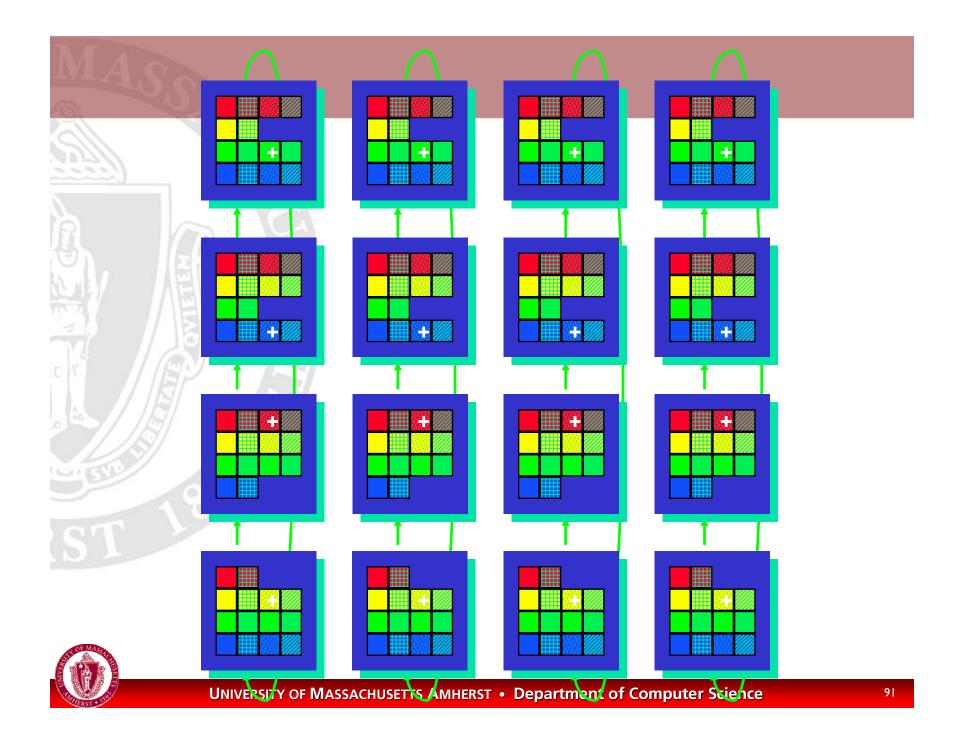


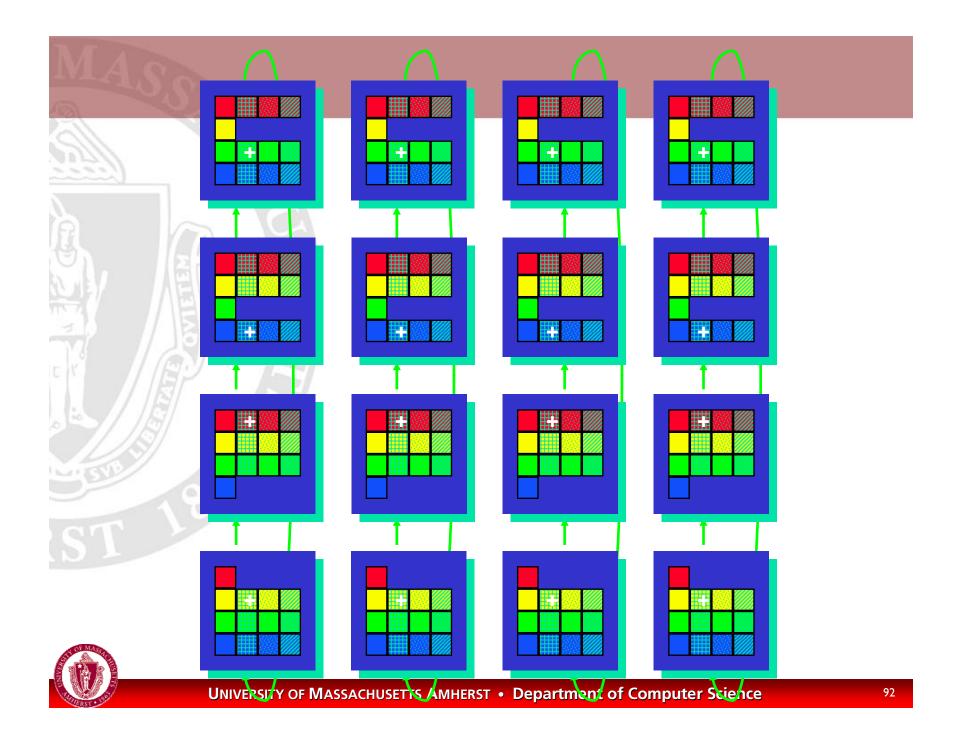


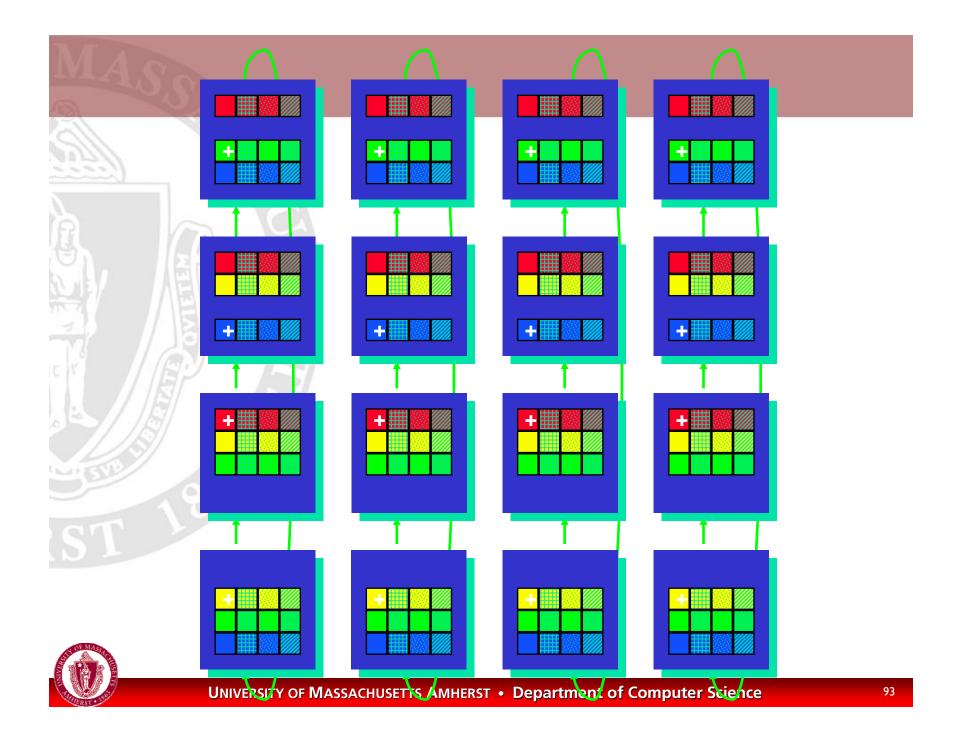
MASS					
ST 1					
	University of MA	SSACHUSETTS AMHERST	<ul> <li>Department of C</li> </ul>	omputer Science	88

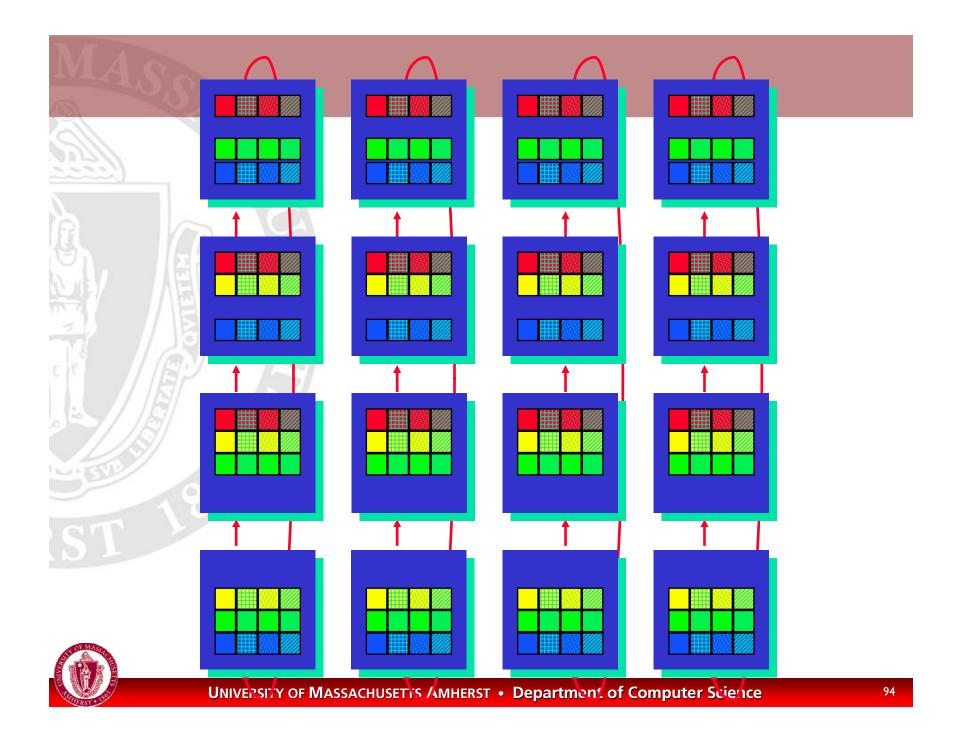


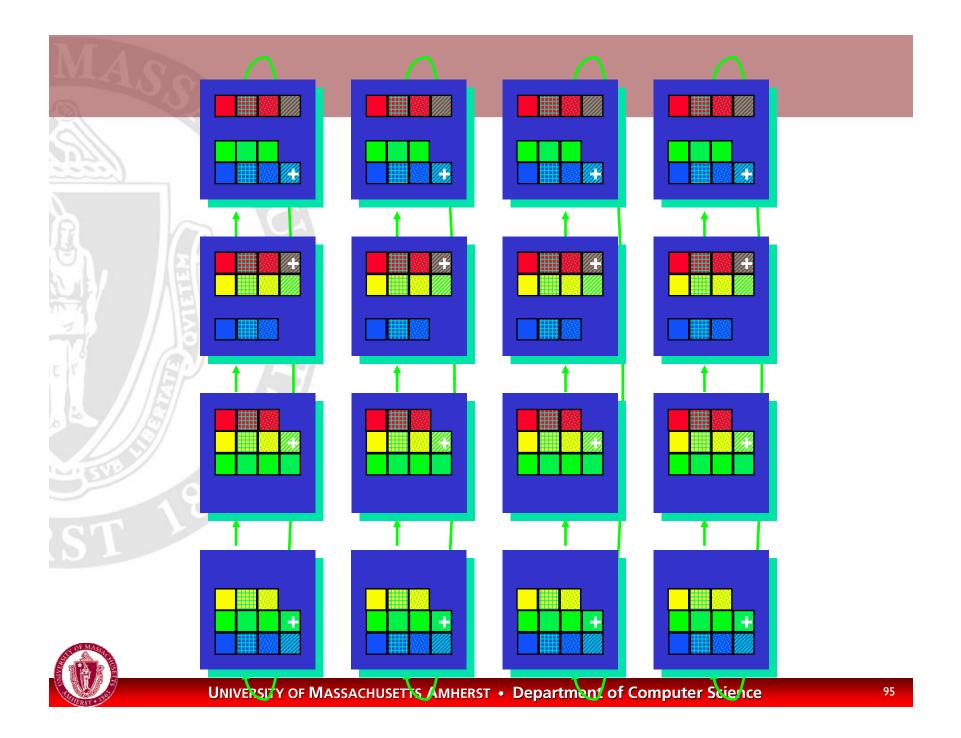


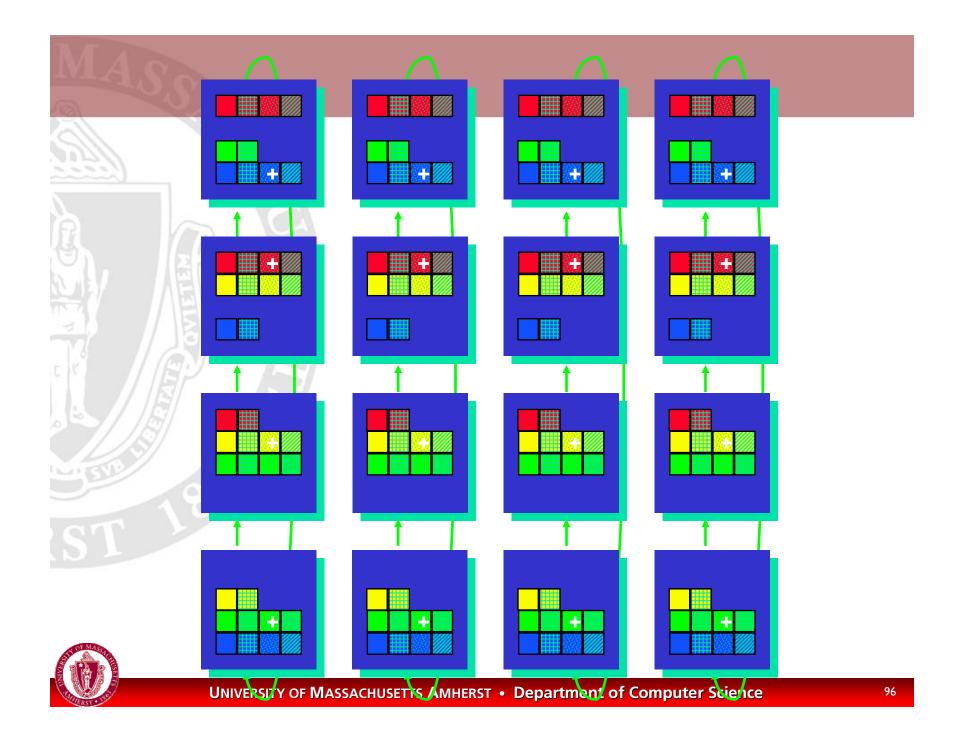


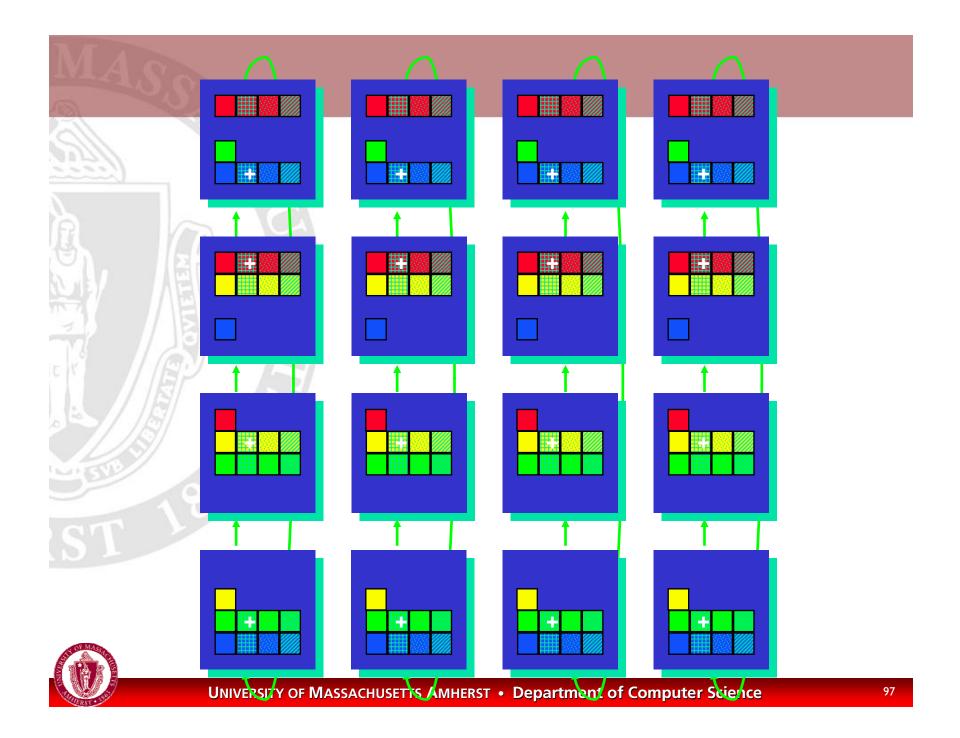


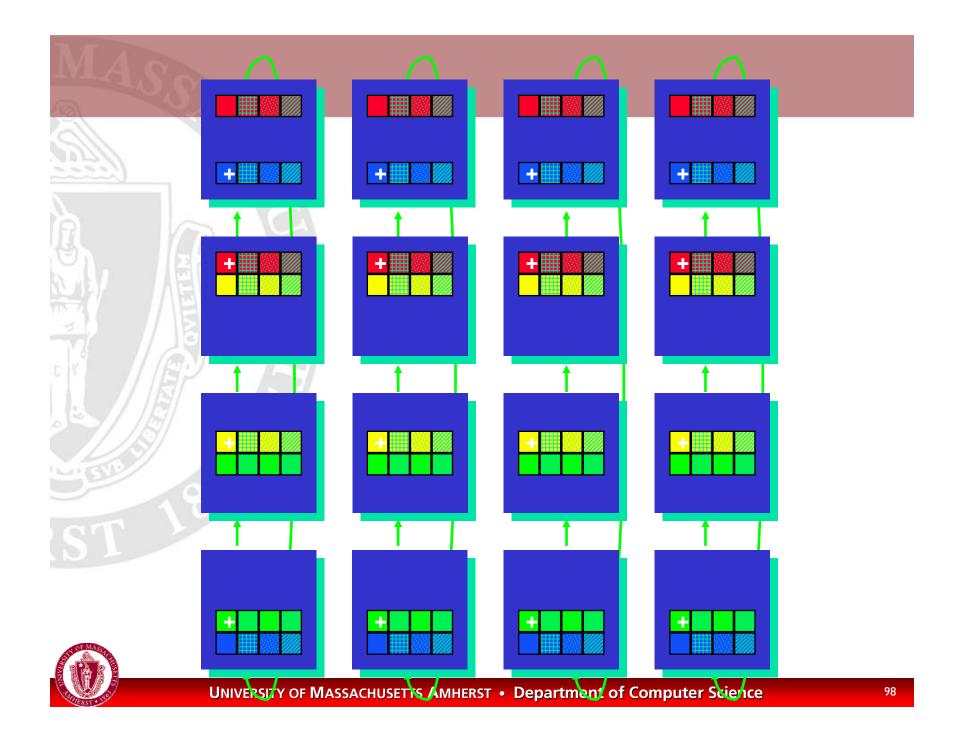


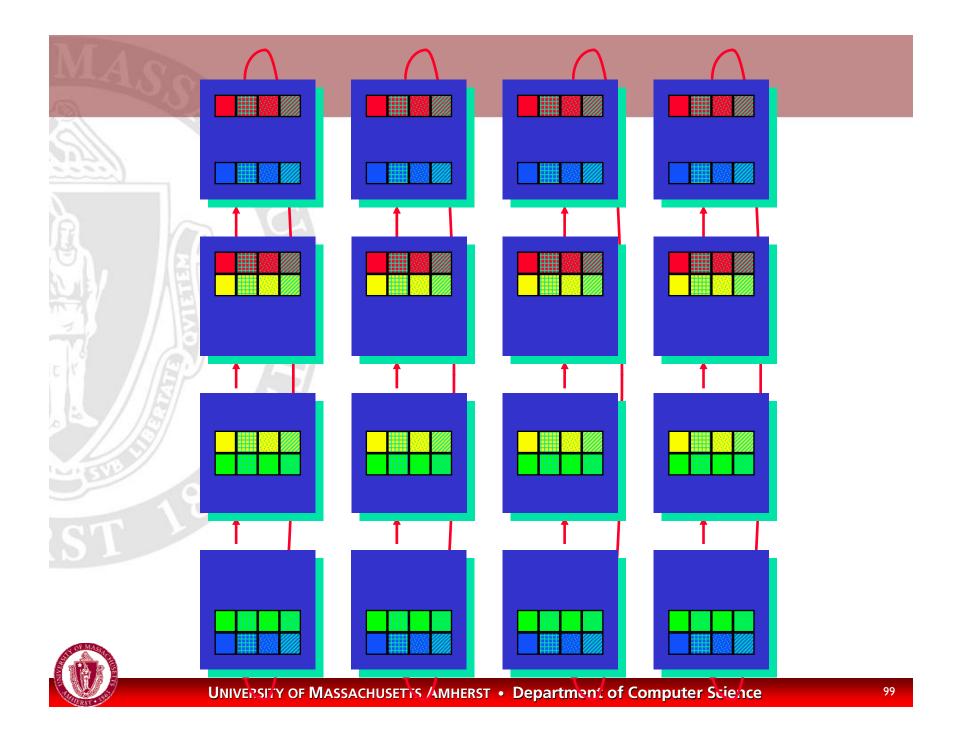


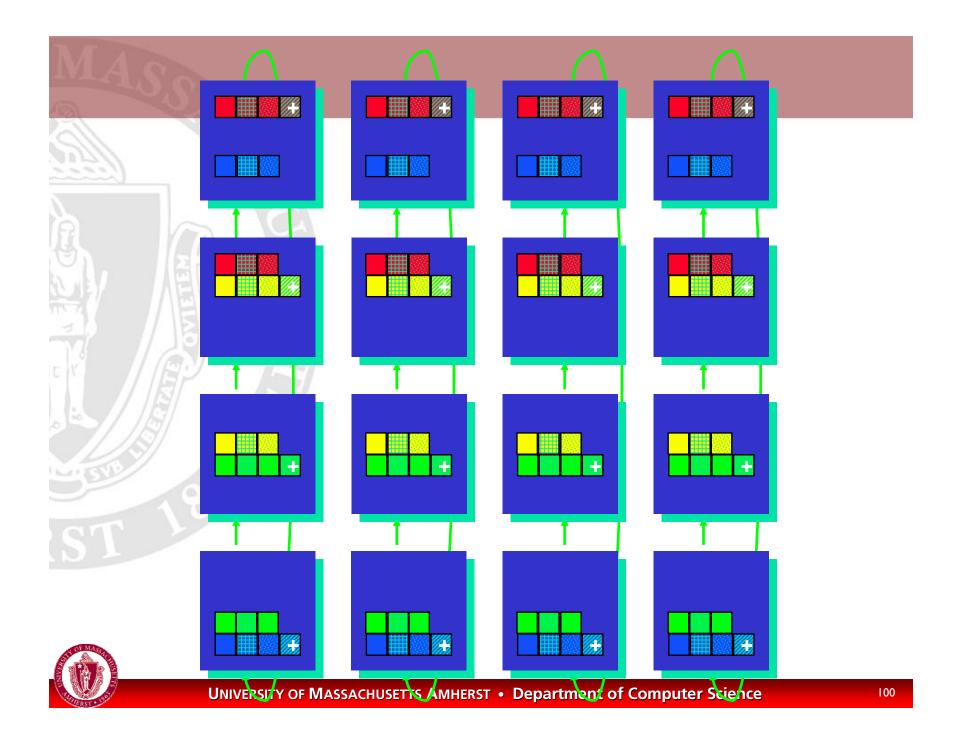


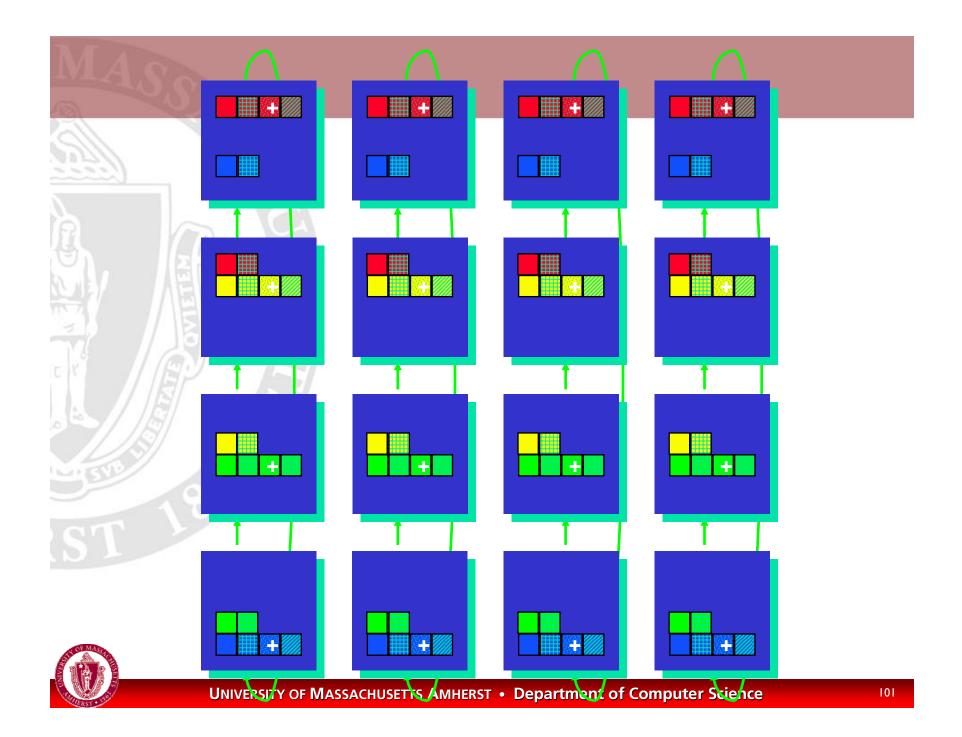


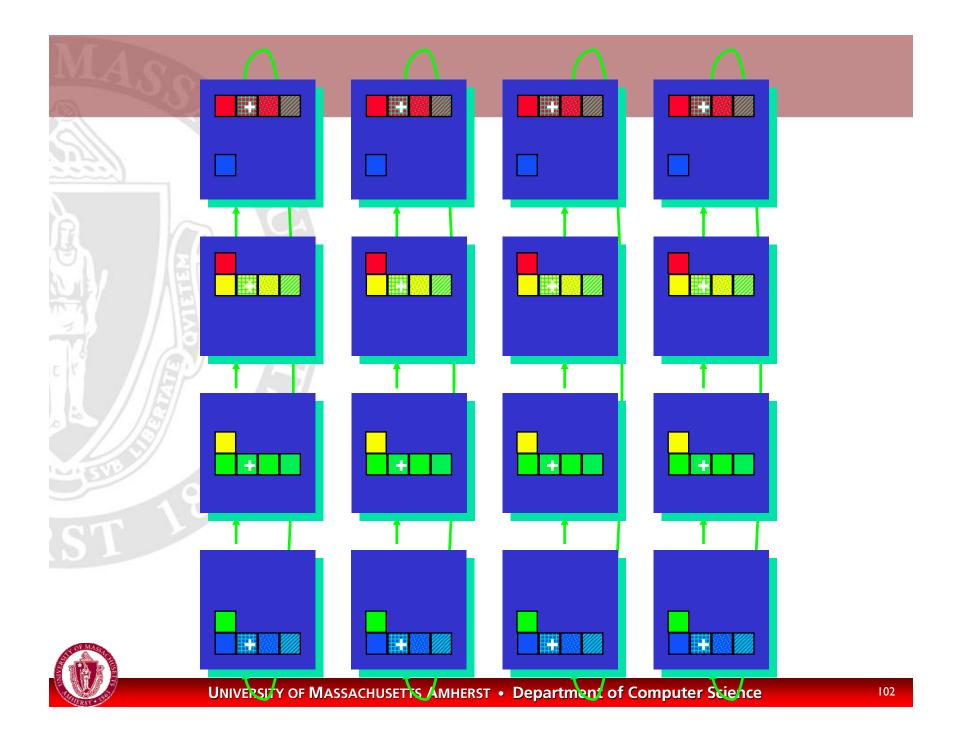


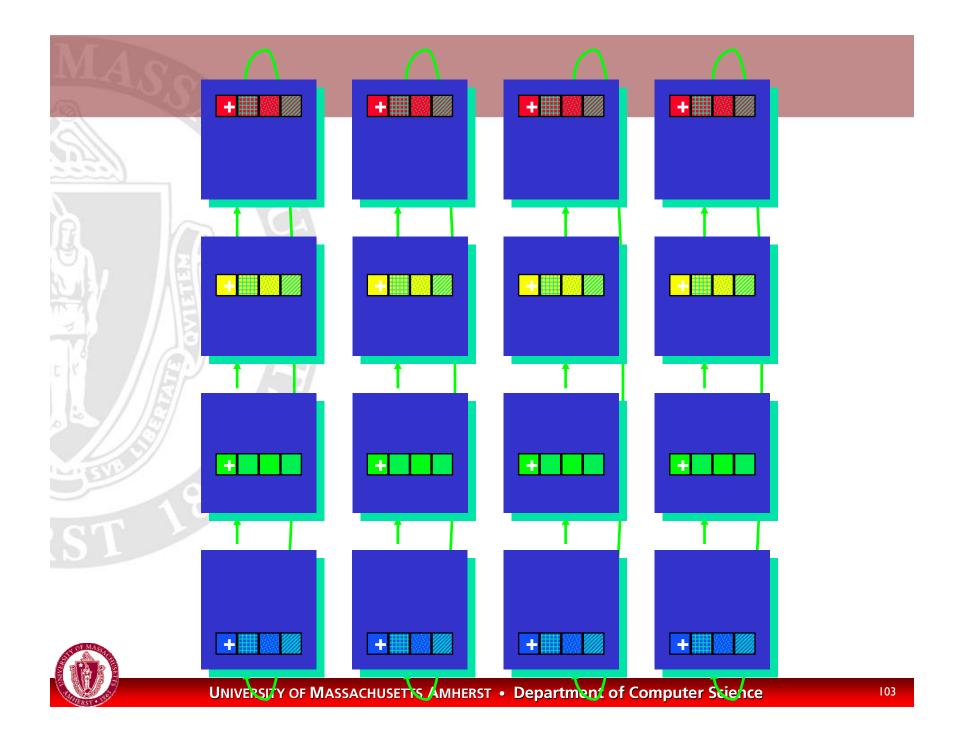


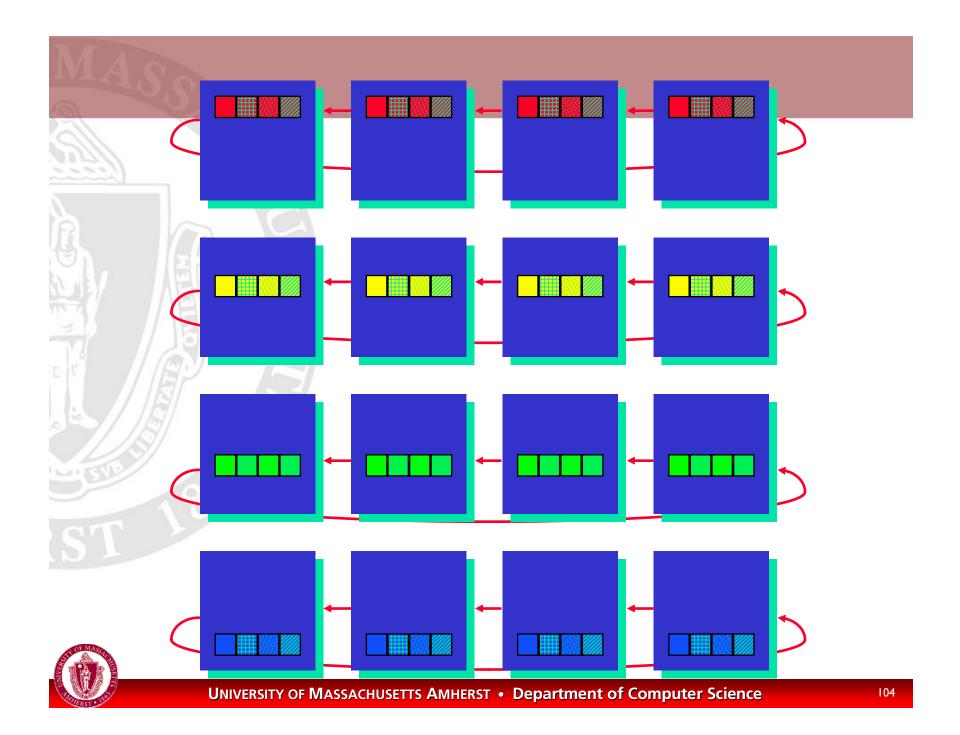


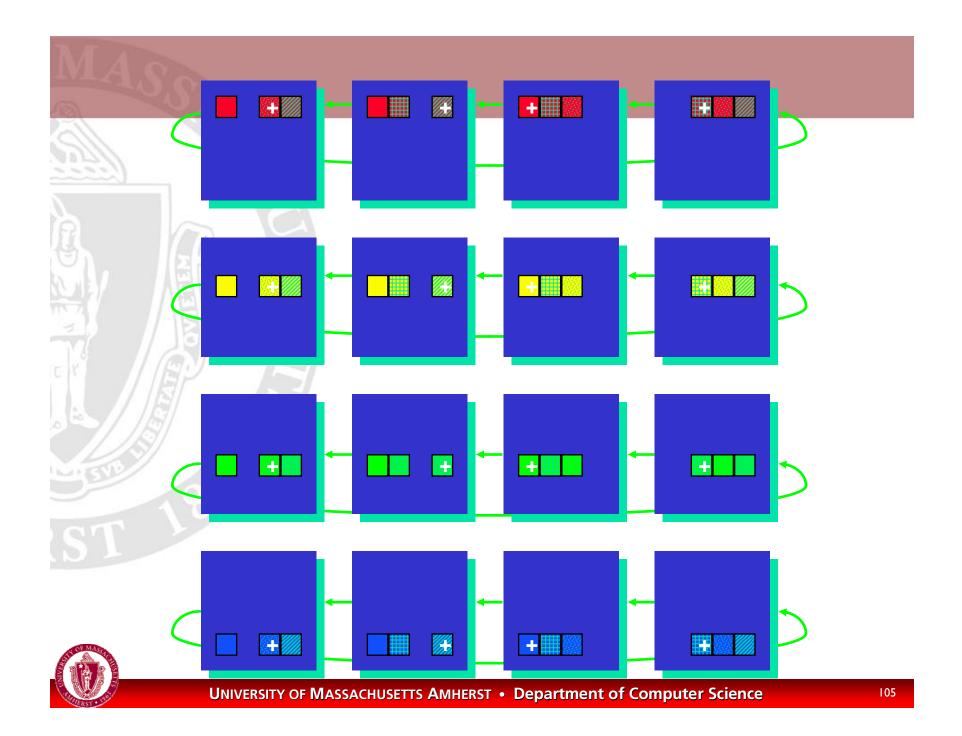




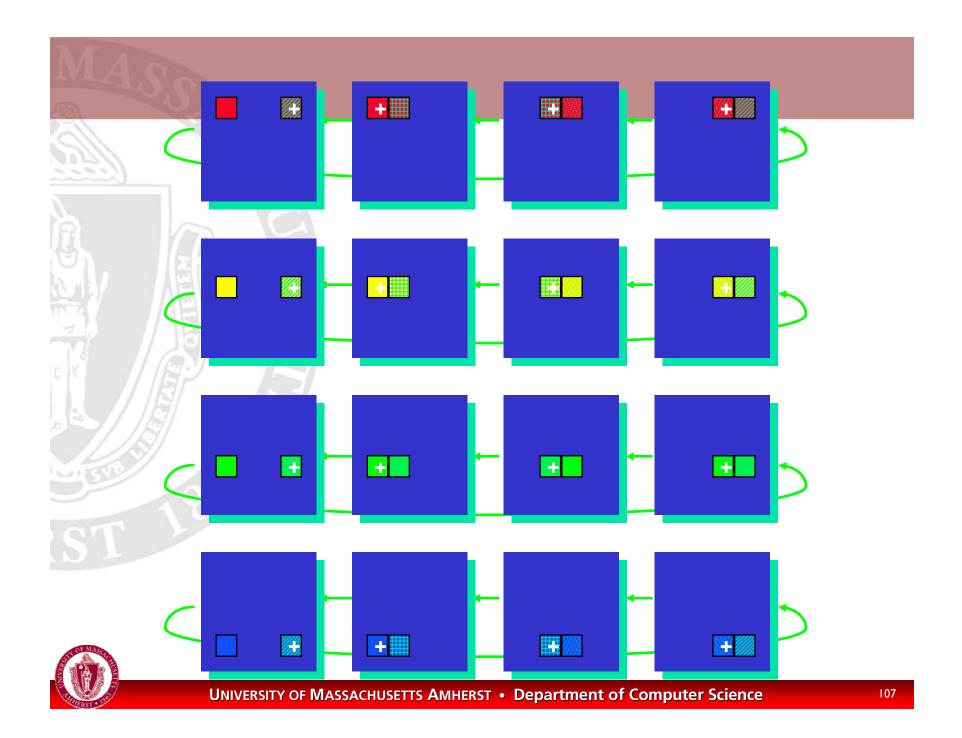


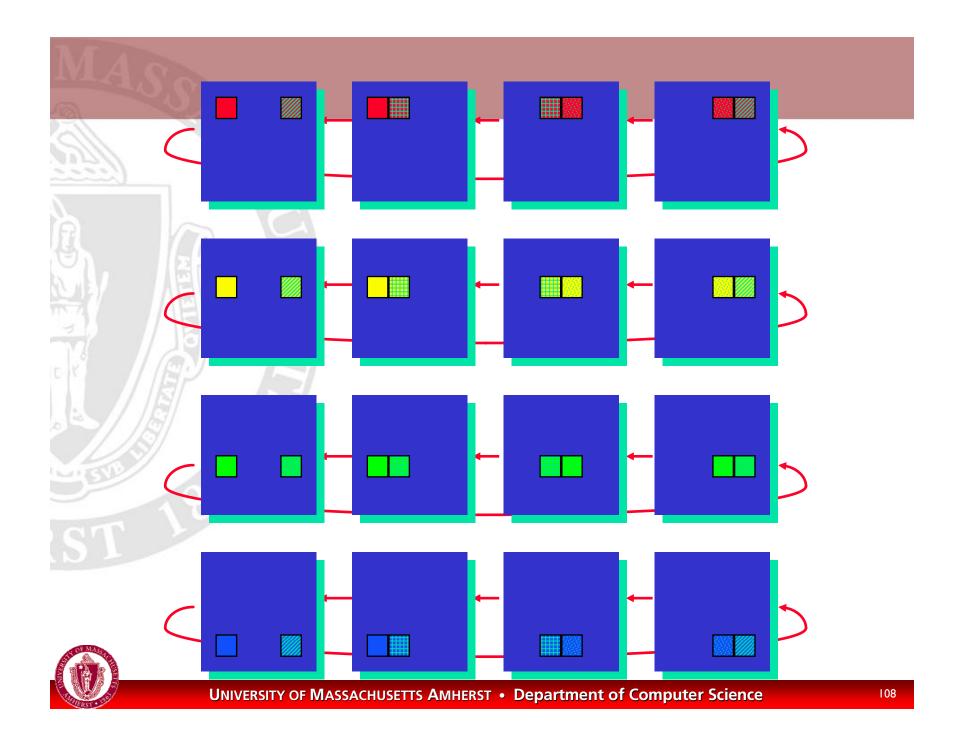


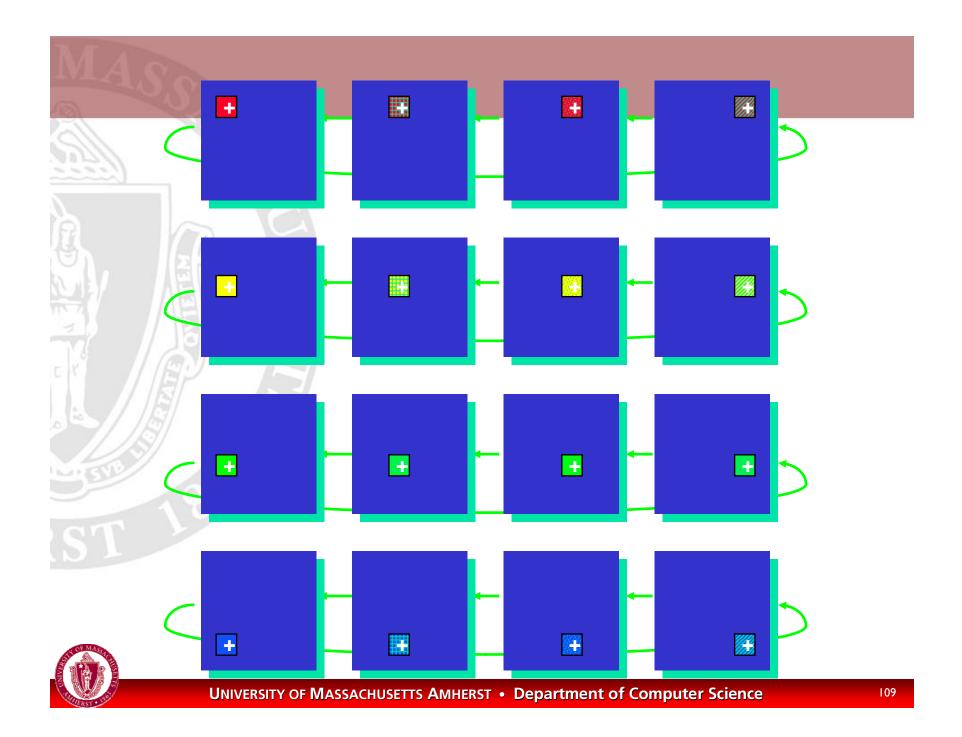


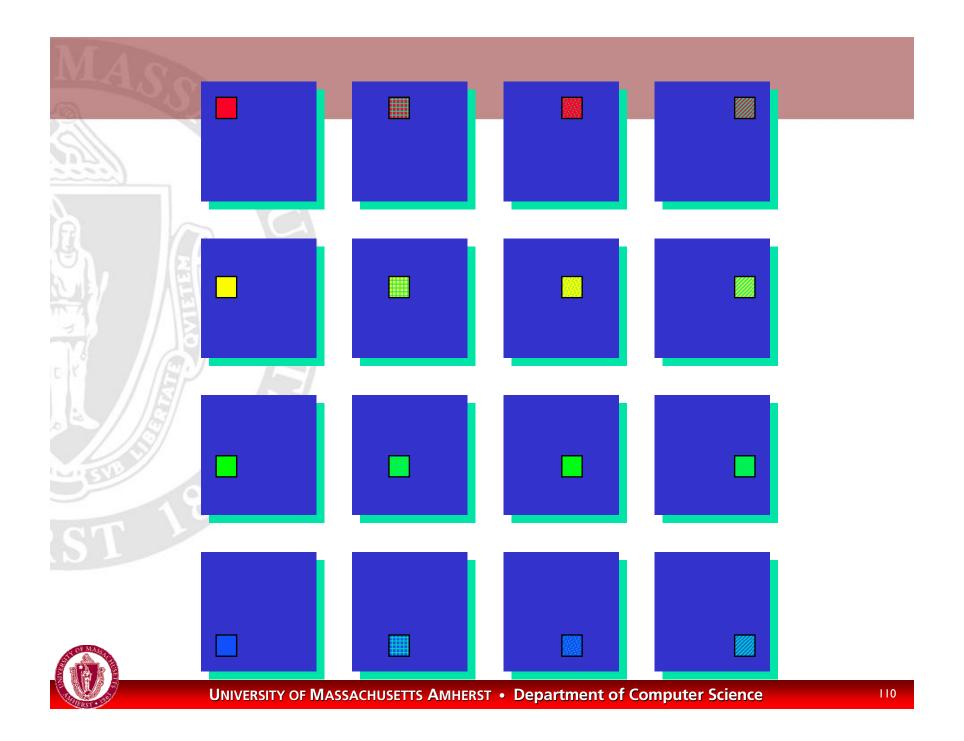












### Summarv

#### From root task to all others

- MPI\_Bcast
  - (buf, cnt, type, root, comm)
- MPI\_Scatter (sendbuf, cnt, type, recvbuf, recvcnt, type, root, comm)

#### From all tasks to root

 MPI\_Gather (sendbuf, cnt, type, recvbuf, recvcnt, type, root, comm)

#### From all to all

- MPI\_Alltoall (sendbuf, cnt, type, recvbuf, recvcnt, type, comm)
- MPI\_Allgather (sendbuf, cnt, type, recvbuf, recvcnt, type, comm)

#### Reductions

- MPI\_Reduce (sendbuf, recvbuf, cnt, type, op, root, comm)
- MPI\_Scan (sendbuf, recvbuf, cnt, type, op, comm)



# **Communicators**

#### Can duplicate:

- MPI\_Comm\_dup (old, &new)
- Can then split communicators:
  - MPI\_Comm\_split (old, color, key, &new)
  - Same color ⇒ same communicator
  - Key determines rank in new

