## Slide 1

# Parallel Implementation Techniques for Embedded DSP Systems
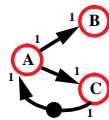
**S. Sriram**

**Prof. Edward A. Lee**

## Slide 2

# Parallelism in Embedded DSP Systems

- **Parallelism : concurrency at the system level**
  - target a system consisting of a mix of dedicated parts such as FFT chips and programmable DSPs

- **Embedded : low cost, dedicated multiprocessor**
  - Examples — multimedia: set-top boxes, multimedia workstations, communications: digital cell phones

- **Motivation**
  - high throughput applications demand processing power
  - use of commodity programmable parts: attractive alternative to ASICs
  - advantages of software solutions
  - silicon technology: multi-DSP chips available from number of companies

## Slide 3

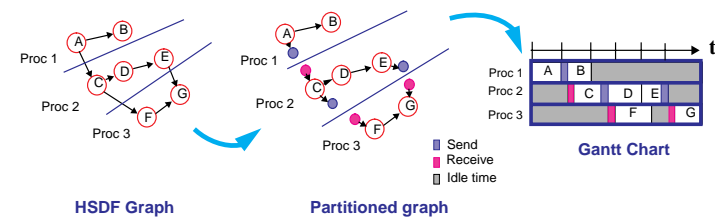# Parallel Implementation of DSP Algorithms Static

- **Issues: scheduling, interprocessor communication (IPC) & synchronization overhead, hard real-time requirement**
  - computation on unbounded data streams

- **Our strategy**
  - Restricted application domain: **Synchronous Dataflow and extensions**



  - Well-defined methodology:
    - Compilation from dataflow graphs
    - Extensive use of compile-time scheduling techniques
  - Given this methodology optimize hardware architecture and parallel implementation:
    - Reduce IPC overhead: Ordered Transactions scheme
    - Reduce synchronization overhead
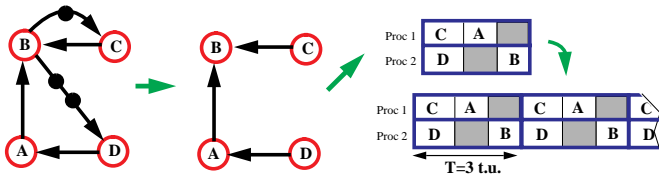
## Slide 4

# Scheduling

- **Scheduling homogeneous SDF graphs**
  - Assigning actors to processor: **Assignment**
  - Determining the order of execution of actors on a processor: **Ordering**
  - Determining when an actor actually fires: **Firing times**
- **Dynamic (run time) vs. Static (compile time) strategies**
- **Use execution time estimates**
  - **Fully Static** : all three scheduling steps performed at compile time, assuming execution times estimates are precise



HSDF Graph          Partitioned graph
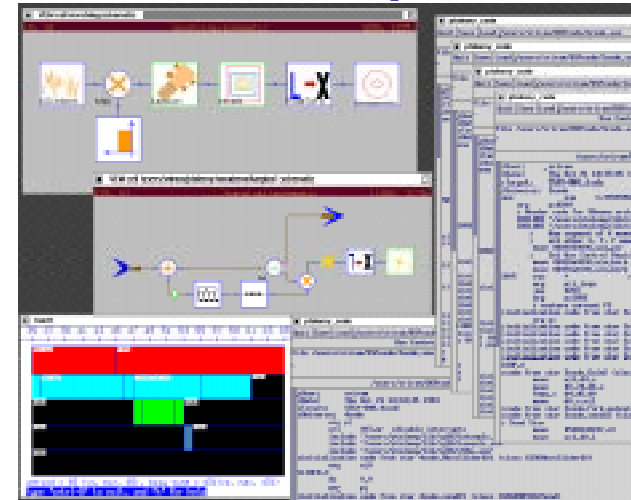
## Blocked Multiprocessor Schedules

- **Homogeneous SDF graph converted to an Acyclic Precedence Graph (APG) by removing edges with delays**
  - Intra-iteration precedences ignored during scheduling



- **Minimize $T$: classical MP scheduling from an APG**
  - Optimal scheduling under resource constraints is intractable (NP-Hard)
  - Several heuristics exist: list scheduling [Hu 61], [Sih 92], [Sarkar 89], ...
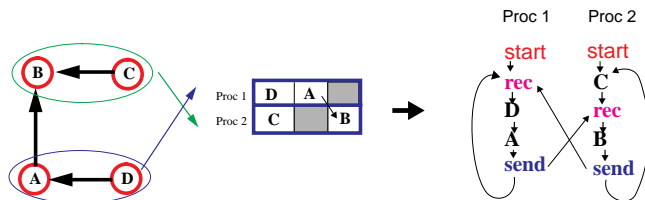  - Unfolding (increased blocking factor) and retiming transformations

## Implementation in Ptolemy

**Block Diagram ➜ SDF Graph ➜ Homogeneous APG ➜ Parallel Schedule ➜ Multiprocessor code**
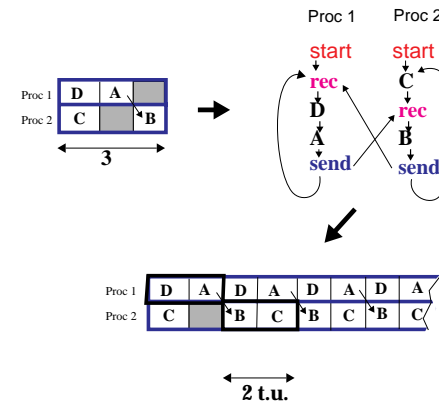
## Self-timed Scheduling

- **Fully-static schedule assumes knowledge of exact actor execution times - not always practical:**
  - compilation from high-level language, error handling, unpredictable execution times due to instruction-level parallelism

- **Model followed in Ptolemy**
  - reasonably good estimates of execution times known at compile time
  - construct fully-static schedule, ignore exact timing information
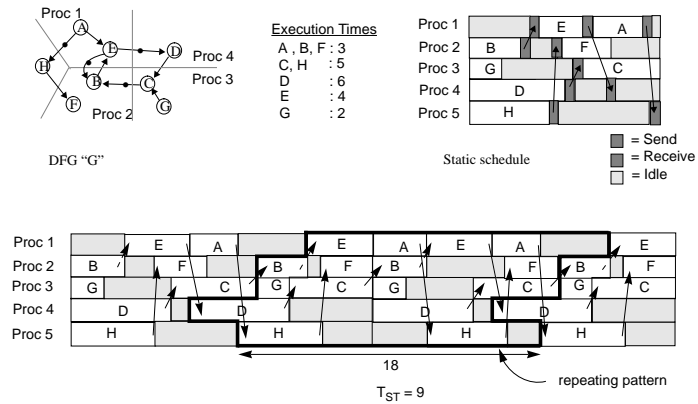


- **Larger run time overhead compared to fully-static sched.**

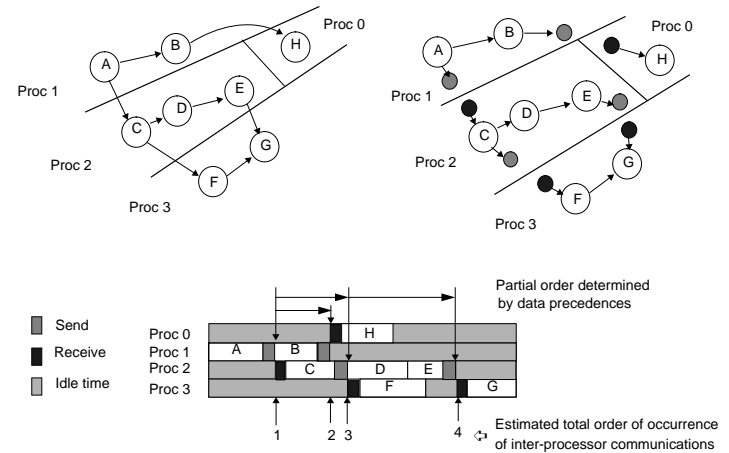## Communication Pattern

- **Attempt to predict run time inter-processor communication pattern and use this information to optimize parallel implementation**
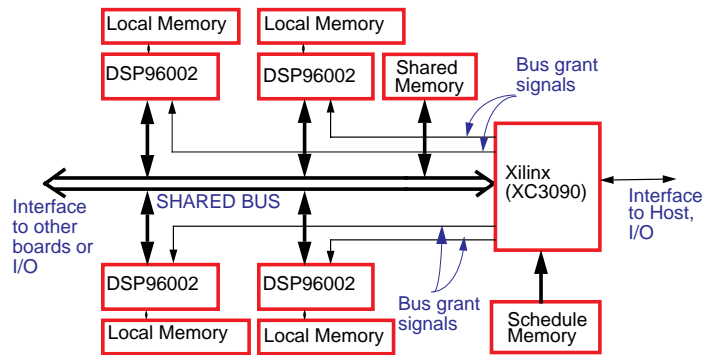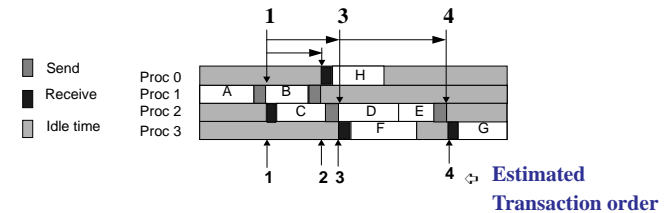
## Communication Pattern

Proc 1

Execution Times
A , B, F : 3
C, H : 5
D : 6
E : 4
G : 2

DFG "G"

Static schedule

Proc 1
Proc 2
Proc 3
Proc 4
Proc 5

E | A
B | F
G | C
D
H

■ = Send
■ = Receive
□ = Idle

Proc 1
Proc 2
Proc 3
Proc 4
Proc 5

18

$T_{ST} = 9$

repeating pattern

---

## Ordered Transactions

Proc 0
Proc 1
Proc 2
Proc 3

Partial order determined
by data precedences

■ Send
■ Receive
□ Idle time

Proc 0
Proc 1
Proc 2
Proc 3

A | B | H
C | D | E
F | G

1    2 3              4   Estimated total order of occurrence
                          of inter-processor communications

---

## Ordered Memory Access Architecture

Local Memory     Local Memory

DSP96002     DSP96002     Shared Memory

Bus grant signals

SHARED BUS

Interface to other boards or I/O

Xilinx (XC3090)

Interface to Host, I/O

DSP96002     DSP96002

Bus grant signals

Local Memory     Local Memory

Schedule Memory

- **Low overhead IPC (3 instruction cycles)**
- **No need for explicit synchronization**
- **Performance degrades if execution times vary at run time (or if compile estimates are bad), computations are correct**

---

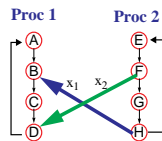## Choosing a Transaction Order

- **Transaction order imposes run time constraints absent in the unconstrained self-timed schedule**
  - **Partial order due to precedence constraints:** $2 \gg 1, 4 \gg 3 \gg 1$
  - **Any total order is a valid transaction order**

1          3          4

■ Send
■ Receive
□ Idle time

Proc 0
Proc 1
Proc 2
Proc 3

A | B | H
C | D | E
F | G

1      2 3        4   **Estimated Transaction order**

- **Naive transaction order derived from one block of the schedule: not always the best choice**
  - **Can efficiently determine a transaction order that is "optimal"**

## Minimizing Synchronization in self-timed schedules

- **Self-timed scheduling: each inter-processor communication point is also a synchronization point**
  - sender needs to check for buffer overflow
  - receiver needs to check for buffer empty

- **Compile time analysis of schedule can reduce this overhead**
  - **Sender synchronizations are eliminated by sizing buffers appropriately**
  - **Remove redundant receiver synchronizations: synch x2 is redundant**



  - **Perform transformations on the schedule: introduce new synchronization points**

## Conclusions

- **Discussed mechanism for constructing parallel schedules from SDF graphs**

- **Discussed how compile time scheduling can be effectively employed for SDF applications**

- **Discussed parallel code generation methodology in Ptolemy**

- **Presented the ordered transactions approach: hardware architecture optimized for the self-timed strategy employed in Ptolemy**

- **Described minimization of synchronization costs by means of compile-time analysis**