# FastForward for Efficient Pipeline Parallelism

John Giacomoni, Tipp Moseley,
and Manish Vachharajani
University of Colorado

## Parallelizing Sequential Tasks:

Need to efficiently support the parallelization of tasks with partially or totally ordered data. Pipeline-parallel organizations avoid contention and stalls caused by intra-thread synchronization other organizations.



## Hardware Performance Issues

### Cache Thrashing (Communication Channel)



UMA        NUMA        Careful resource scheduling may be required for performance.

(AMD)

### Weak Memory Models



Weak memory models can exhibit unexpected behavior;
In this example, a write passing a write may result in a read of a stale data value.

## The FastForward Solution

```
1: enqueue_lamport(...) {            1: enqueue_fastforward(...) {
2:   if(NEXT(head) == tail) {        2:   if(NULL != buf[head]) {
3:     // Handle full queue.         3:     // Handle full queue.
4:   }                               4:   }
5:   buf[head] = data;              5:   buf[head] = data;
6:   head = NEXT(head);             6:   head = NEXT(head);
7: }                                7: }

9: dequeue_lamport(...) {            9: dequeue_fastforward(...) {
10:                                 10:   data = buf[tail];
11:   if (head == tail) {           11:   if (NULL == data) {
12:     // Handle empty queue.       12:     // Handle empty queue.
13:   }                             13:   }
14:   data = buf[tail]              14:
15:   tail = NEXT(tail)             15:   tail = NEXT(tail)
16: }                              16: }
```

Decoupling at the cache coherence layer can eliminate cache thrashing, hide non-uniform memory access issues, and support weak memory models.

## Proof Sketch

The references [1] prove that "in the program order of the consumer, the consumer dequeues values in the same order that they were enqueued in the producer's program order," for strong to weakly ordered consistency models, showing that FastForward works even on relaxed memory models.

[1] J. Giacomoni, T. Moseley, and M. Vachharajani. FastForward for efficient pipeline parallelism. Technical Report CU-CS-1028-07, Univerity of Colorado at Boulder, 2007.

## Performance Comparison



2 Threads        2 Threads        3 Threads
On-die           Cross-die        Cross-die
Dual-processor Dual-core 2.0 GHz AMD Opteron

## Conclusions

1) Decoupling communicating threads at the cache layer on ccNUMA machines may yield significant performance improvements.

2) FastForward provides an efficient point-to-point communication primitive ideally suited to pipeline-parallel applications.
   a) Consistent performance.
   b) Software only solution.
   c) Correct under strong to weakly ordered consistency models.
   d) May provide performance improvements to other streaming parallel organizations.