

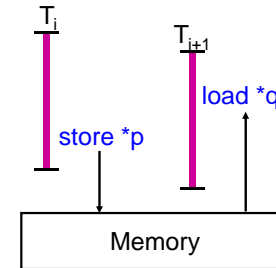
Lecture 28(b)

Compiler Optimizations for Thread-Level Speculation

From the paper: *"Compiler Optimization of Scalar Value Communication Between Speculative Threads"*, by Antonia Zhai, Christopher B. Colohan, J. Gregory Steffan and Todd C. Mowry. ASPLOS, 2002.

Carnegie Mellon

Speculation



☞ good when $p \neq q$

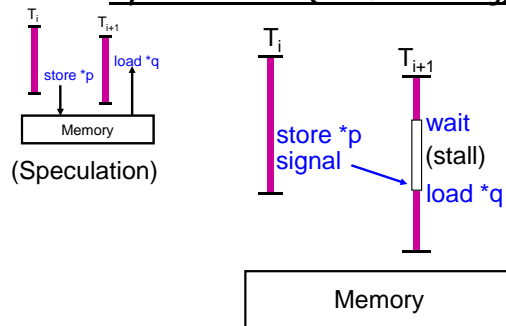
15-745: Optimizing TLS

2

Carnegie Mellon

Todd C. Mowry

Synchronization (and Forwarding)



☞ good when $p == q$

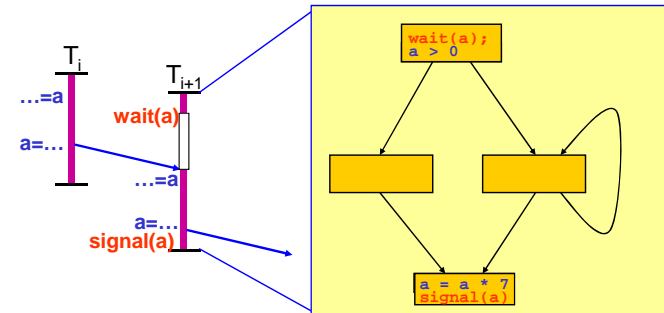
15-745: Optimizing TLS

3

Carnegie Mellon

Todd C. Mowry

Synchronizing Scalars

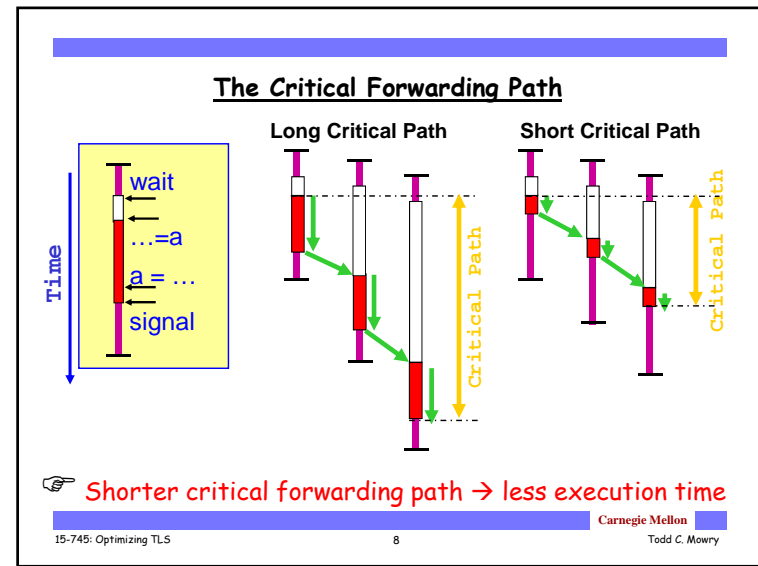
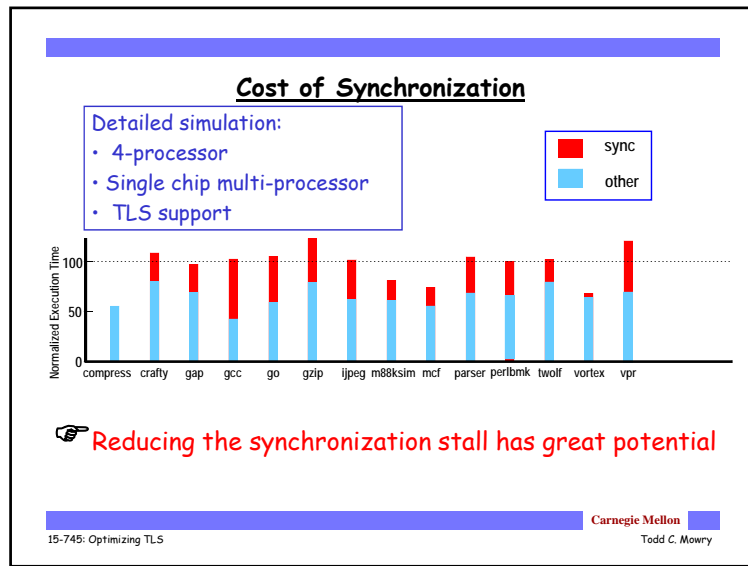
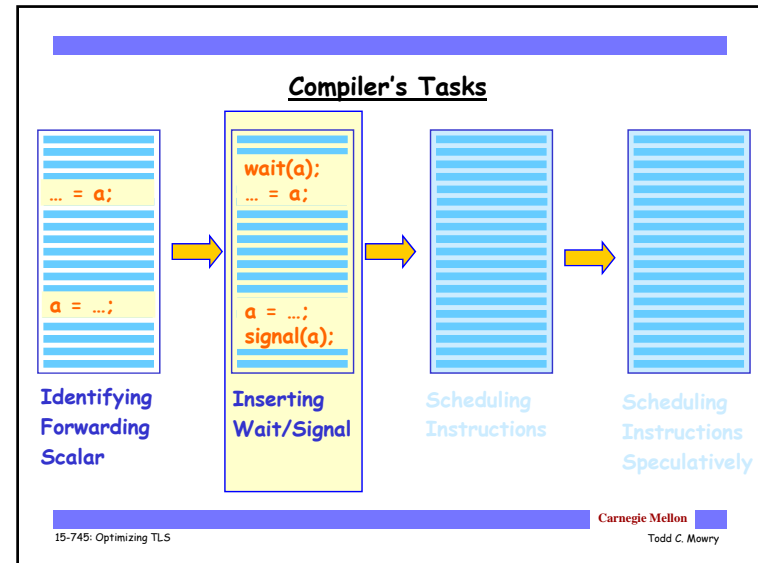
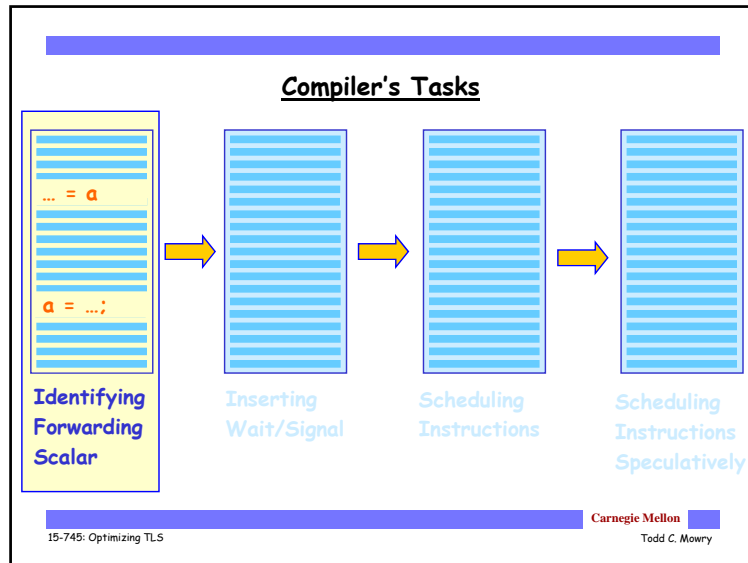


15-745: Optimizing TLS

4

Carnegie Mellon

Todd C. Mowry



Reducing the Critical Forwarding Path

Instruction scheduling can reduce critical forwarding path

15-745: Optimizing TLS Carnegie Mellon
9 Todd C. Mowry

Compiler's Tasks

15-745: Optimizing TLS Carnegie Mellon
Todd C. Mowry

Related Work and Contribution

Related work:

- Multiscalar instruction scheduling [Vijaykumar, Thesis '98]
 - Moving instructions backward one basic block at a time
 - Evaluated under the context of Multiscalar

Our contributions:

- A robust instruction scheduling algorithm
 - Deals with larger threads
 - Handles complex control flow
- Control and data dependence speculation
 - Extends our algorithm to accommodate speculative scheduling
 - Evaluates with detailed simulation
- Comparison with hardware techniques that reduce critical path

15-745: Optimizing TLS Carnegie Mellon
11 Todd C. Mowry

Scheduling Instructions

Dataflow analysis
Handles complex control flow

Define two dataflow analyses

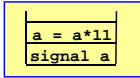
- Stack**
Find the instructions to compute the forwarded value?
- Earliest**
Find the earliest node to compute the forwarded value?

15-745: Optimizing TLS Carnegie Mellon
Todd C. Mowry

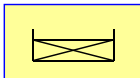
Computation Stack

☞ Stores the instructions to compute a forwarded value

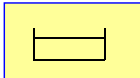
Associating a stack with every node for every forwarded scalar



We know how to compute the forwarded value



We don't know how to compute the forwarded value

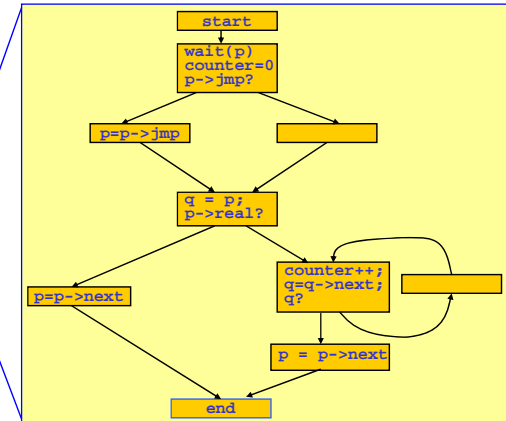


We have not evaluated this node

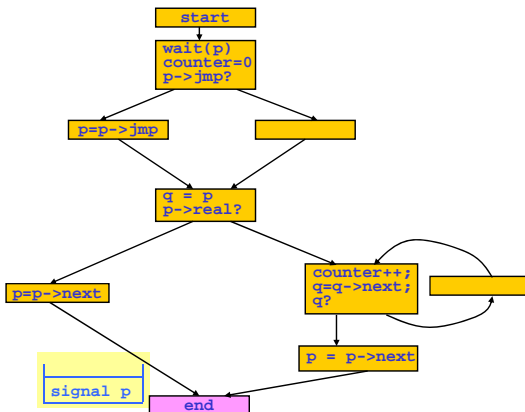
A Simplified Example from GCC

```

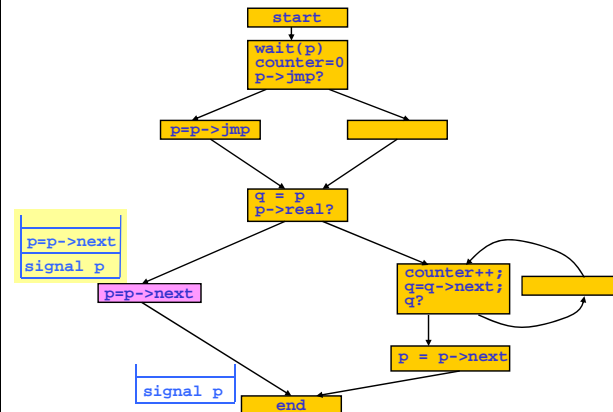
do {
  ...
} while(p);
  
```

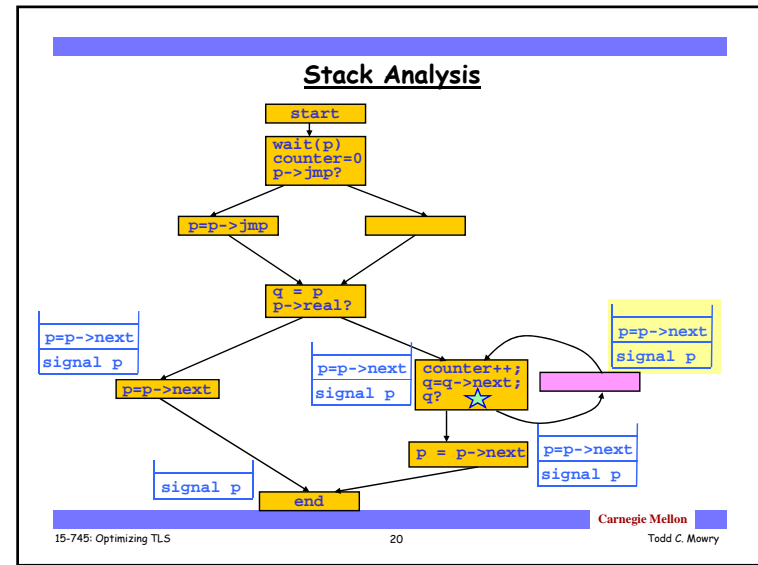
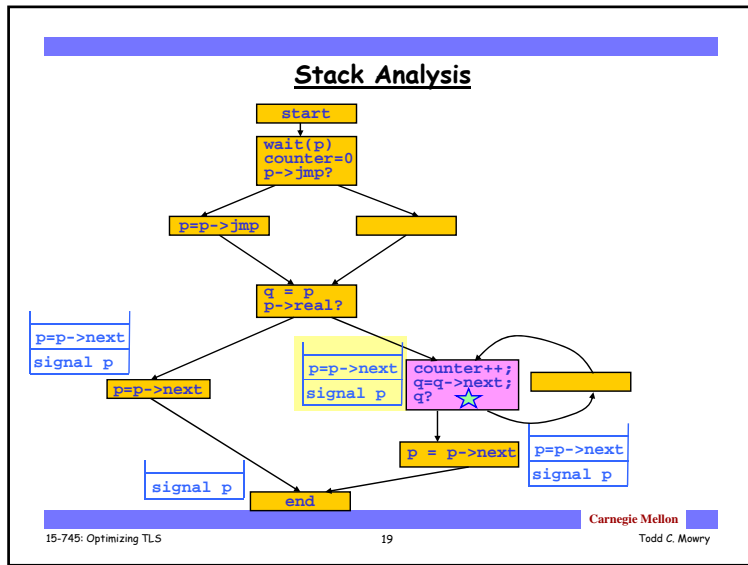
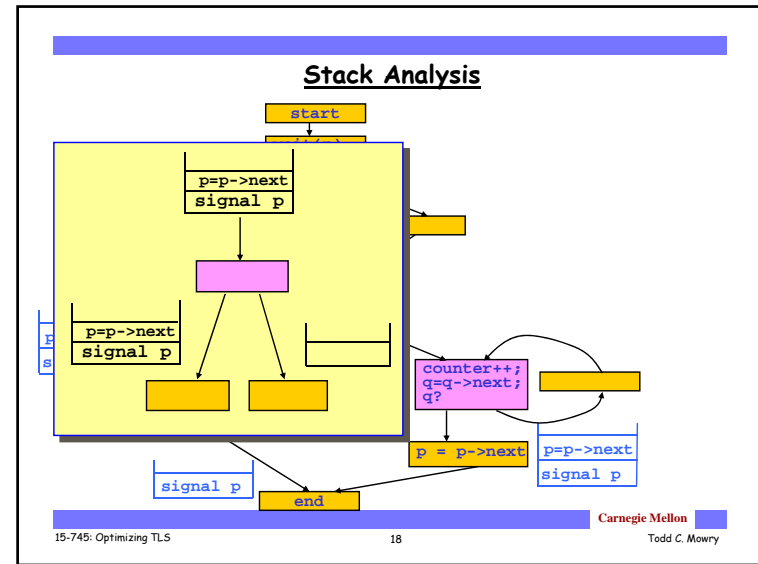
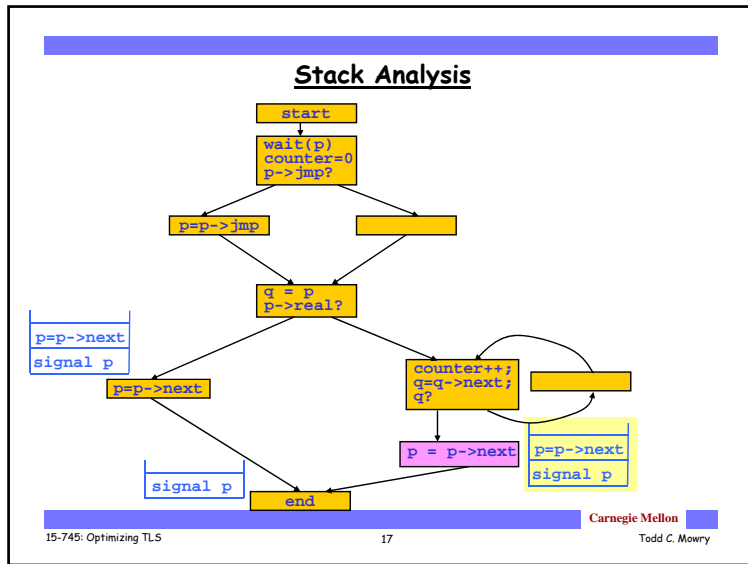


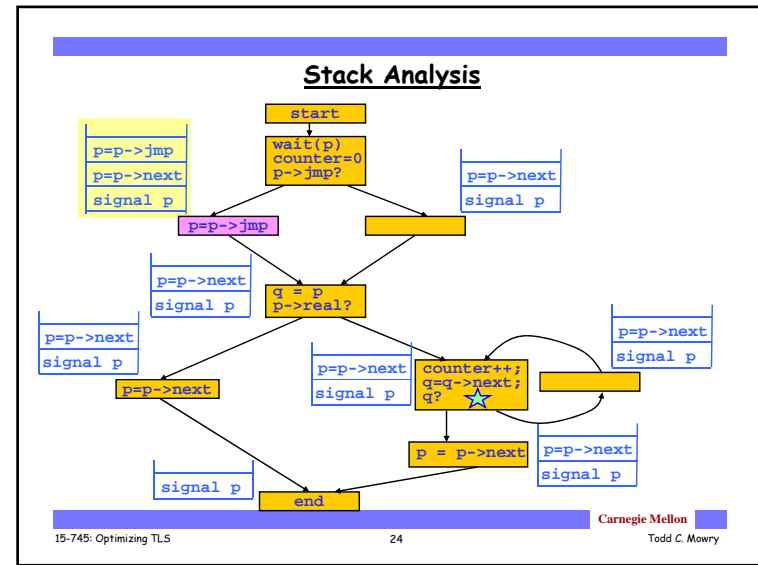
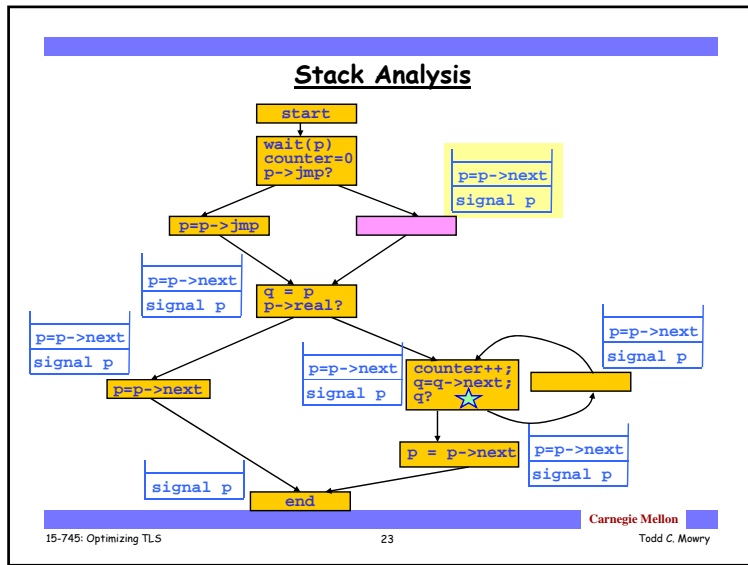
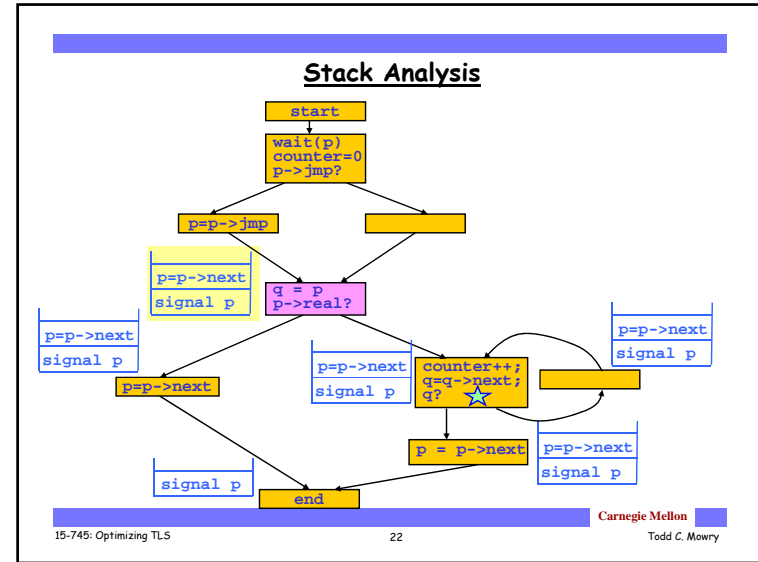
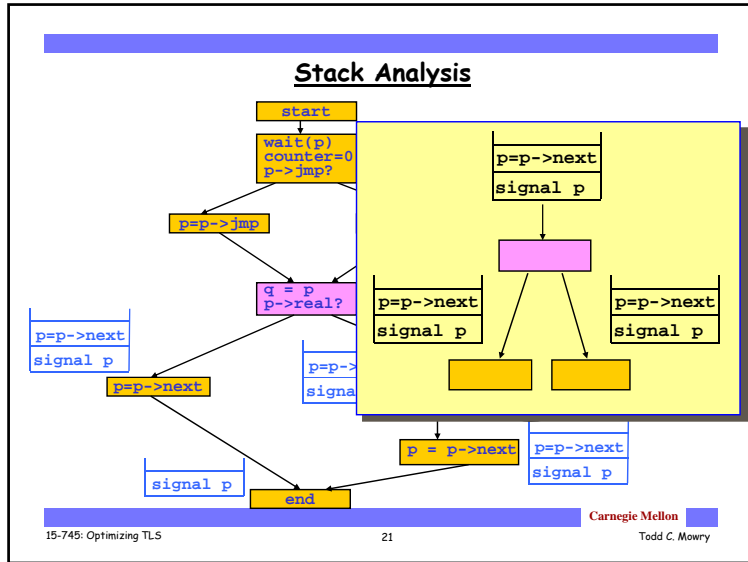
Stack Analysis



Stack Analysis







Scheduling Instructions

Dataflow analysis

Handles complex control flow

Define two dataflow analyses

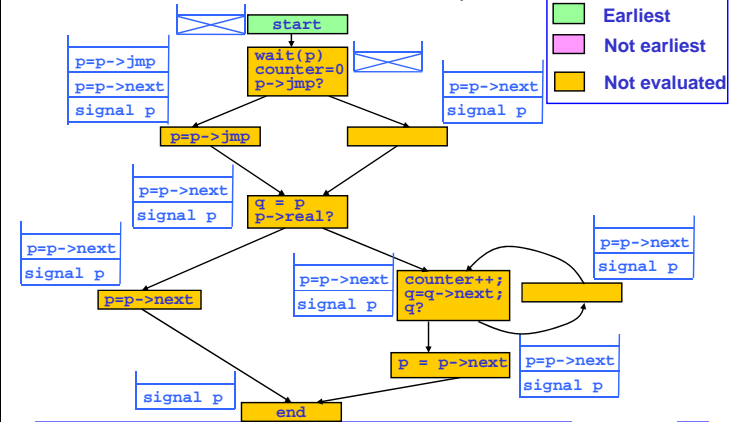
Stack

Find the instructions to compute the forwarded value?

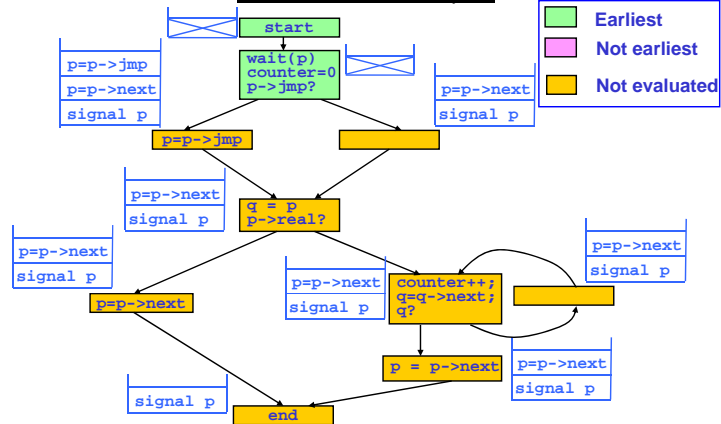
Earliest

Find the earliest node to compute the forwarded value?

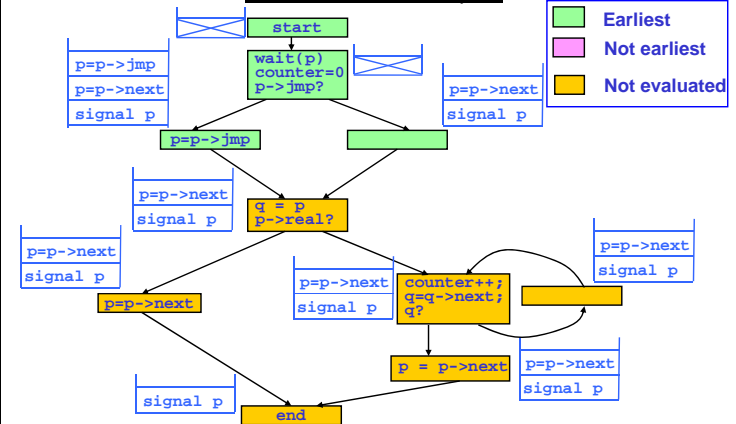
The Earliest Analysis

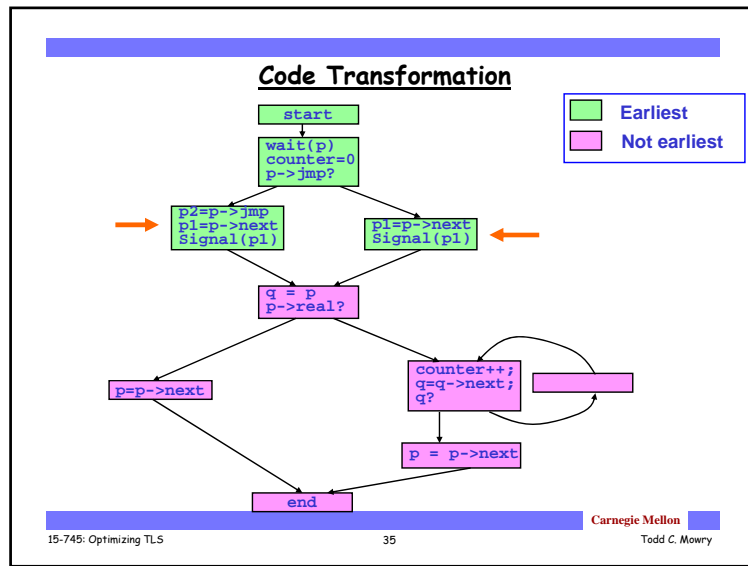
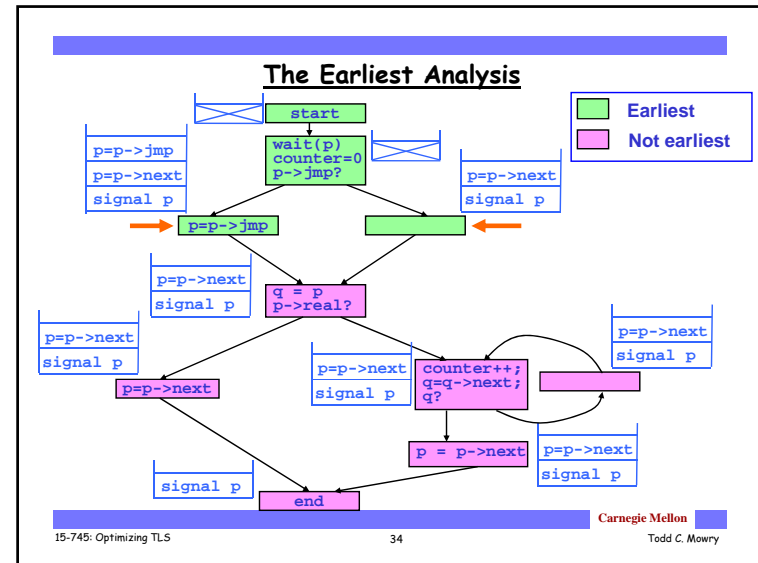
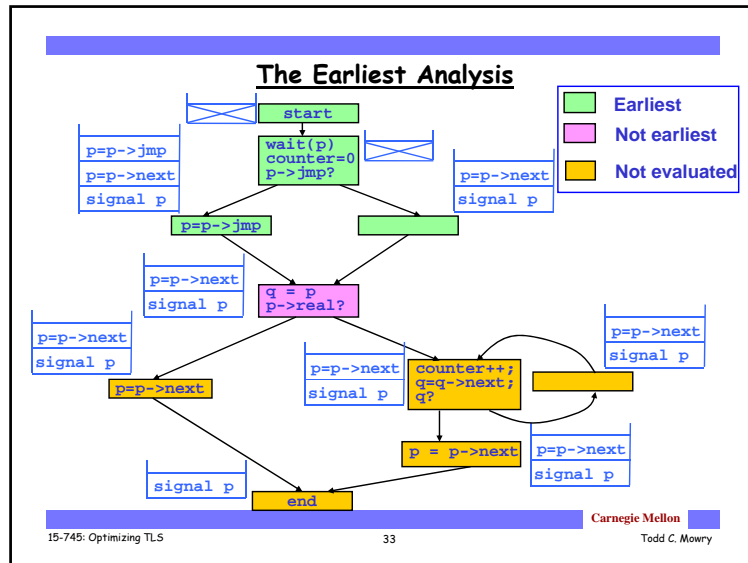


The Earliest Analysis



The Earliest Analysis





Experimental Framework

Benchmarks

- from SPECint95 and SPECint2000, -O3 optimization

Underlying architecture

- 4-processor, single-chip multiprocessor
- speculation supported by coherence

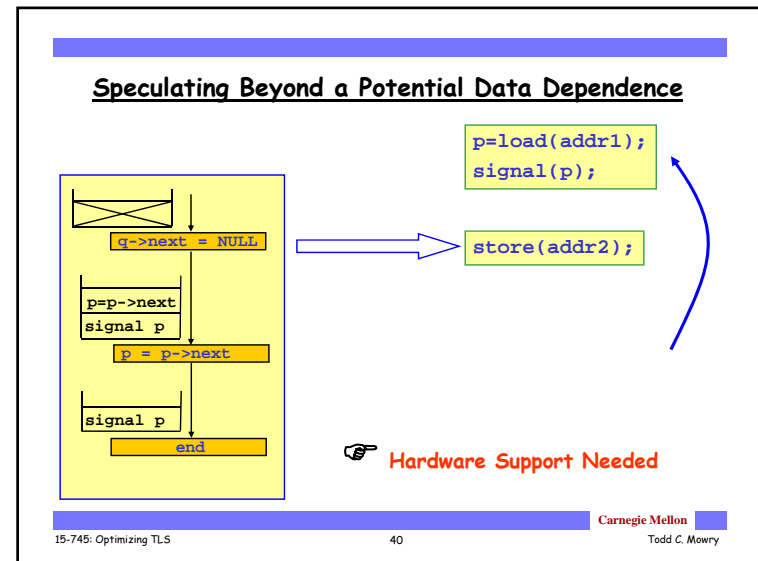
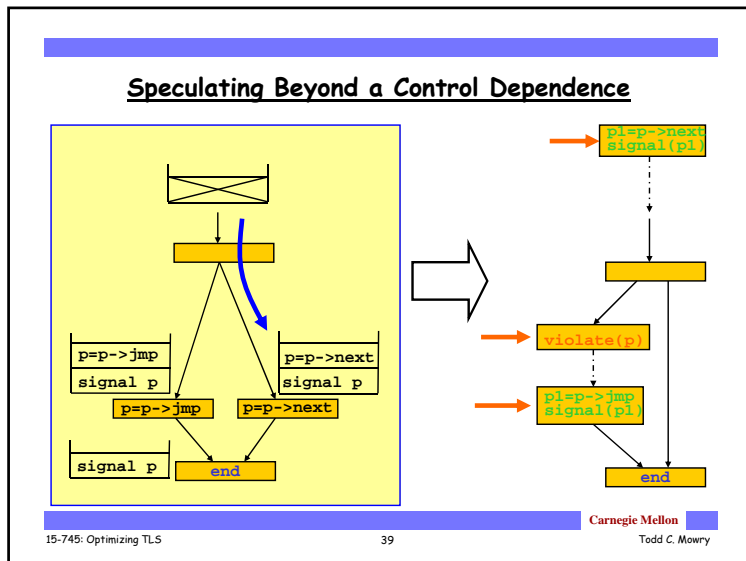
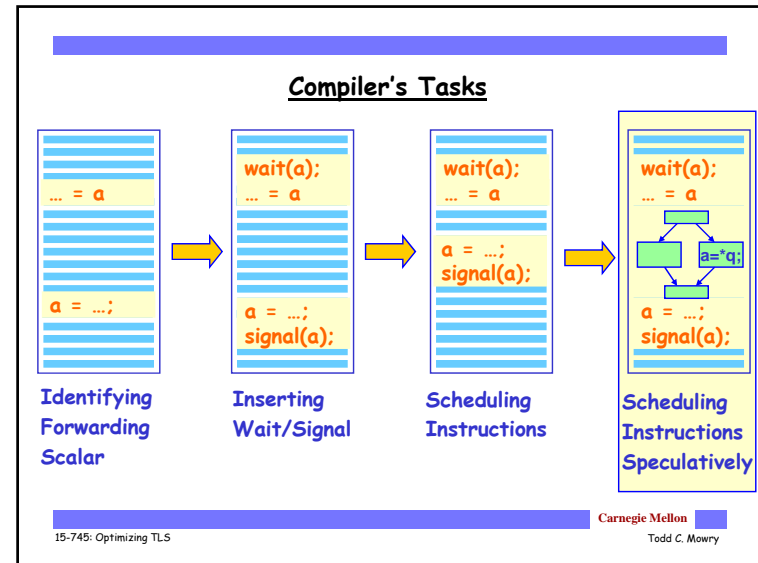
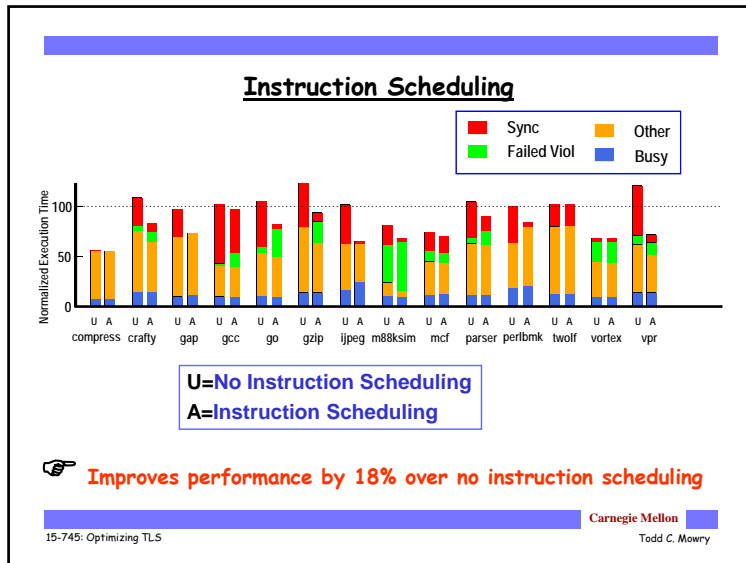
Simulator

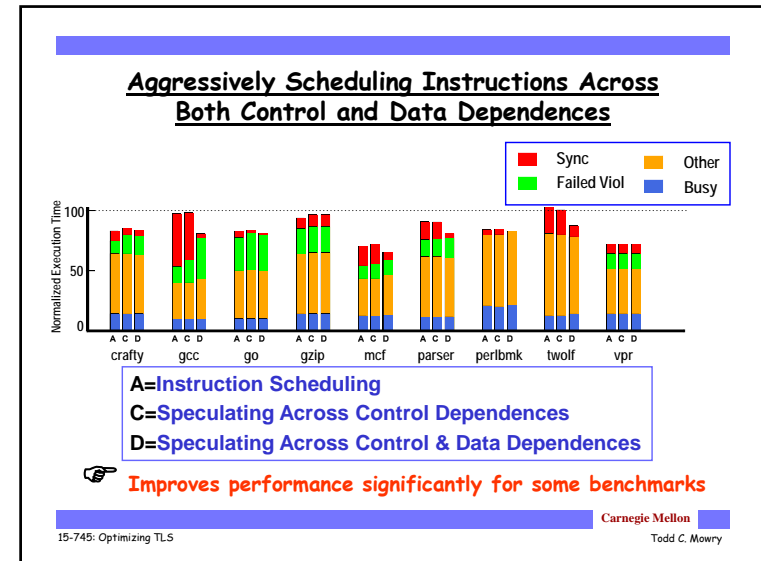
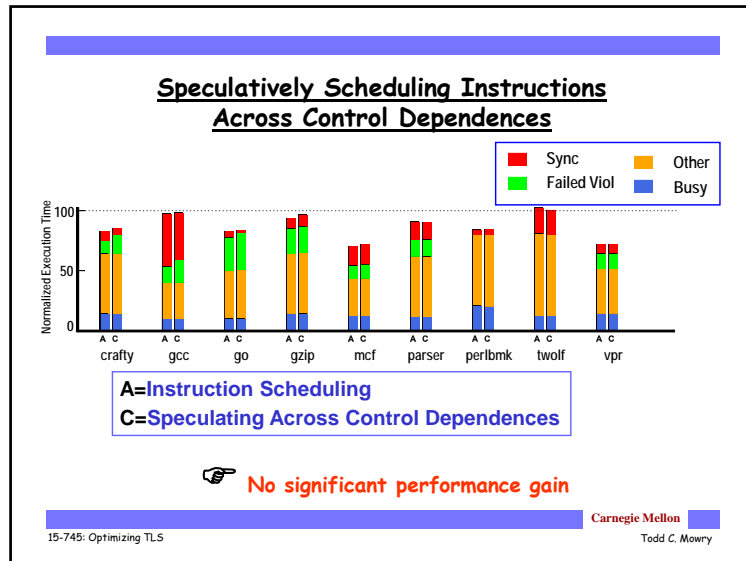
- superscalar, similar to MIPS R10K
- models all bandwidth and contention

The diagram shows a 4-processor, single-chip multiprocessor architecture. It consists of four processors (P) connected to a central crossbar, which is connected to a cache (C). The processors are arranged in a 2x2 grid, and the crossbar is a central box with four ports. The cache is a single box at the bottom.

detailed simulation!

15-745: Optimizing TLS 36 Carnegie Mellon Todd C. Mowry





- ### Hardware Optimization to Reduce Synchronization
- Hardware optimization techniques [Steffan et al, HPCA'02]
- Avoid synchronization:
 - use the value from a hardware value predictor
 - Reduce synchronization stalls:
 - prioritize computation of forwarded value
- Hardware optimization impact [Steffan et al, HPCA'02]
- No compiler optimization: **Effective**
 - With compiler optimization: **Negligible**
- Carnegie Mellon
Todd C. Mowry

- ### Conclusions
- #### Instruction scheduling for reducing synchronization
- Is effective in reducing critical forwarding path
 - Performance improved by 18%
 - Is beneficial to handle complex control flow, such as inner loops
 - Improved GCC by 3%
 - Gives additional benefit with speculative instruction scheduling
 - Our robust instruction scheduling algorithm can be easily extended to accommodate this
 - One biggest benefactor is GCC, performance improved by 18%
 - Reduces the importance of additional hardware optimization
- 👉 **Critical forwarding path can be addressed by the compiler**
- Carnegie Mellon
Todd C. Mowry