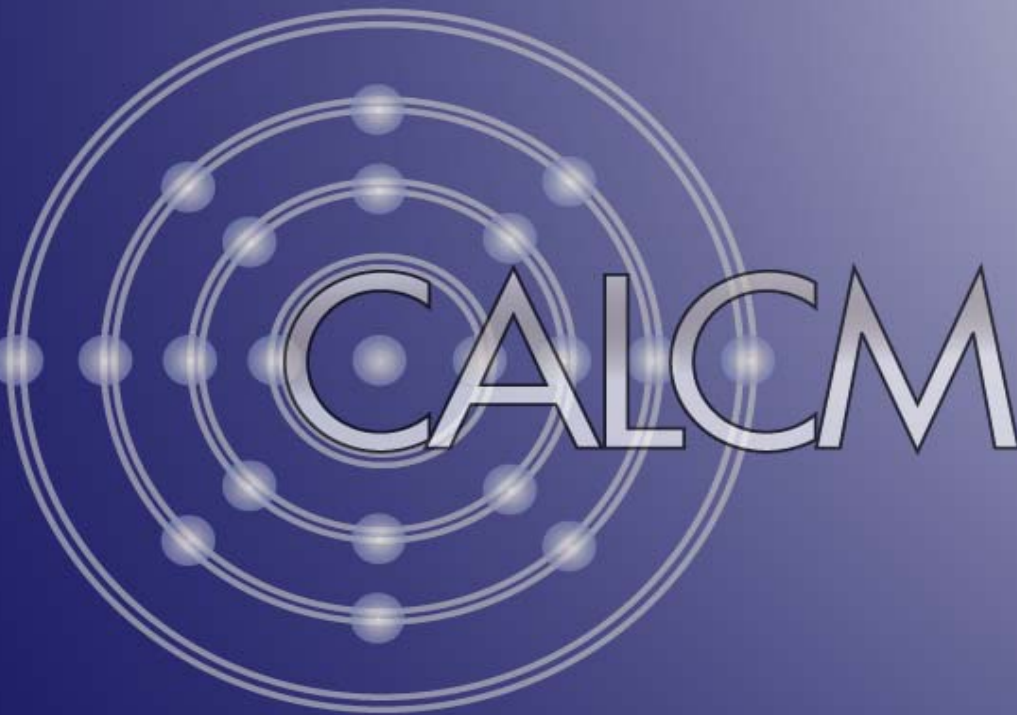


# Execution Slack and Criticality

Nikos Hardavellas



The Computer Architecture Lab at Carnegie Mellon  
<http://www.ece.cmu.edu/CALCM>



# Execution Slack & Criticality: Context

Computer Architecture Lab at Carnegie Mellon

- Out of order processors are highly parallel
  - overlap computation
- Technological constraints
  - wire delay, power, complexity
- Non-uniform designs
  - clustering, multi-frequency FU's, memory hierarchies
- Criticality & slack
  - what is it?



# Execution Slack & Criticality: Why do I care?

Computer Architecture Lab at Carnegie Mellon

- Exploit to guide control policies
  
- Examples here do :
  1. Resource Arbitration
  2. Misspeculation Reduction
  3. Power Management
  4. Cache Management



1. **Focusing processor policies via critical-path prediction**

Framework to identify instruction criticality and design to increase instruction scheduling locality in a cluster

2. **Slack: maximizing performance under technological constraints**

Framework to identify variability in instruction execution criticality and application in slowing down execution to save power

3. **Locality vs. Criticality**

Is sacrificing locality for criticality a good idea?

4. **Non-vital loads**

Identify instructions that can tolerate longer cache hit latencies and remove them from higher-level caches



- Overview of each technique
- Present results
- Present pros & cons



1. **Focusing processor policies via critical-path prediction**

Framework to identify instruction criticality and design to increase instruction scheduling locality in a cluster

2. **Slack: maximizing performance under technological constraints**

Framework to identify variability in instruction execution criticality and application in slowing down execution to save power

3. **Locality vs. Criticality**

Is sacrificing locality for criticality a good idea?

4. **Non-vital loads**

Identify instructions that can tolerate longer cache hit latencies and remove them from higher-level caches



- Out of order processors are highly parallel
- Critical-path analysis identifies dominating instructions
  - Micro-architectural critical path
- Focus on optimizing critical instructions
  - Better arbitration of scarce resources
  - Reduce misspeculation



- Dependence-graph model of micro-architectural critical path
- Hardware predictor to approximate instruction criticality
- Resource arbitration
- Misspeculation reduction

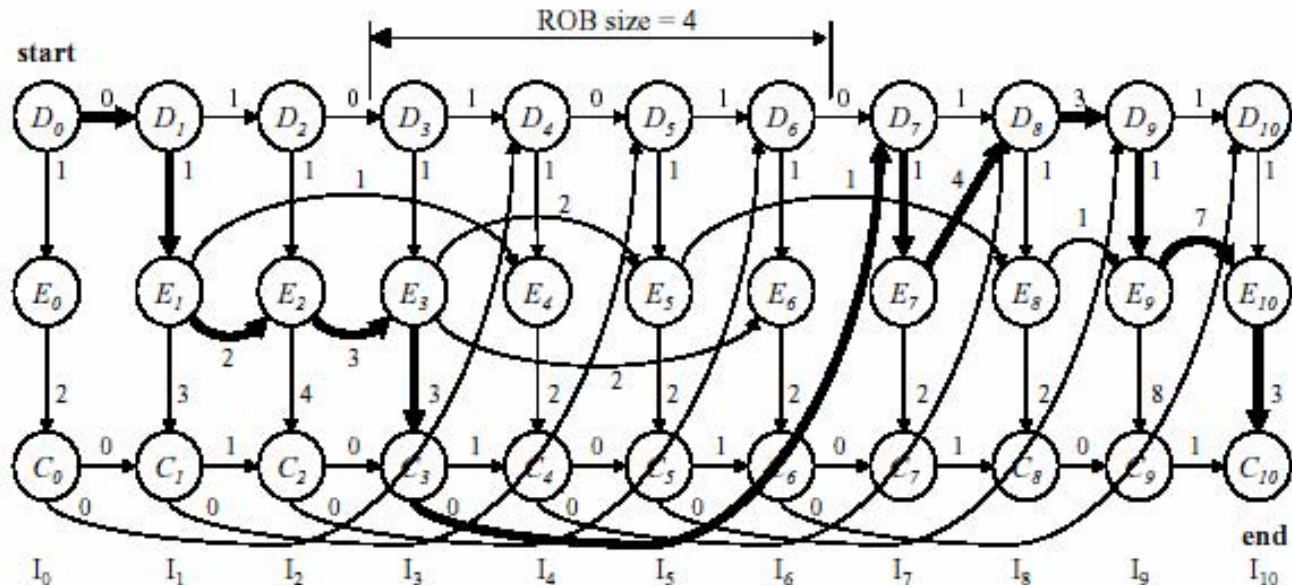


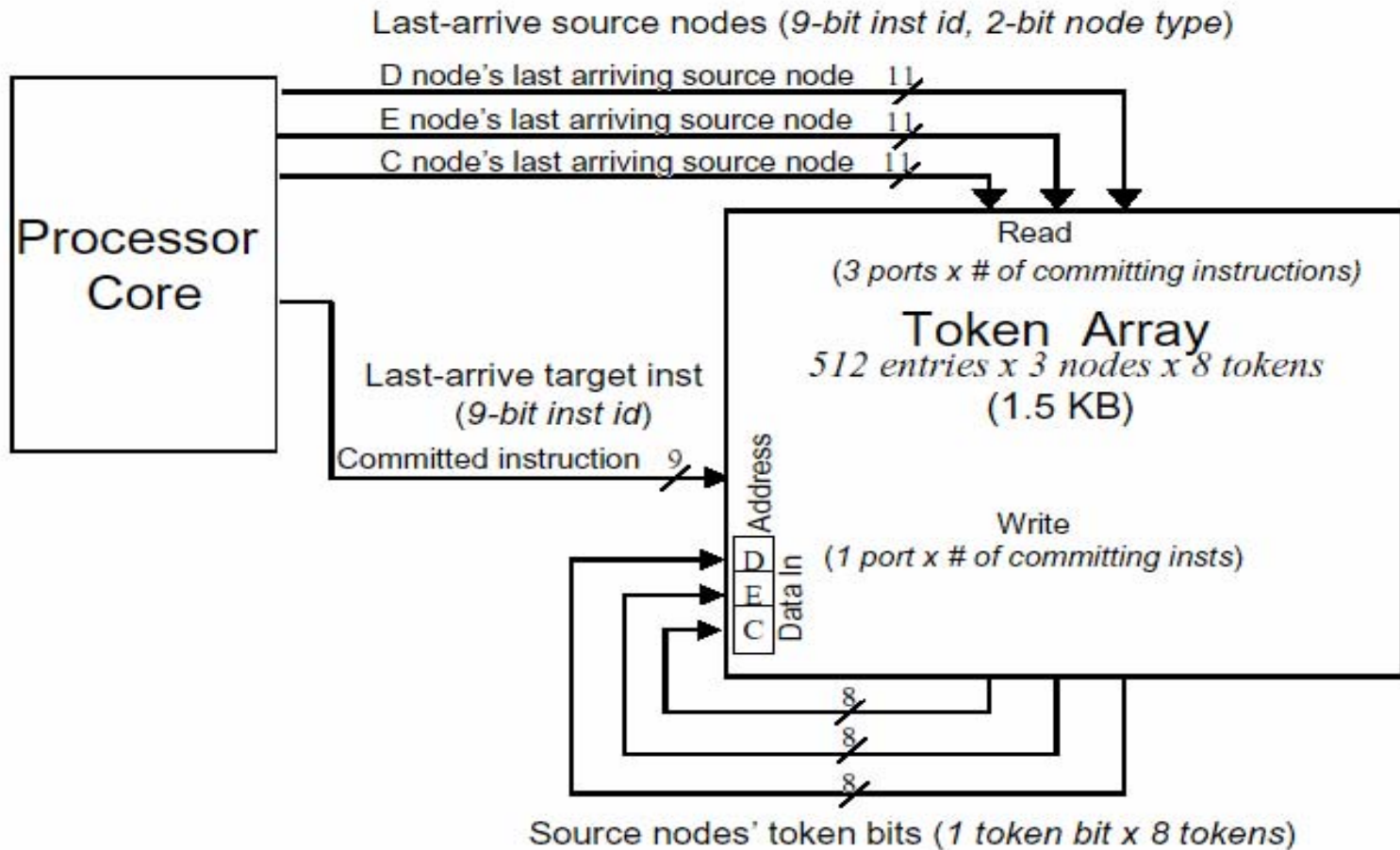
name	constraint modeled	edge	
<i>DD</i>	In-order dispatch	$D_{i-1} \rightarrow D_i$	
<i>CD</i>	Finite re-order buffer	$C_{i-w} \rightarrow D_i$	$w =$ size of the re-order buffer
<i>ED</i>	Control dependence	$E_{i-1} \rightarrow D_i$	inserted if $i - 1$ is a mispredicted branch
<i>DE</i>	Execution follows dispatch	$D_i \rightarrow E_i$	
<i>EE</i>	Data dependences	$E_j \rightarrow E_i$	inserted if instruction $j$ produces an operand of $i$
<i>EC</i>	Commit follows execution	$E_i \rightarrow C_i$	
<i>CC</i>	In-order commit	$C_{i-1} \rightarrow C_i$	

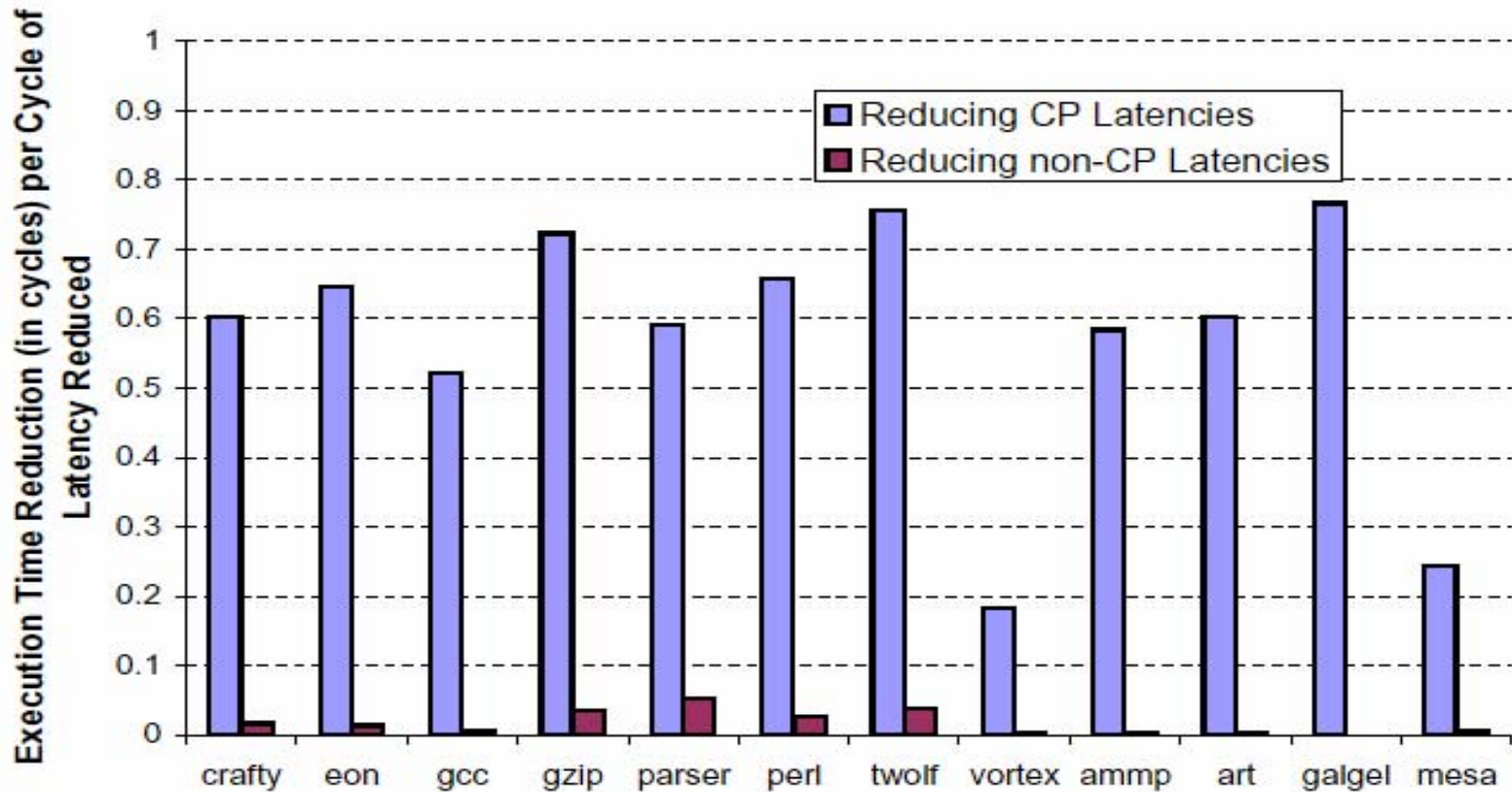
Table 1: Dependences captured by the critical-path model, grouped by the target of the dependence.

$I_0$ :  $r5=0$   
 $I_1$ :  $r3=ld[r2]$   
 L1:  $I_2$ :  $r1=r3*6$   
 $I_3$ :  $r6=ld[r1]$   
 $I_4$ :  $r3=r3+1$   
 $I_5$ :  $r5=r6+r5$   
 $I_6$ :  $cmp R6,0$   
 $I_7$ :  $br L1$   
 $I_8$ :  $r5=r5+100$   
 $I_9$ :  $r0=r5/3$   
 $I_{10}$ :  $Ret r0$

Dynamic Inst. Trace

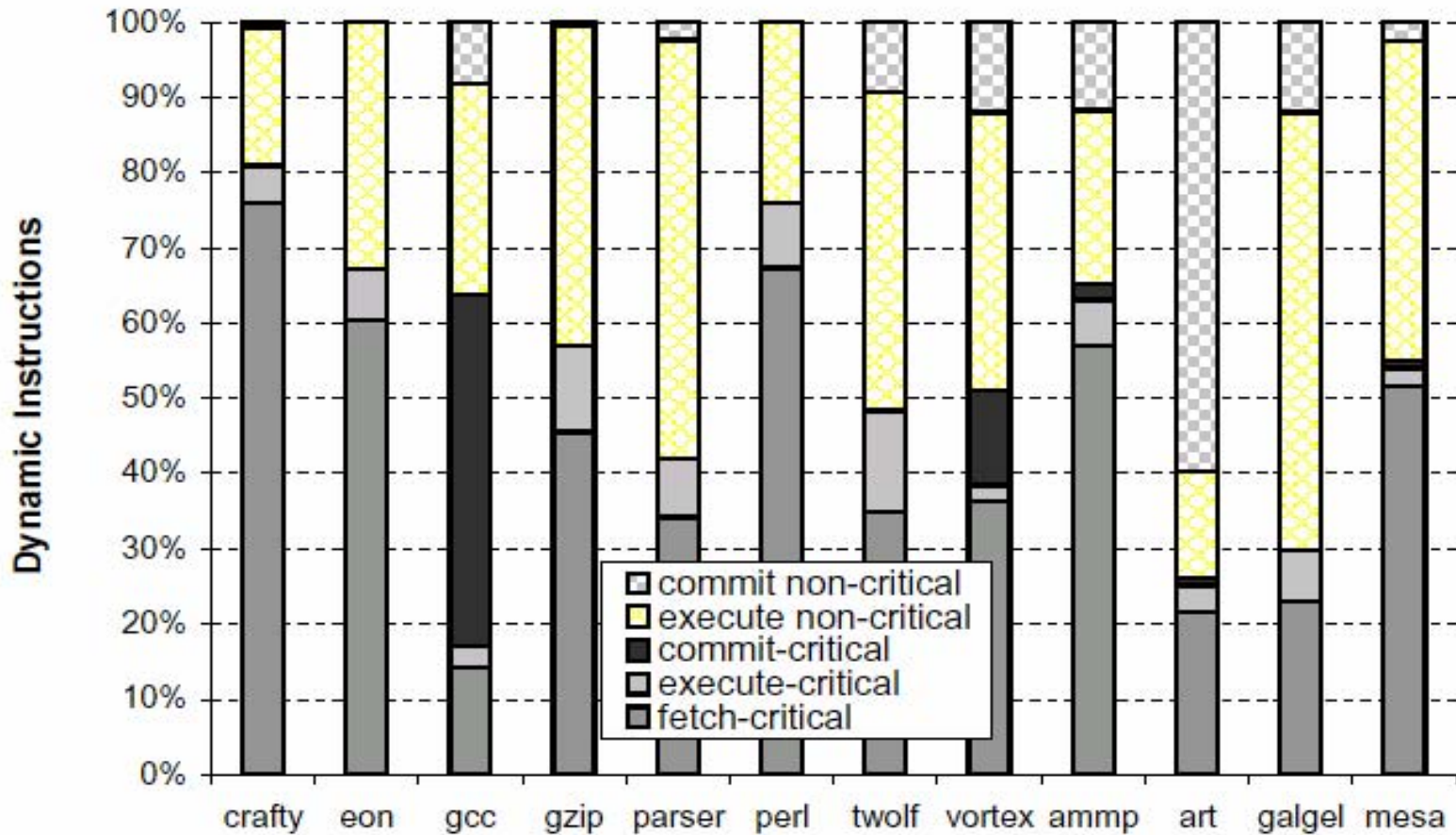




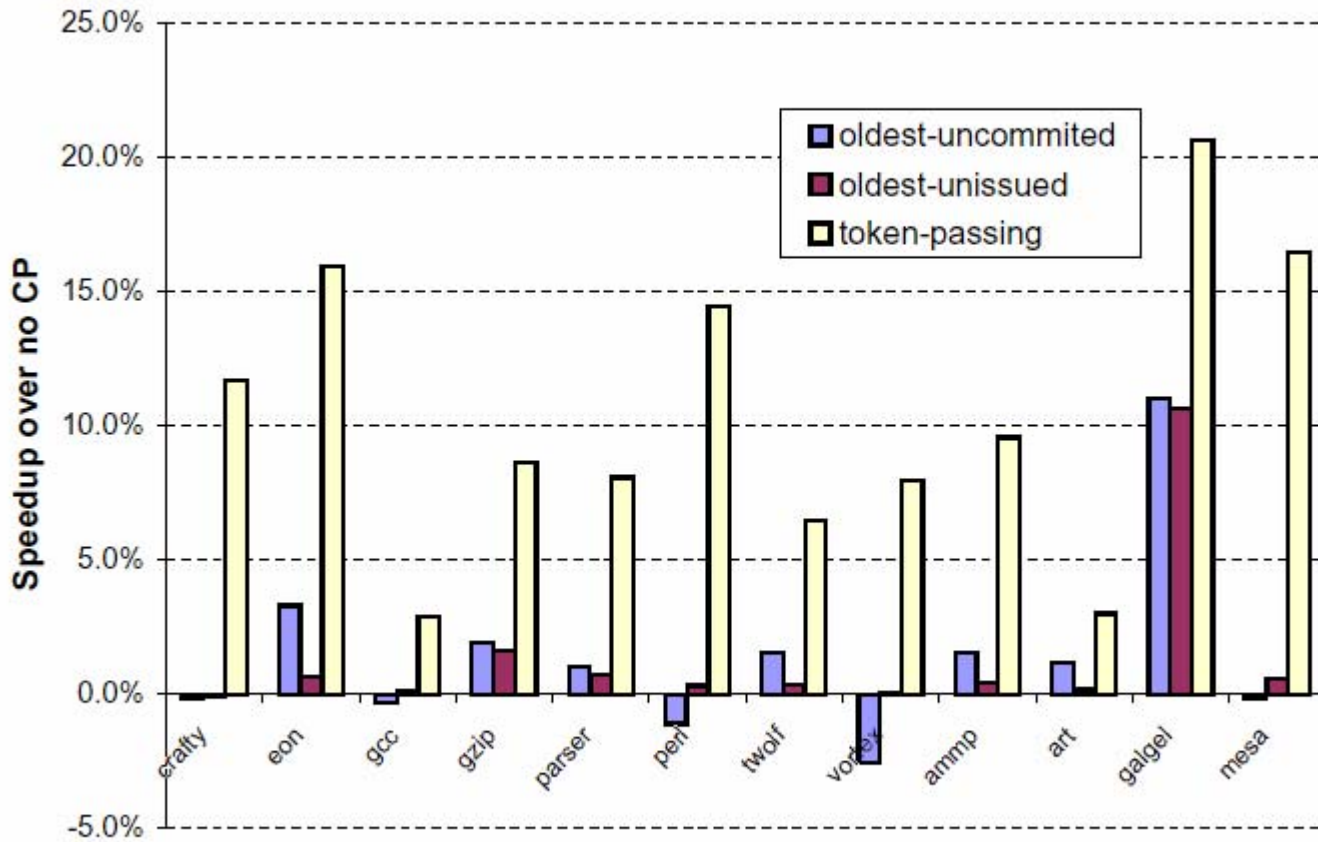


**(a)** Validation of the critical-path model.

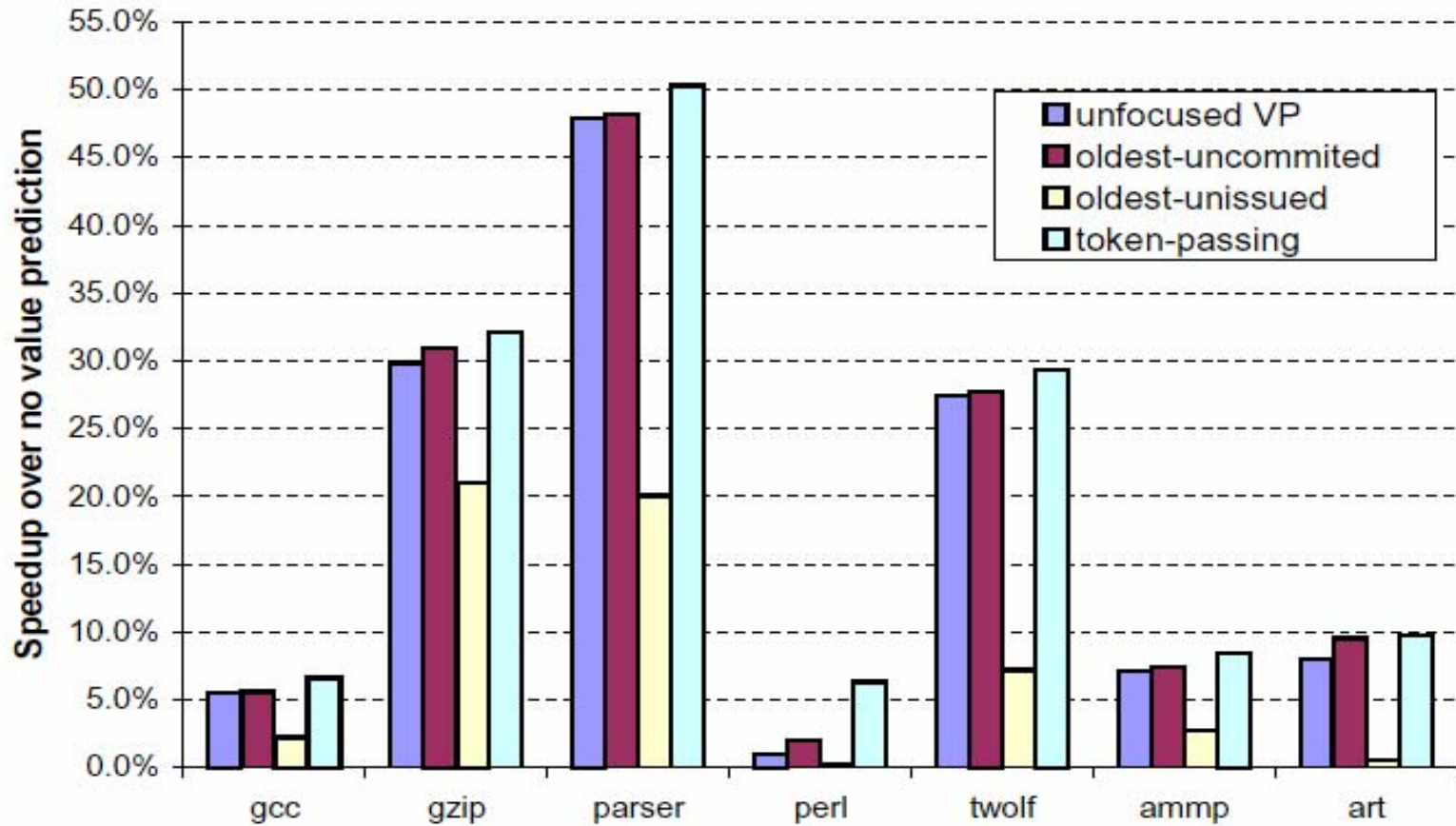




(b) Breakdown of the dynamic instruction count.



(b) Comparison to heuristics-based predictors.



**(b)** Speedup of focused value prediction.



## Prons

- Multipurpose framework
- Adaptive
- Practical algorithmic approach
- Yields interesting results

## Cons

- Wrong path instructions, criticality at memory system ?
- Coarse
- Classification imbalance
- Practical hardware realization



1. **Focusing processor policies via critical-path prediction**

Framework to identify instruction criticality and design to increase instruction scheduling locality in a cluster

2. **Slack: maximizing performance under technological constraints**

Framework to identify variability in instruction execution criticality and application in slowing down execution to save power

3. **Locality vs. Criticality**

Is sacrificing locality for criticality a good idea?

4. **Non-vital loads**

Identify instructions that can tolerate longer cache hit latencies and remove them from higher-level caches

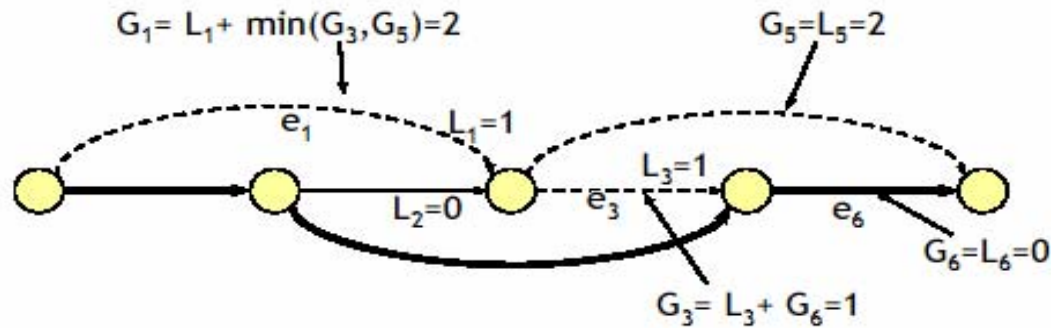
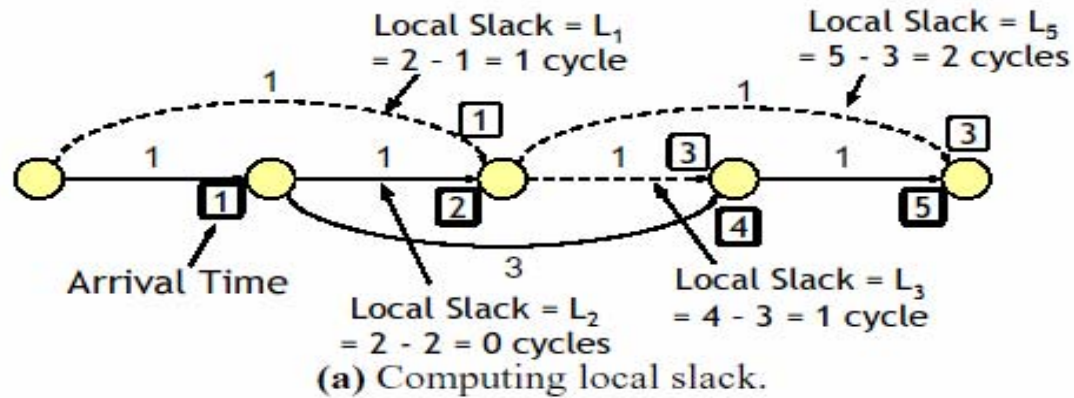


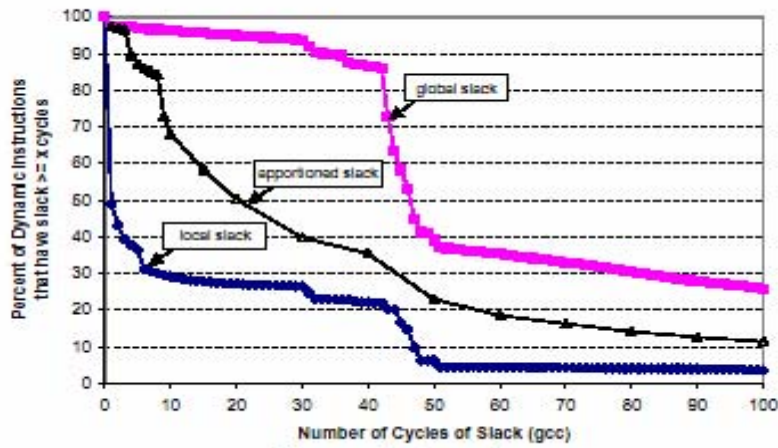


- Technological constraints
  - wire delay, power, complexity
- Non-uniform designs
  - clustering, multi-frequency FU's
- Control policies
- Slack

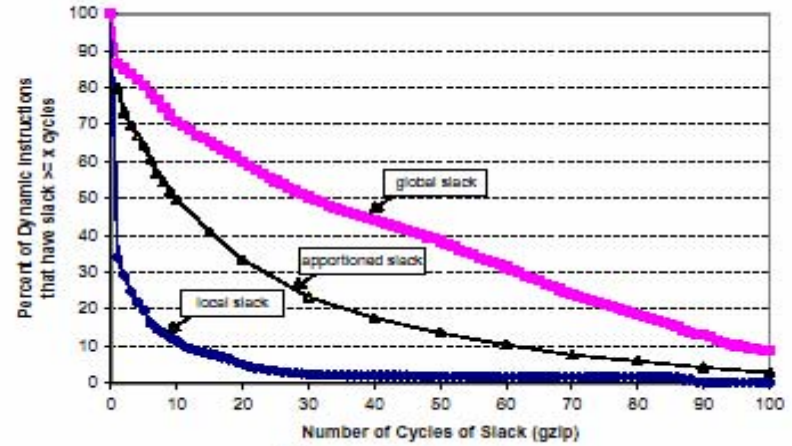


- Model and characterize slack
  - local, global, apportioned
  
- Slack prediction
  - explicit, implicit
  
- Application in power saving hardware
  - fast/slow pipeline

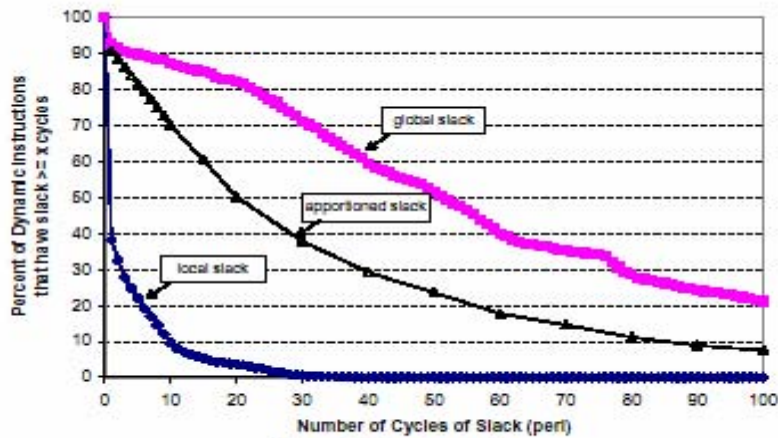




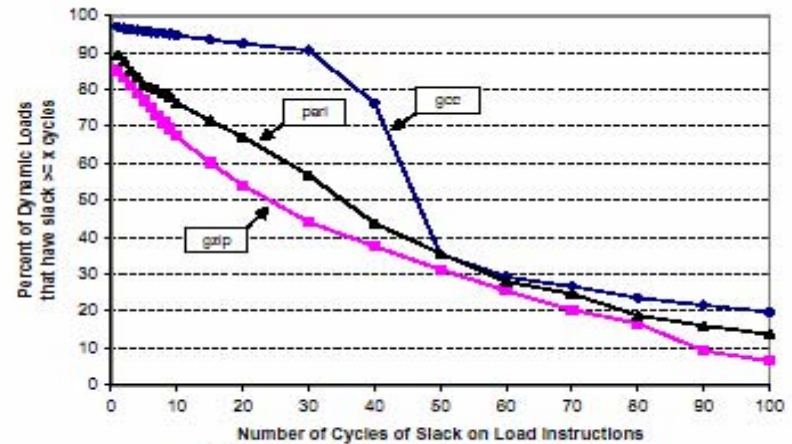
(a) gcc slack results.



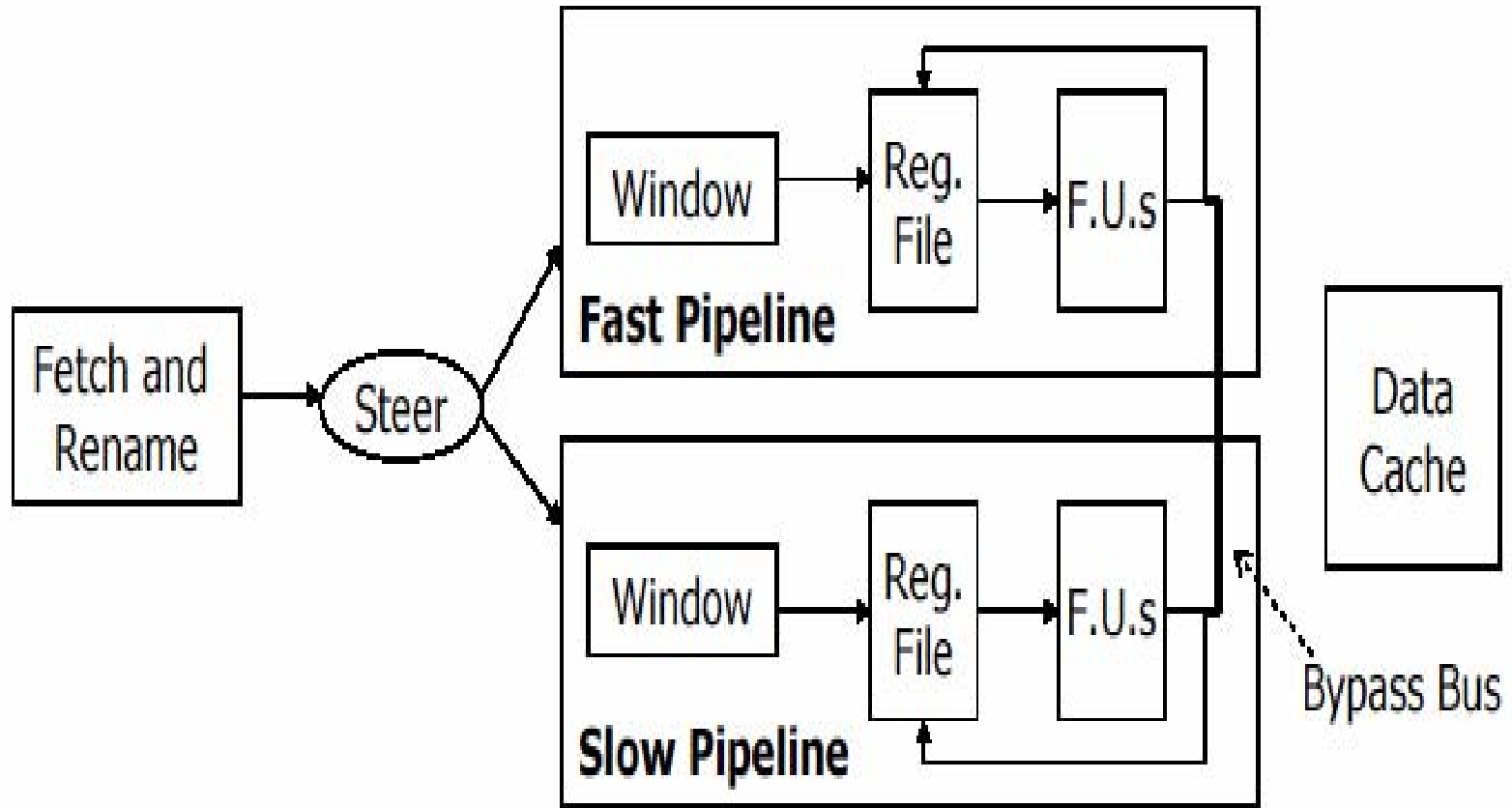
(b) gzip slack results.



(c) perl slack results.



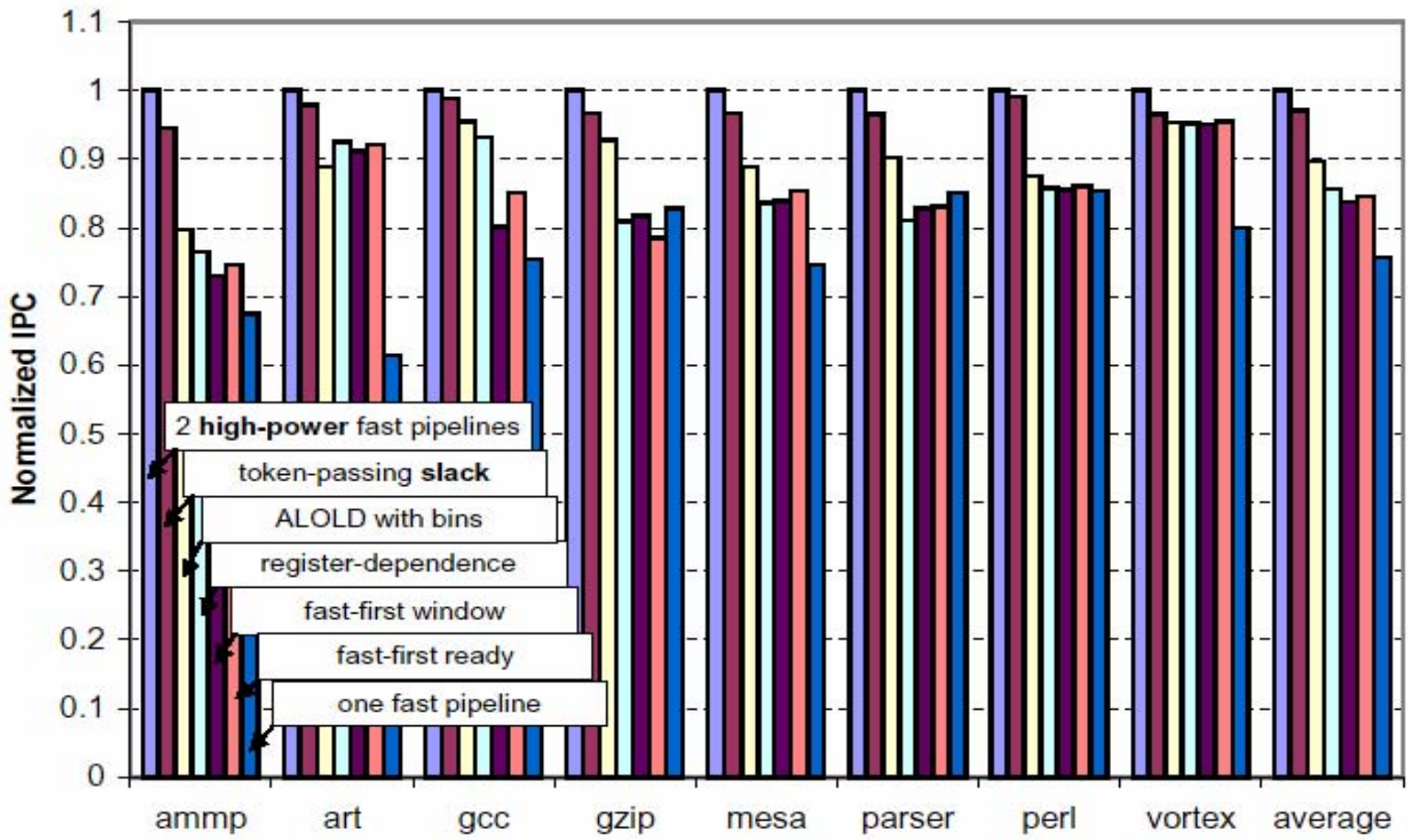
(d) All slack apportioned to loads.





# Paper 2: Performance Impact

Computer Architecture Lab at Carnegie Mellon





## Prons:

- Easy to reason
- Adaptive
- Fine grained

## Cons

- Sampling constraints
- Criticality at memory system ?
- Weakly defined apportioned slack





1. **Focusing processor policies via critical-path prediction**

Framework to identify instruction criticality and design to increase instruction scheduling locality in a cluster

2. **Slack: maximizing performance under technological constraints**

Framework to identify variability in instruction execution criticality and application in slowing down execution to save power

3. **Locality vs. Criticality**

Is sacrificing locality for criticality a good idea?

4. **Non-vital loads**

Identify instructions that can tolerate longer cache hit latencies and remove them from higher-level caches





- Caches exploit spatial and temporal locality
- Locality based schemes ignore nature of misses
- Is criticality a strong enough property to warrant a change in memory hierarchies?
- Classify loads as critical
  - feed mispredicted branch, feed into missed load, few independent instructions following

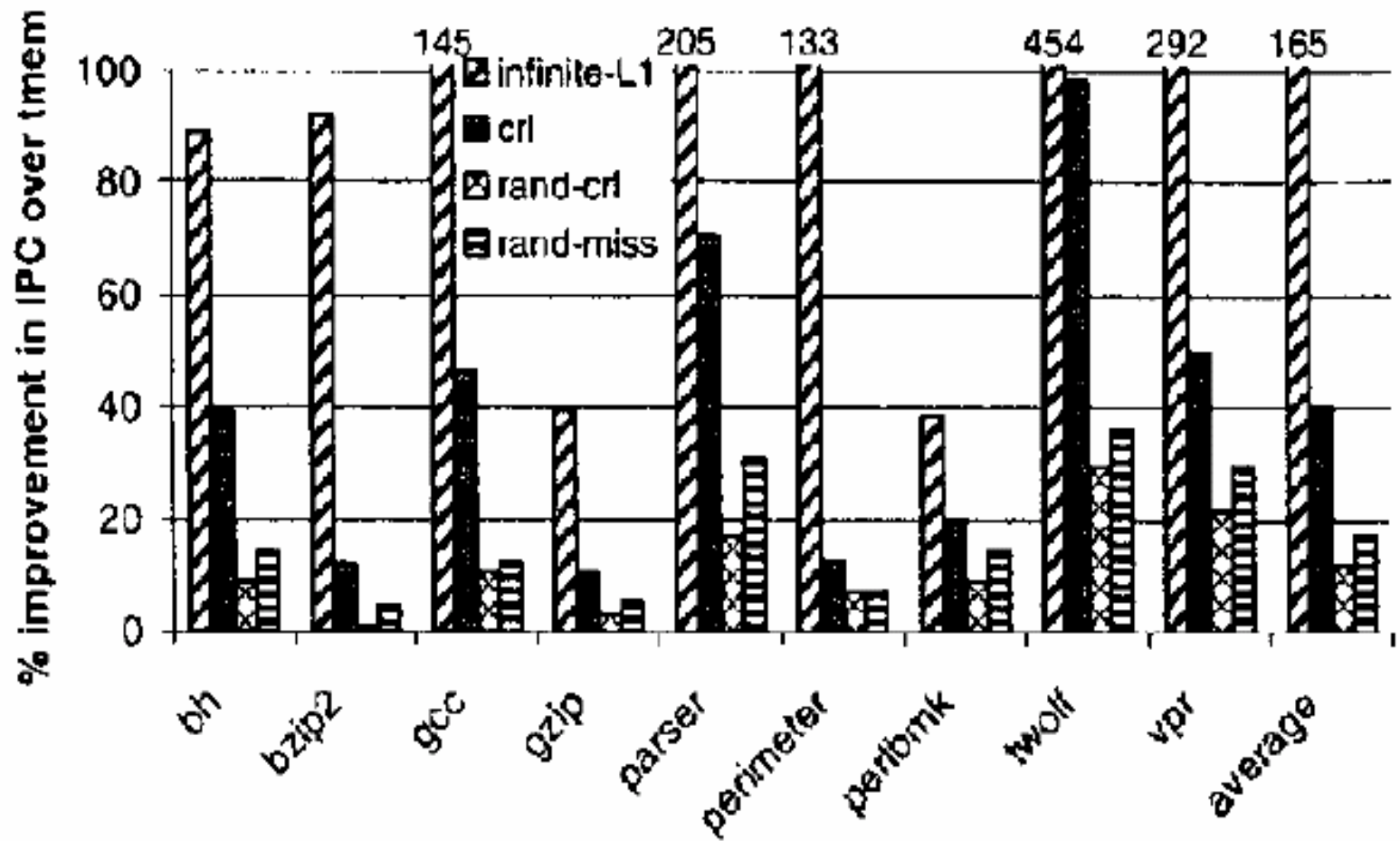


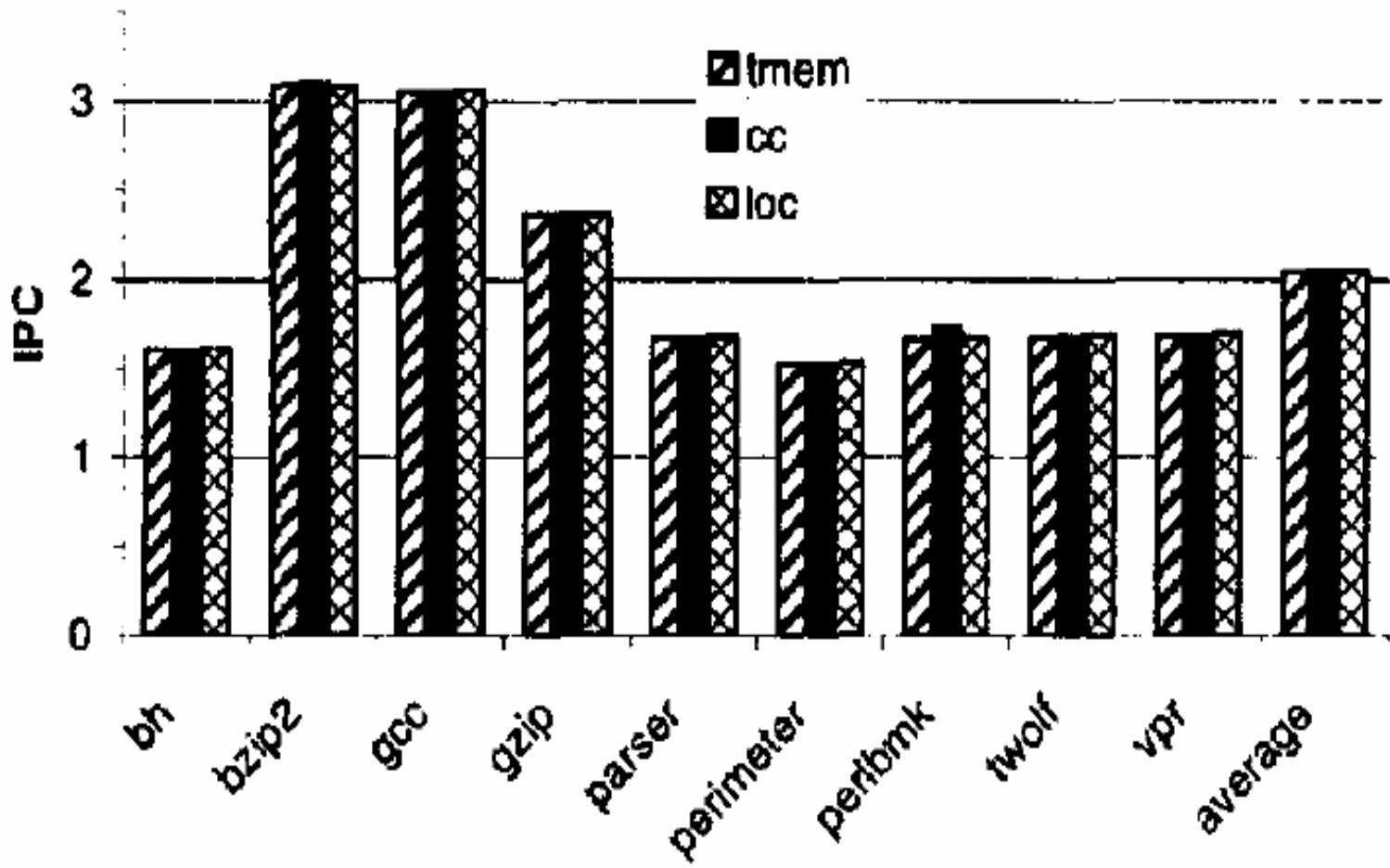
- Provides an answer to the question posed
  - Potential is there
  - Implementations fail

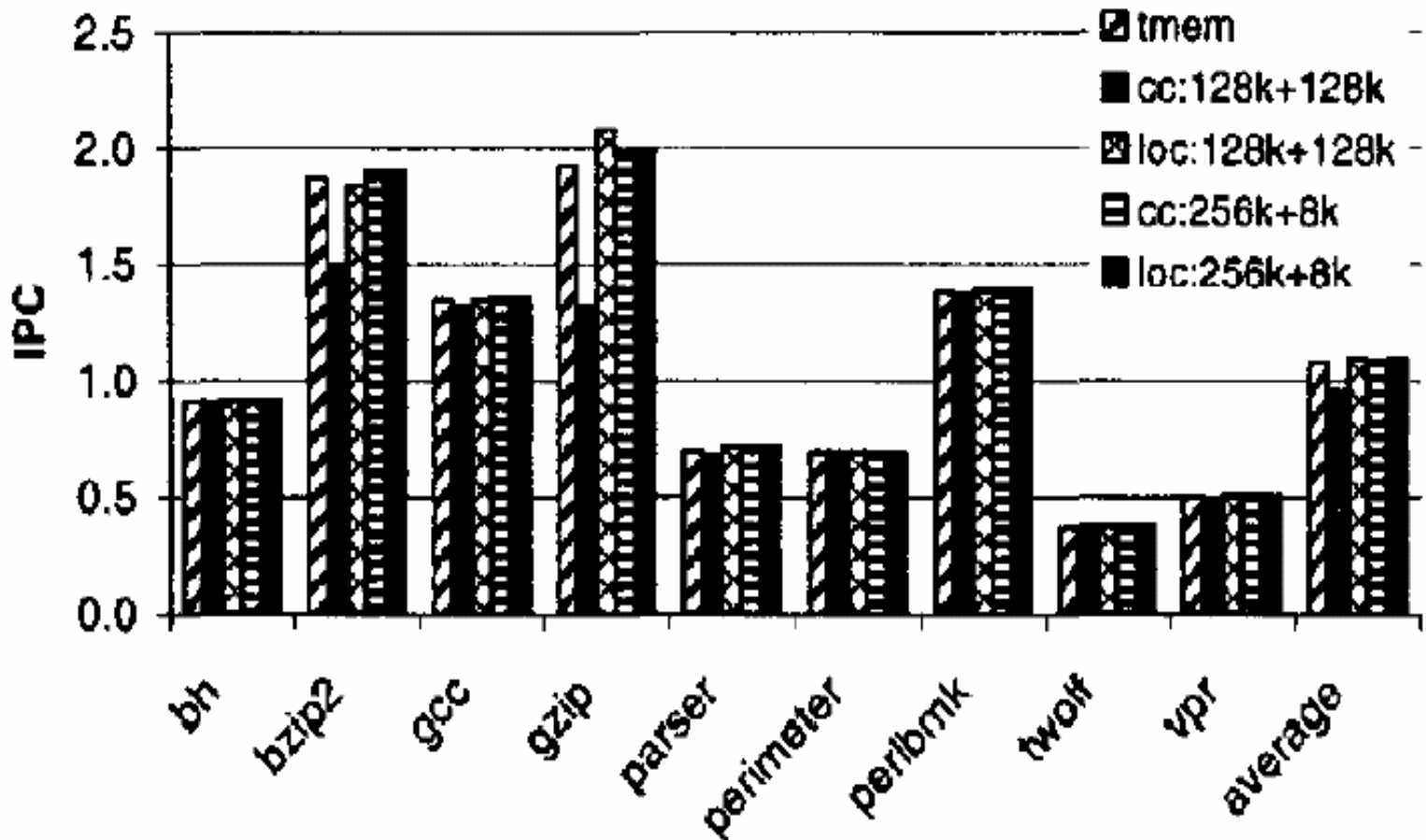


# Paper 3: L1 Criticality Potential

Computer Architecture Lab at Carnegie Mellon



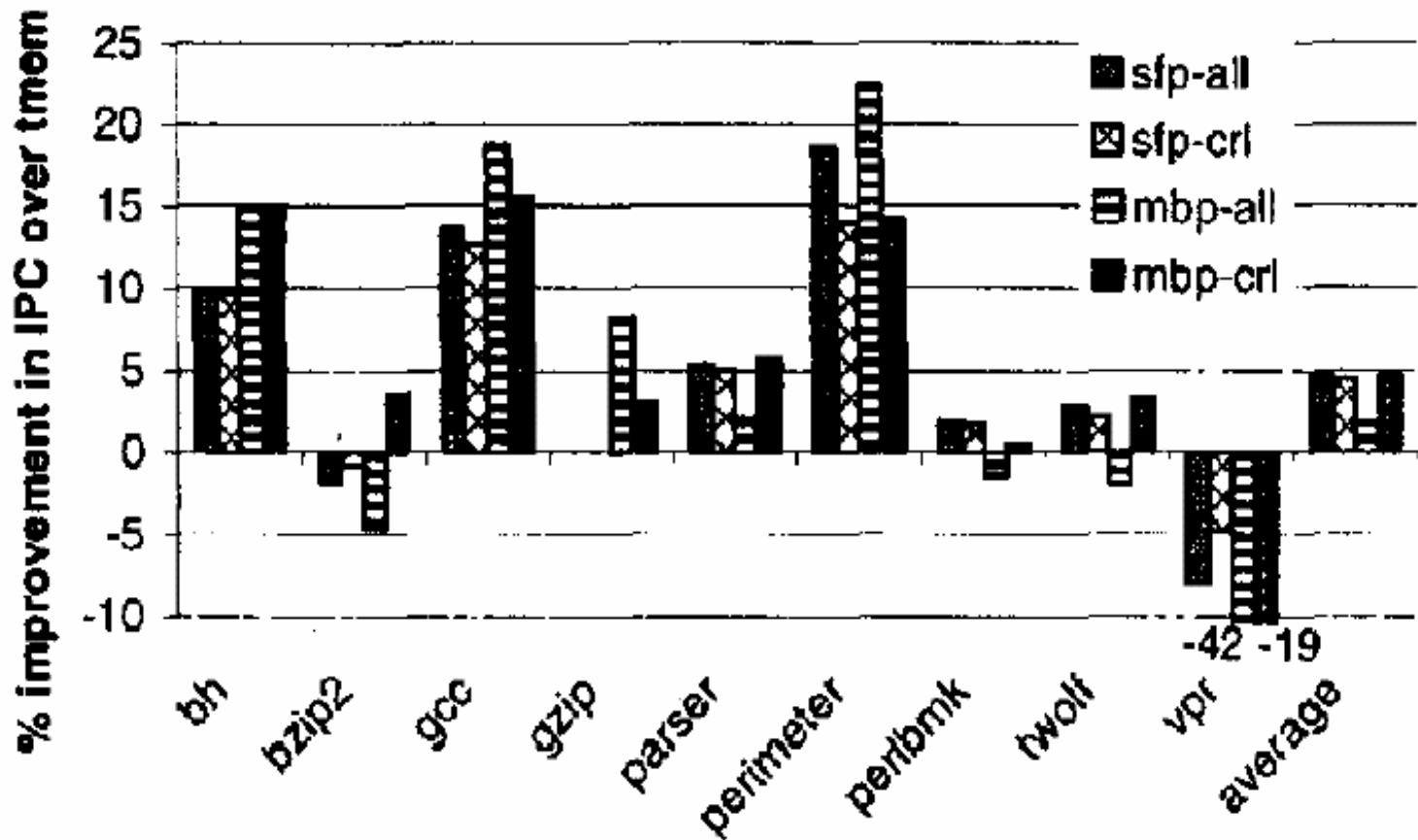






# Paper 3: Criticality Guided Prefetching

Computer Architecture Lab at Carnegie Mellon





## Prons

- Raises important question
- Provides a limit study: the potential is there

## Cons

- Is classification really complete?
- Coverage?
- Accuracy?



1. **Focusing processor policies via critical-path prediction**

Framework to identify instruction criticality and design to increase instruction scheduling locality in a cluster

2. **Slack: maximizing performance under technological constraints**

Framework to identify variability in instruction execution criticality and application in slowing down execution to save power

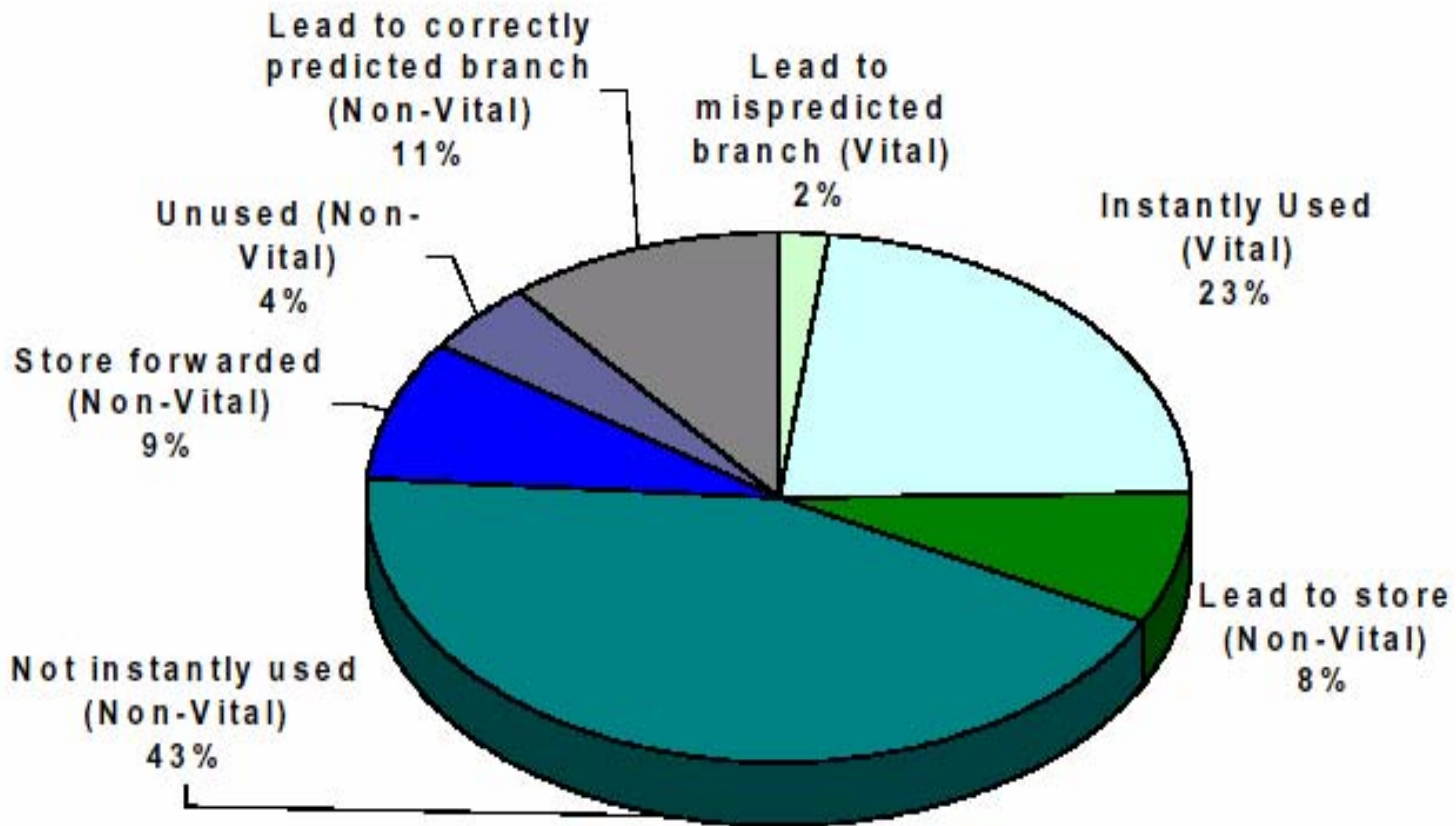
3. **Locality vs. Criticality**

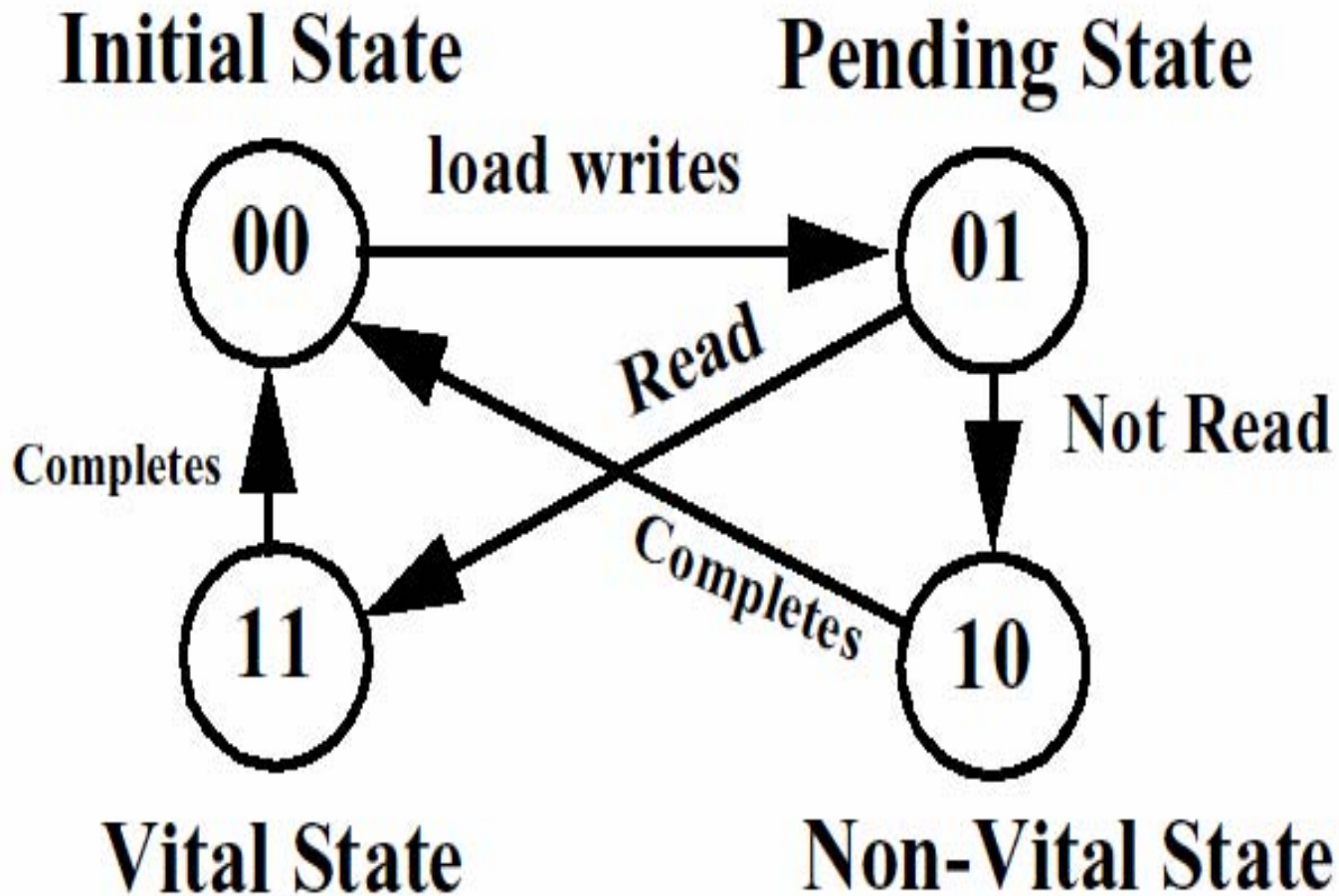
Is sacrificing locality for criticality a good idea?

4. **Non-vital loads**

Identify instructions that can tolerate longer cache hit latencies and remove them from higher-level caches



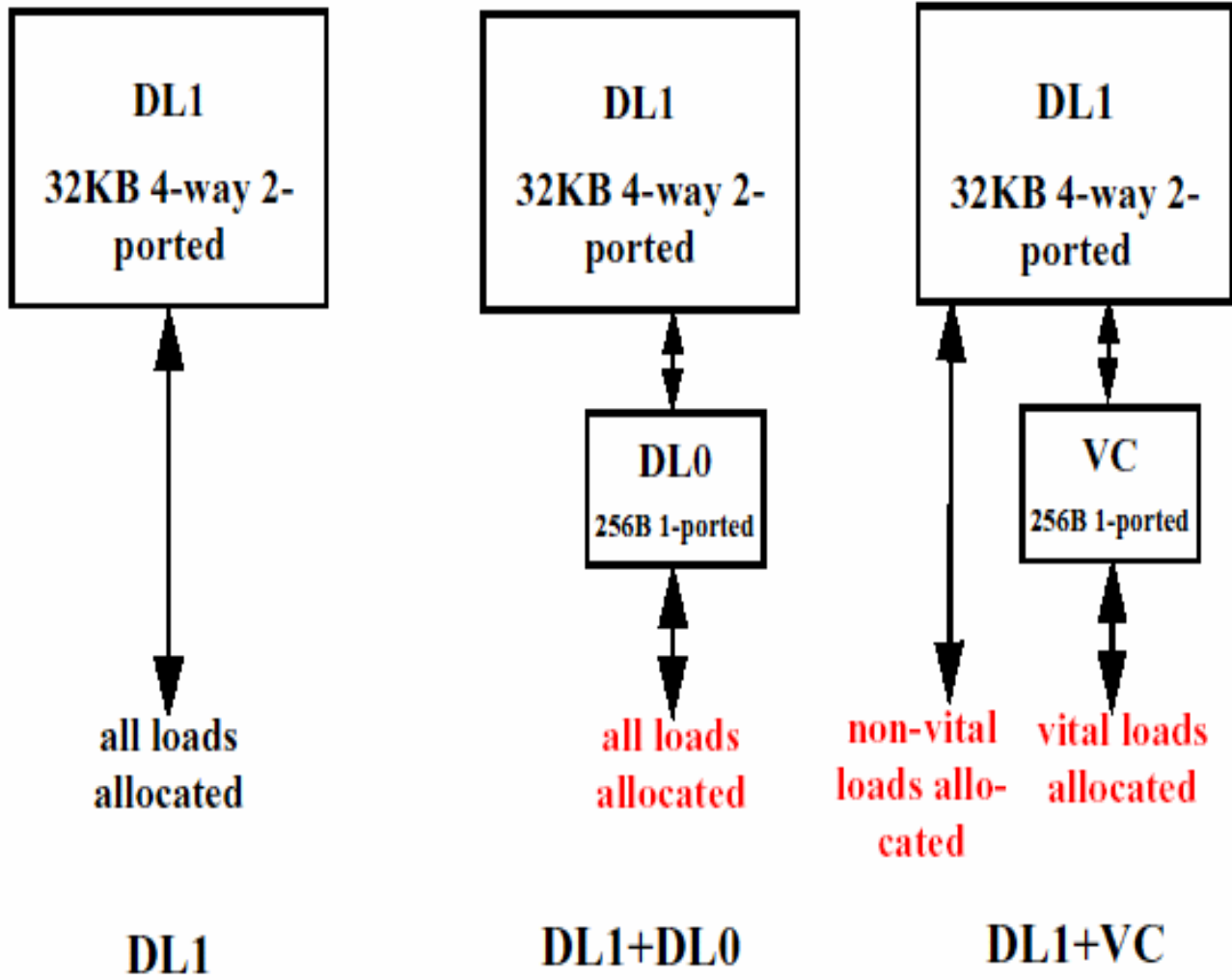


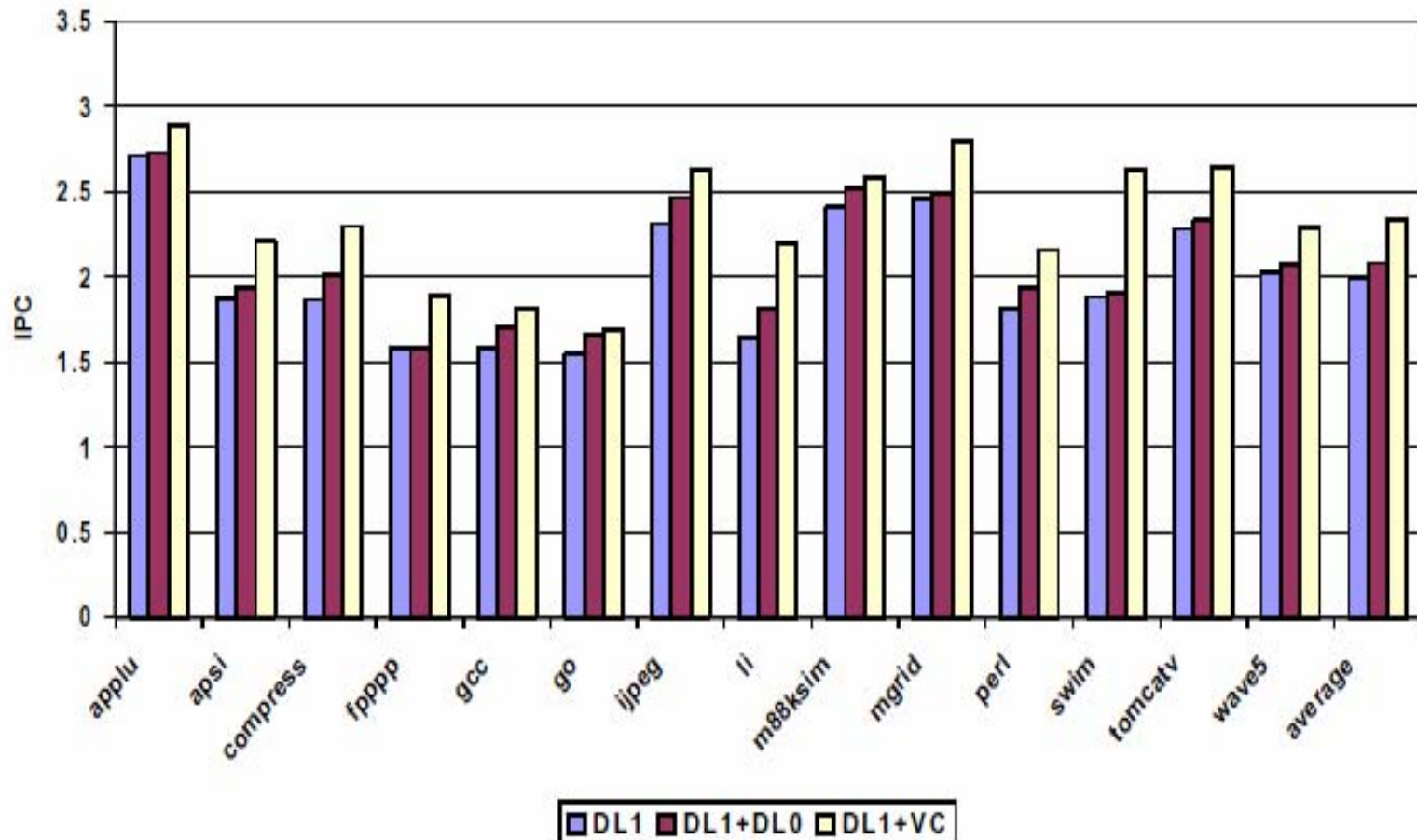




# Paper 4: Cache Hierarchy

Computer Architecture Lab at Carnegie Mellon







## Prons

- Vital Cache improves performance
- Shows more promise than previous paper

## Cons

- Does not explore slack
- Still need an optimality study



- Criticality and Slack can be exploited to guide policies
- Correct classification is key
- Are hardware implementations practical?
- Criticality at the memory system is still challenging
  - Should we give up hope?