
RENO: a RENAME-based instruction Optimizer

Vlad Petric, Tingting Sha, Amir Roth

University of Pennsylvania

RENO...

Is a modified MIPS R10K renamer

Implements dynamic versions of static optimizations

Unifies these previously proposed optimizations

- **Move Elimination** [“Unified Renaming”, Jourdan+]
- **Register Allocation** [“Register Integration”, Petric+]
- **Common Subexpression Elimination** [same]

Adds **Constant Folding**

- Synergistic with other optimizations

Has a simple implementation

Why Optimize in Hardware?

Static optimizations

– file/function boundaries

Dynamic optimizations (RENO)

+ no boundaries

Why Optimize in Hardware?

Static optimizations

- file/function boundaries
- conservative memory info

Dynamic optimizations (RENO)

- + no boundaries
- + can speculate

Why Optimize in Hardware?

Static optimizations

- file/function boundaries
- conservative memory info
- correct on all paths

Dynamic optimizations (RENO)

- + no boundaries
- + can speculate
- + correct on dynamic path

Why Optimize in Hardware?

Static optimizations

- file/function boundaries
- conservative memory info
- correct on all paths
- + removes insns at fetch

Dynamic optimizations (RENO)

- + no boundaries
- + can speculate
- + correct on dynamic path
- removes insns at rename

Why Optimize in Hardware?

Static optimizations	Dynamic optimizations (RENO)
<ul style="list-style-type: none">- file/function boundaries- conservative memory info- correct on all paths+ removes insns at fetch	<ul style="list-style-type: none">+ no boundaries+ can speculate+ correct on dynamic path- removes insns at rename

RENO complements static compiler optimizations

- Removes 23% of instructions from -O4 optimized programs

Overview

RENO framework: physical register sharing

RENO optimizations

- Focus on constant folding

Hardware implementation

Performance evaluation

Physical Register Sharing

RENO examines dynamic instructions “one at a time”

Optimizes away some instructions

How? **physical register sharing**

- Correct value already exists in some physical register
 - 1 Instruction’s map-table output set to that register
 - 2 Instruction itself bypasses execution core

Performance and implementation benefits

- + Reduces dataflow graph latency
- + Amplifies execution core bandwidth
- + Map-table manipulations only

RENO optimization: way to detect sharing opportunities

Conventional Processing

Static insn	MapTable action	Execution dataflow
op1 _, _ → r1		
...		
addi r1, 0 → r2		
...		
op2 _, r2 → _		

The diagram illustrates the execution dataflow for three instructions. The first instruction, 'op1', has two blank operands and a result register 'r1'. The second instruction, 'addi', has 'r1' as the first operand, '0' as the second operand, and 'r2' as the result register. The third instruction, 'op2', has a blank first operand, 'r2' as the second operand, and a blank result register. A green arrow points from the 'r1' result of the first instruction to the 'r1' operand of the second instruction. A purple arrow points from the 'r2' result of the second instruction to the 'r2' operand of the third instruction.

Conventional Processing

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3]	op1 <u>, <u> → p3</u></u>
...		
addi r1, 0 → r2		
...		
op2 <u>, r2 → <u></u></u>		

Conventional Processing

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3]	op1 <u>, <u>, → p3</u></u>
...		↓
addi r1, 0 → r2	r2 := [p9]	addi p3, 0 → p9
...		
op2 <u>, r2 → <u></u></u>		

Conventional Processing

Static insn	MapTable action	Execution dataflow
op1 $_$, $_$ \rightarrow r1	r1 := [p3]	op1 $_$, $_$, \rightarrow p3
...		↓
addi r1, 0 \rightarrow r2	r2 := [p9]	addi p3, 0 \rightarrow p9
...		↓
op2 $_$, r2 \rightarrow $_$	$_$:= []	op2 $_$, p9 \rightarrow $_$

RENO_{ME}: Move Elimination

Static insn	MapTable action	Execution dataflow
op1 <code>_</code> , <code>_</code> \rightarrow <code>r1</code>		
...		
<code>addi</code> <code>r1</code> , <code>0</code> \rightarrow <code>r2</code>		
...		
op2 <code>_</code> , <code>r2</code> \rightarrow <code>_</code>		

Optimization detection: `addi` `_`, `0` \rightarrow `_`

RENO_{ME}: Move Elimination

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3]	op1 <u>, <u>, → p3</u></u>
...		
addi r1, 0 → r2		
...		
op2 <u>, r2 → <u></u></u>		

Optimization detection: addi , 0 →

RENO_{ME}: Move Elimination

Static insn	MapTable action	Execution dataflow	
op1 <u>, <u> → r1</u></u>	r1 := [p3]	<table border="1"><tr><td>op1 <u>, <u>, → p3</u></u></td></tr></table>	op1 <u>, <u>, → p3</u></u>
op1 <u>, <u>, → p3</u></u>			
...			
addi r1, 0 → r2	r2 := [p3] (reuse)		
...			
op2 <u>, r2 → <u></u></u>			

Optimization detection: addi , 0 →

RENO_{ME}: Move Elimination

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3]	op1 <u>, <u>, → p3</u></u>
...		
addi r1, 0 → r2	r2 := [p3] (reuse)	
...		
op2 <u>, r2 → <u< td=""><td><u> := []</u></td><td>op2 <u>, p3 → <u< td=""></u<></u></td></u<></u>	<u> := []</u>	op2 <u>, p3 → <u< td=""></u<></u>

Optimization detection: addi , 0 →

RENO_{ME}: Move Elimination

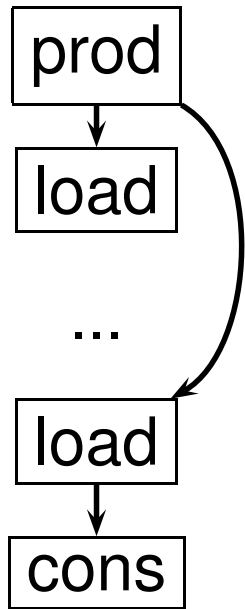
Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3]	op1 <u>, <u>, → p3</u></u>
...		
addi r1, 0 → r2	r2 := [p3] (reuse)	
...		
op2 <u>, r2 → <u< td=""><td><u> := []</u></td><td>op2 <u>, p3 → <u< td=""></u<></u></td></u<></u>	<u> := []</u>	op2 <u>, p3 → <u< td=""></u<></u>

Optimization detection: **addi** , 0 →

+ Latency: **move** removed from dataflow graph

+ Bandwidth: **move** removed from execution core

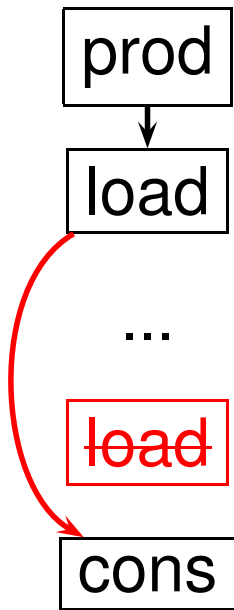
RENO_{CSE/RA}: Register Integration



CSE/RegAlloc using physical register sharing

- RegAlloc: short-circuit stack store-load pairs
 - aka Speculative memory bypassing
- Optimization detection: memoization table (expensive)

RENO_{CSE/RA}: Register Integration



CSE/RegAlloc using physical register sharing

- RegAlloc: short-circuit stack store-load pairs
 - aka Speculative memory bypassing
- Optimization detection: memoization table (expensive)

RENO_{CSE/RA}: Register Integration



CSE/RegAlloc using physical register sharing

- RegAlloc: short-circuit stack store-load pairs
 - aka Speculative memory bypassing
- Optimization detection: memoization table (expensive)

RENO_{CSE/RA}: Register Integration



CSE/RegAlloc using physical register sharing

- RegAlloc: short-circuit stack store-load pairs
 - aka Speculative memory bypassing
- Optimization detection: memoization table (expensive)

RENO_{CF}: Constant Folding

RENO_{ME}: register moves (**addi** , 0 →)

- Not common: compilation artifacts

RENO_{CF}: Constant Folding

RENO_{ME}: register moves (**addi** `_, 0` → `_`)

- Not common: compilation artifacts

RENO_{CF}: register-immediate additions (**addi** `_, 4` → `_`)

+ Surprisingly common: address calculation, etc.

- 12% in SPECint
- 16% in MediaBench

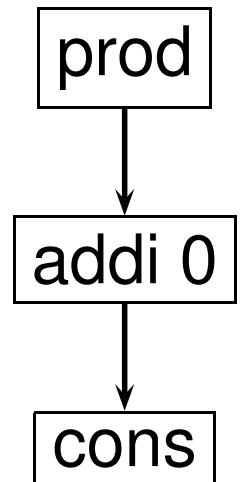
+ Synergy with **CSE/RegAlloc**

- Main benefit of CSE/RegAlloc: eliminating loads
- Address calculations are dataflow “glue” between loads
- Eliminate without expensive table lookups

Physical Register Sharing++

RENO_{ME/CSE/RA}

- Instructions don't compute "new" results
- Physical register sharing sufficient



Physical Register Sharing++

RENO_{ME/CSE/RA}

- Instructions don't compute "new" results
- Physical register sharing sufficient

prod

~~addi 0~~

cons



Physical Register Sharing++

RENO_{ME/CSE/RA}

- Instructions don't compute "new" results
- Physical register sharing sufficient

prod

~~addi 0~~

cons

RENO_{CF}

- Instructions compute "new" results
- Physical register sharing insufficient

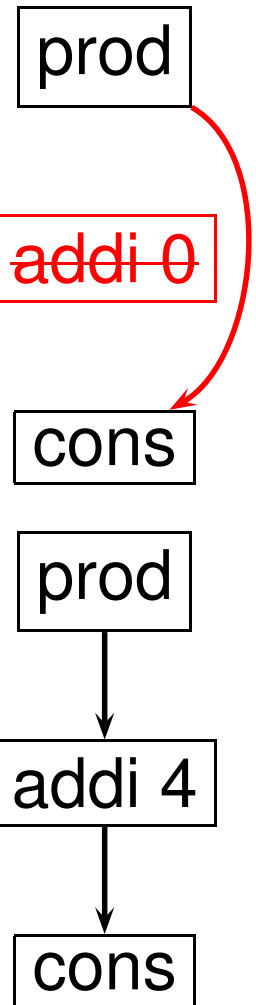
Physical Register Sharing++

$RENO_{ME/CSE/RA}$

- Instructions don't compute "new" results
- Physical register sharing sufficient

$RENO_{CF}$

- Instructions compute "new" results
- Physical register sharing insufficient



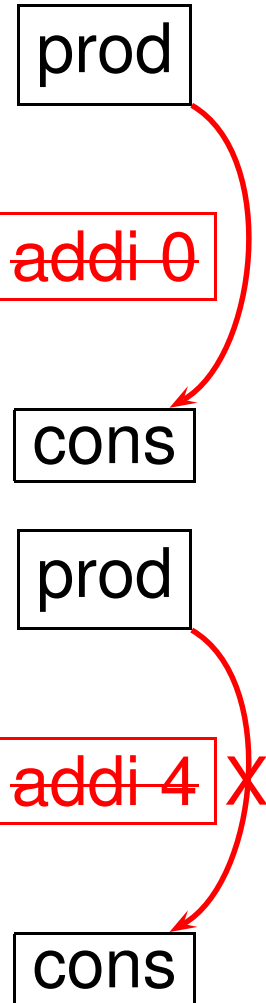
Physical Register Sharing++

RENO_{ME/CSE/RA}

- Instructions don't compute "new" results
- Physical register sharing sufficient

RENO_{CF}

- Instructions compute "new" results
- Physical register sharing insufficient



Physical Register Sharing++

RENO_{ME/CSE/RA}

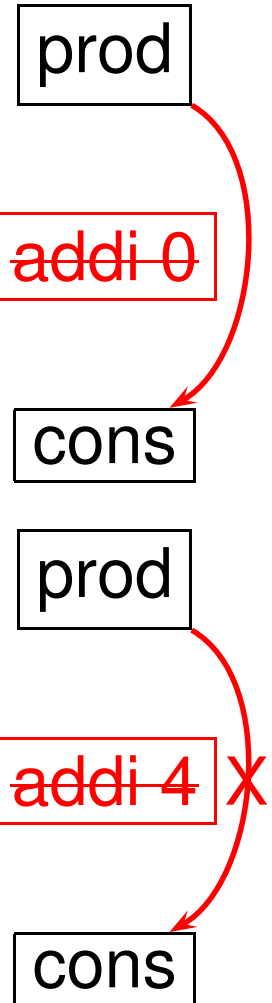
- Instructions don't compute "new" results
- Physical register sharing sufficient

RENO_{CF}

- Instructions compute "new" results
- Physical register sharing insufficient

Solution: extended renaming semantics

- Extended map table: LREG := [PREG, **DISP**]
- **+DISP** implicitly fused to any consumer



Physical Register Sharing++

RENO_{ME/CSE/RA}

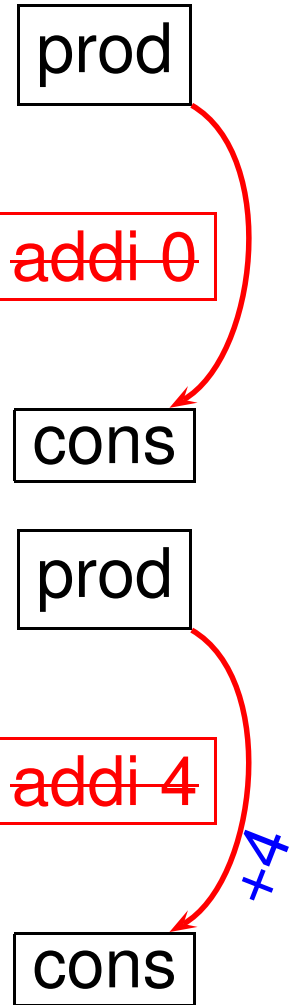
- Instructions don't compute "new" results
- Physical register sharing sufficient

RENO_{CF}

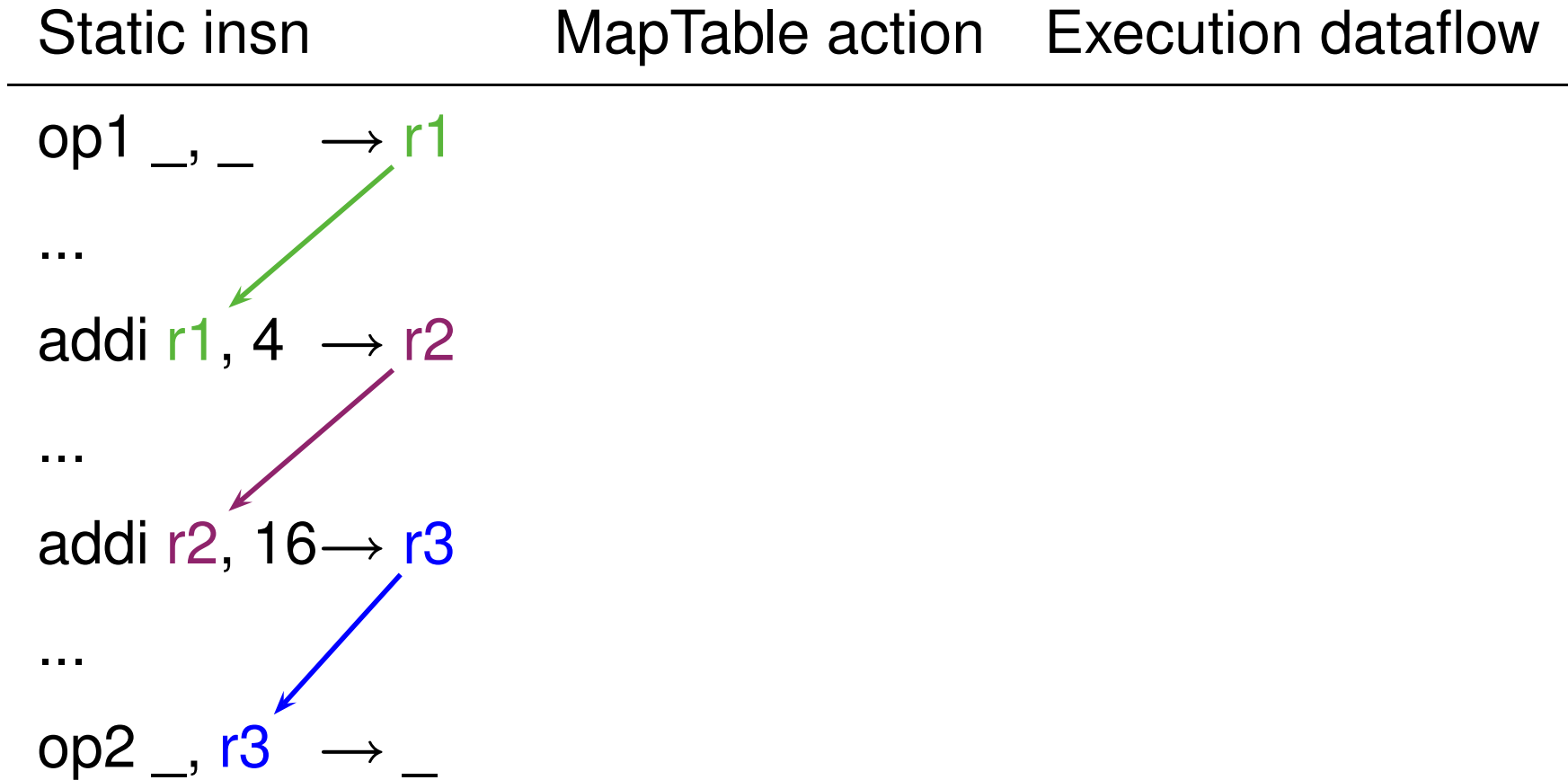
- Instructions compute "new" results
- Physical register sharing insufficient

Solution: extended renaming semantics

- Extended map table: LREG := [PREG, **DISP**]
- **+DISP** implicitly fused to any consumer



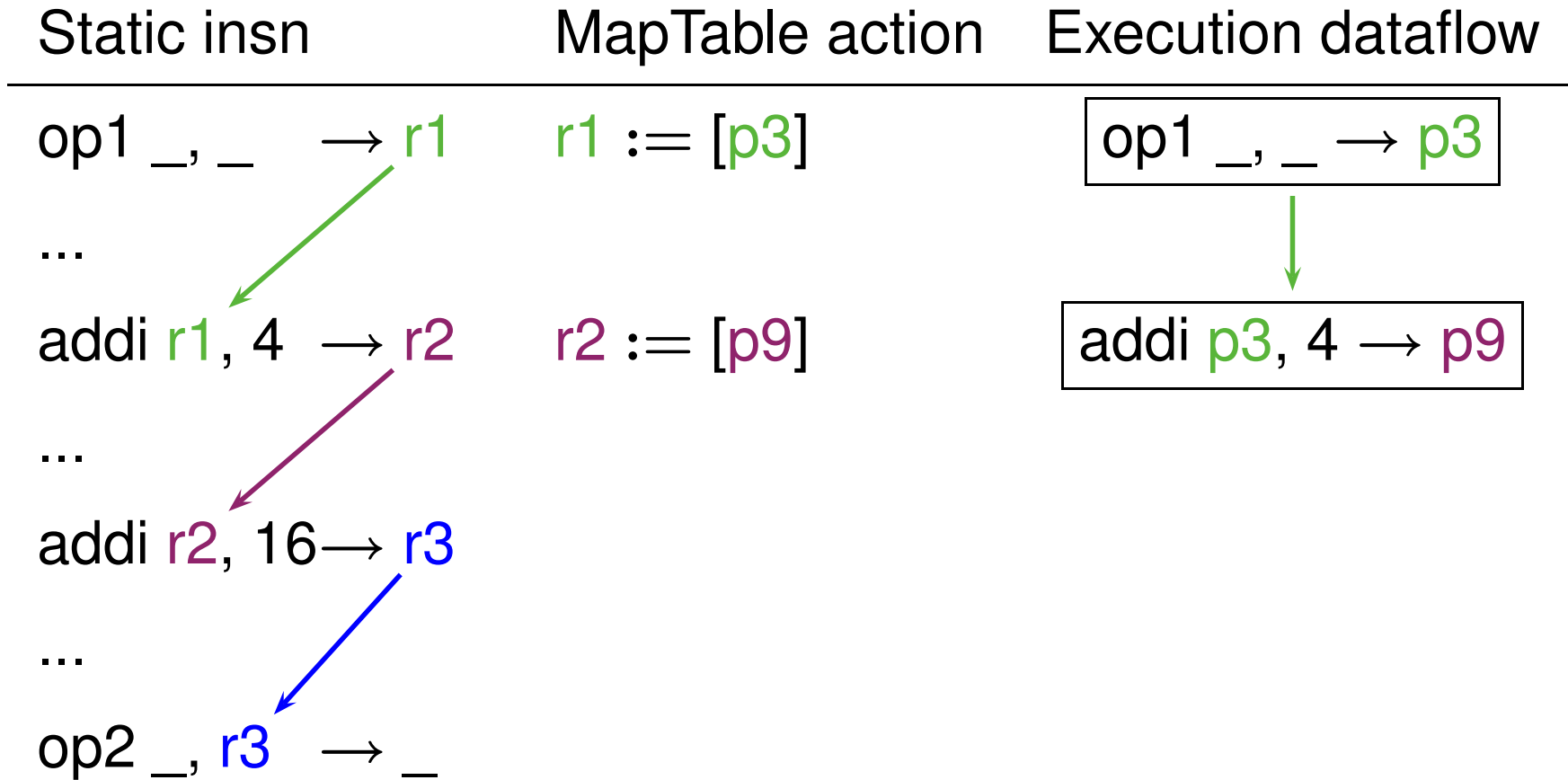
Conventional Processing



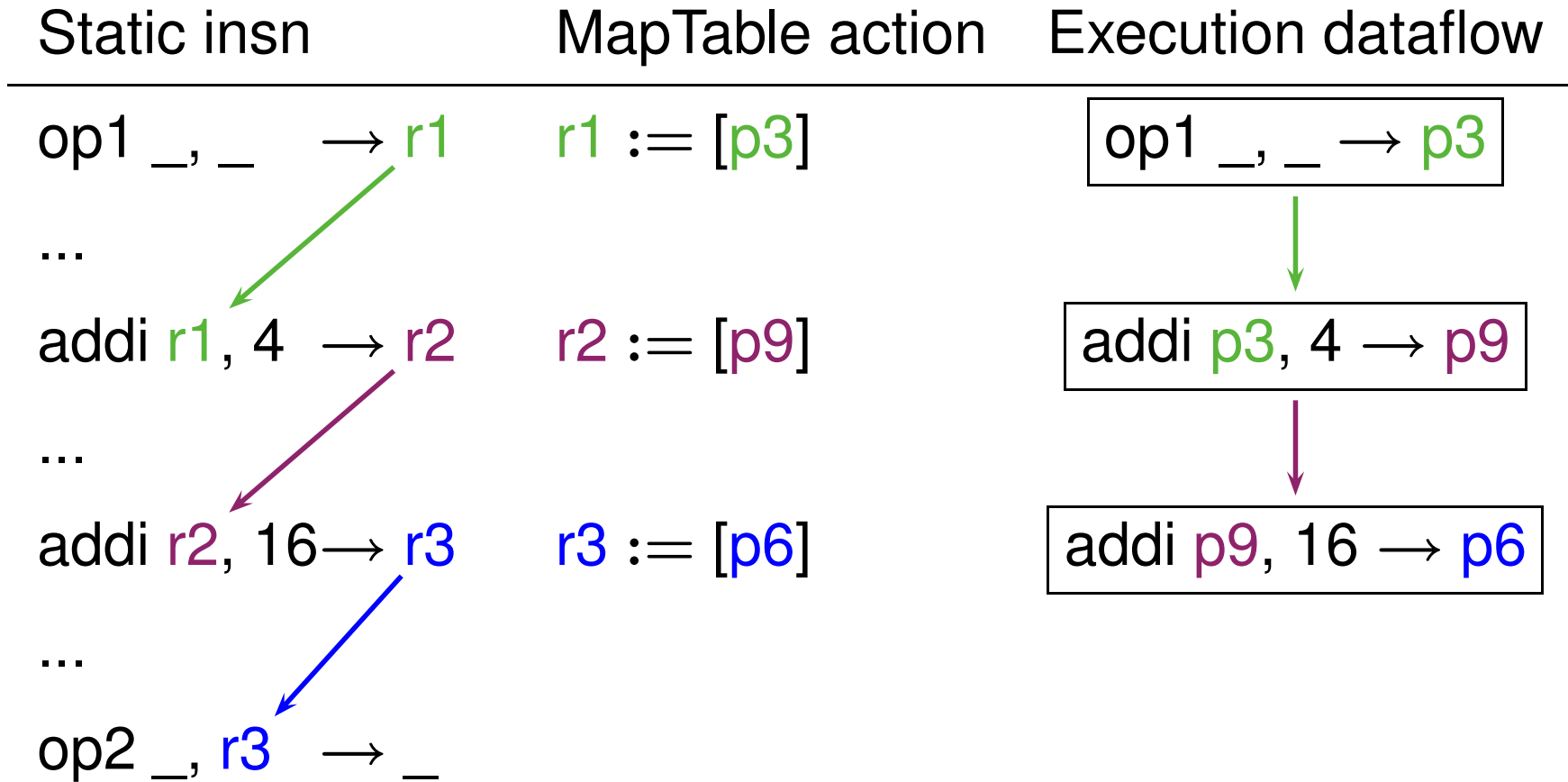
Conventional Processing

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3]	op1 <u>, <u> → p3</u></u>
...		
addi r1, 4 → r2		
...		
addi r2, 16 → r3		
...		
op2 <u>, r3 → <u></u></u>		

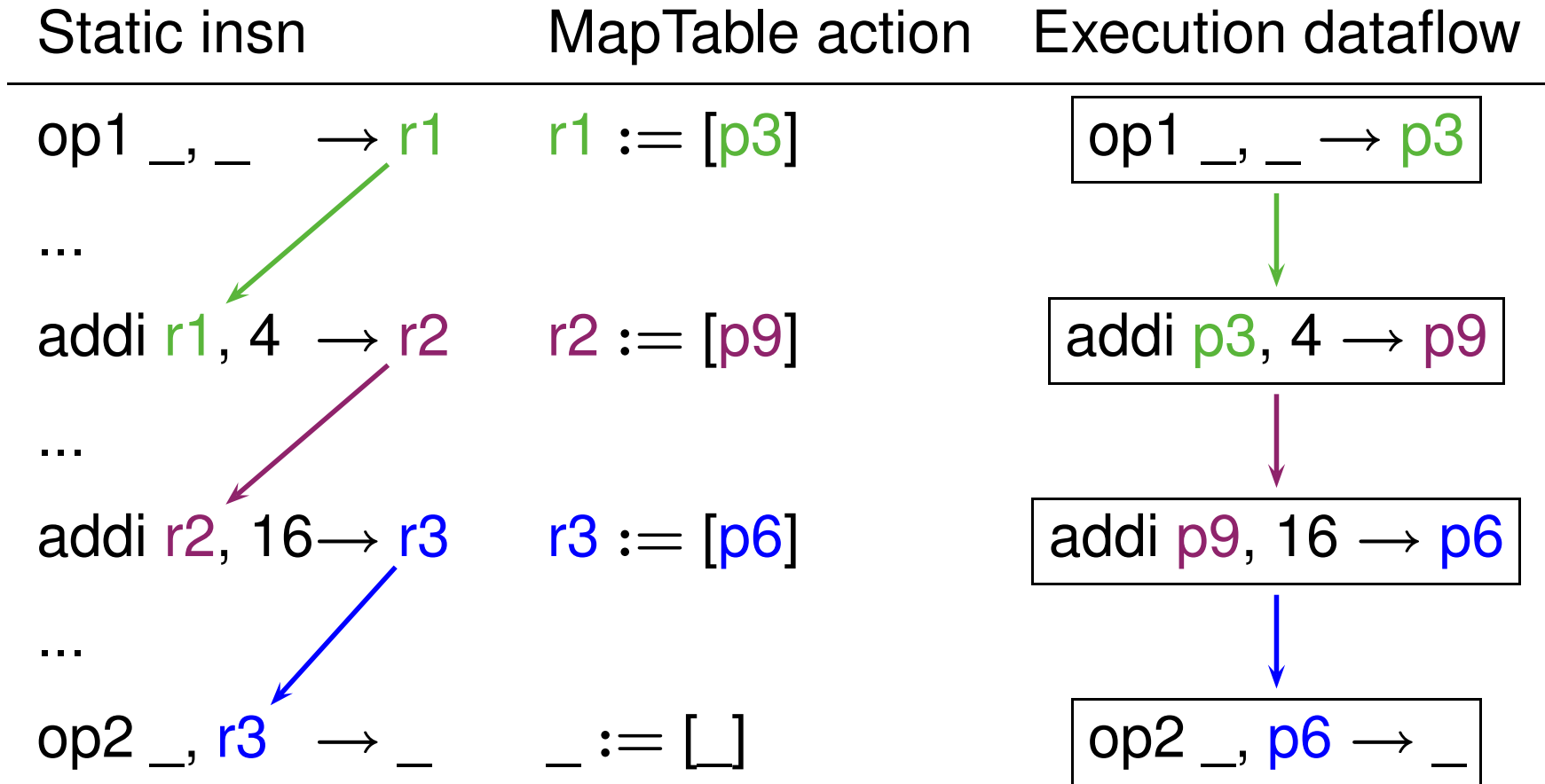
Conventional Processing



Conventional Processing



Conventional Processing



RENO_{CF}: Constant Folding

Static insn	MapTable action	Execution dataflow
op1 _, _ → r1		
...		
addi r1, 4 → r2		
...		
addi r2, 16 → r3		
...		
op2 _, r3 → _		

```
graph TD; I1[op1 _, _ → r1] -- green --> I2[addi r1, 4 → r2]; I2 -- purple --> I3[addi r2, 16 → r3]; I3 -- blue --> I4[op2 _, r3 → _];
```

RENO_{CF}: Constant Folding

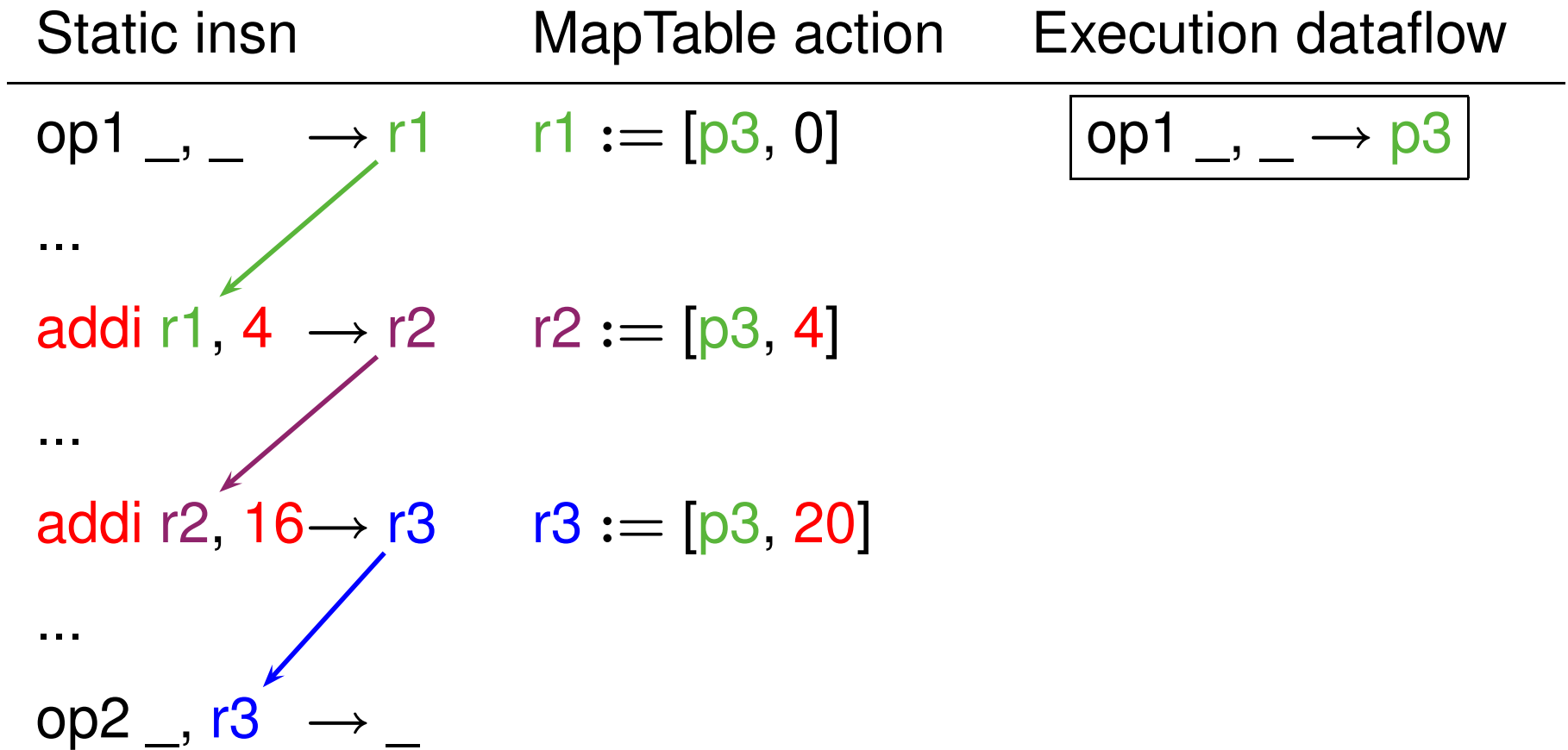
Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3, 0]	op1 <u>, <u> → p3</u></u>
...		
addi r1, 4 → r2		
...		
addi r2, 16 → r3		
...		
op2 <u>, r3 → <u></u></u>		

RENO_{CF}: Constant Folding

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3, 0]	op1 <u>, <u> → p3</u></u>
...		
addi r1, 4 → r2	r2 := [p3, 4]	
...		
addi r2, 16 → r3		
...		
op2 <u>, r3 → <u></u></u>		

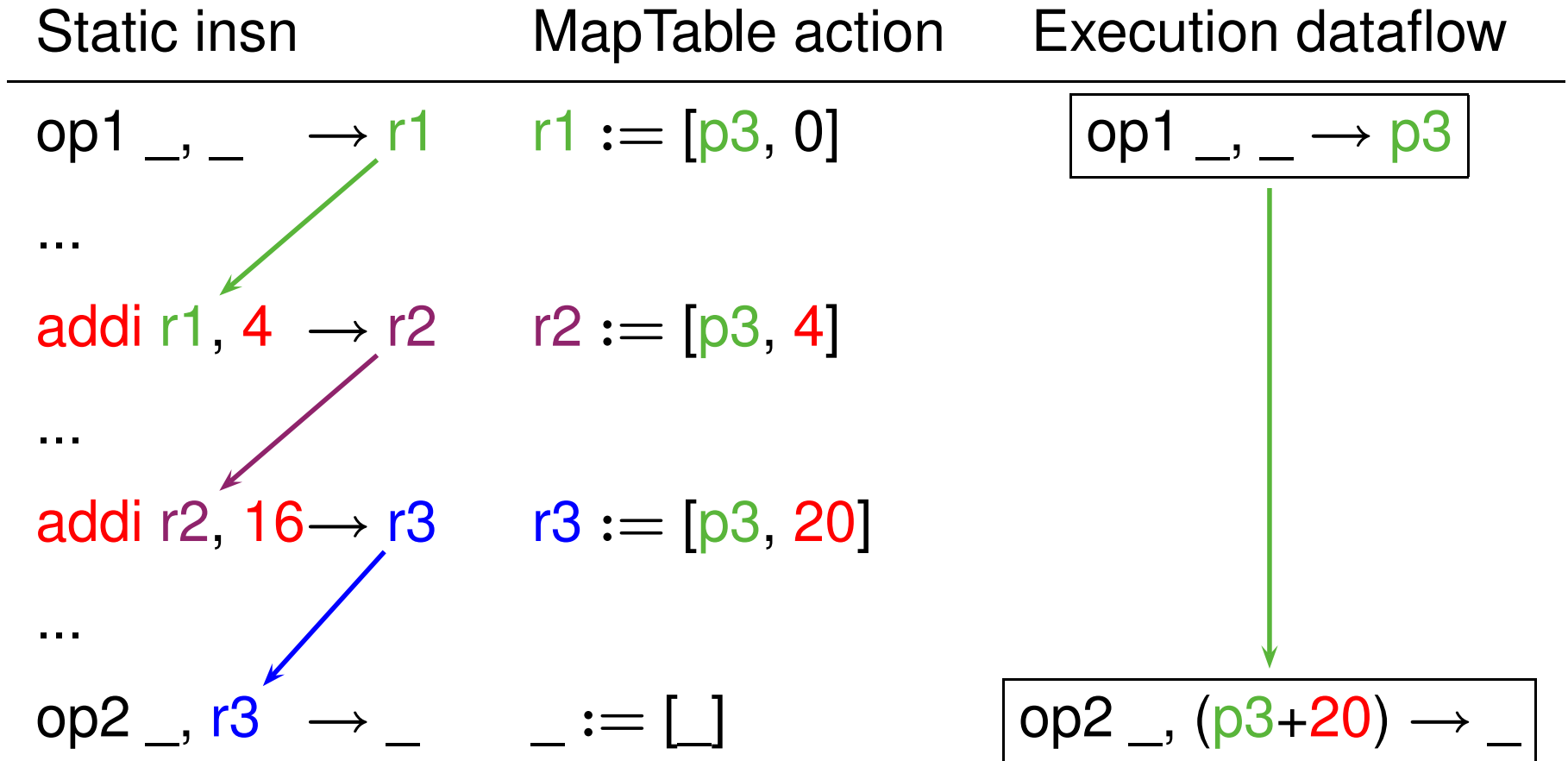
RENO_{CF}: Constant Folding

Static insn	MapTable action	Execution dataflow
op1 <u>, <u> → r1</u></u>	r1 := [p3, 0]	op1 <u>, <u> → p3</u></u>
...		
addi r1, 4 → r2	r2 := [p3, 4]	
...		
addi r2, 16 → r3	r3 := [p3, 20]	
...		
op2 <u>, r3 → <u></u></u>		



Renamer accumulates displacements

RENO_{CF}: Constant Folding



Renamer accumulates displacements

Consumer performs fused **addi-op2** operation

RENO Implementation

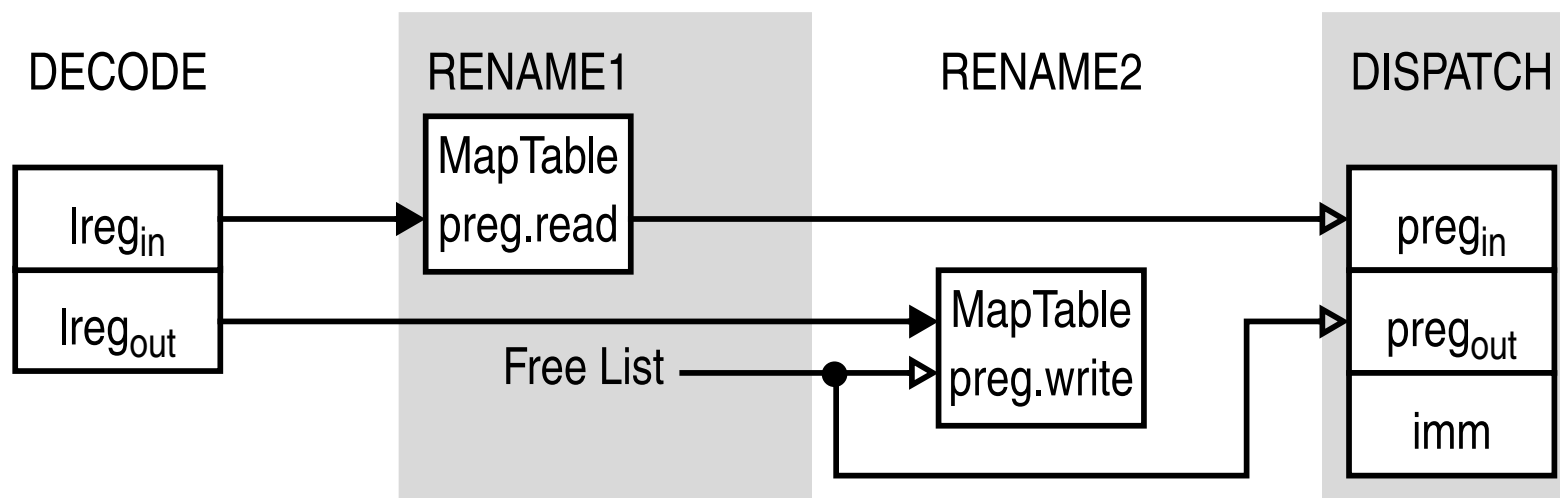
Rename

- $\text{RENO}_{ME/CSE/RA}$: physical register sharing
- RENO_{CF} : physical register sharing++

Execute

- $\text{RENO}_{ME/CSE/RA}$: nothing
- RENO_{CF} : fused “addi-X” functional units

Conventional Renaming

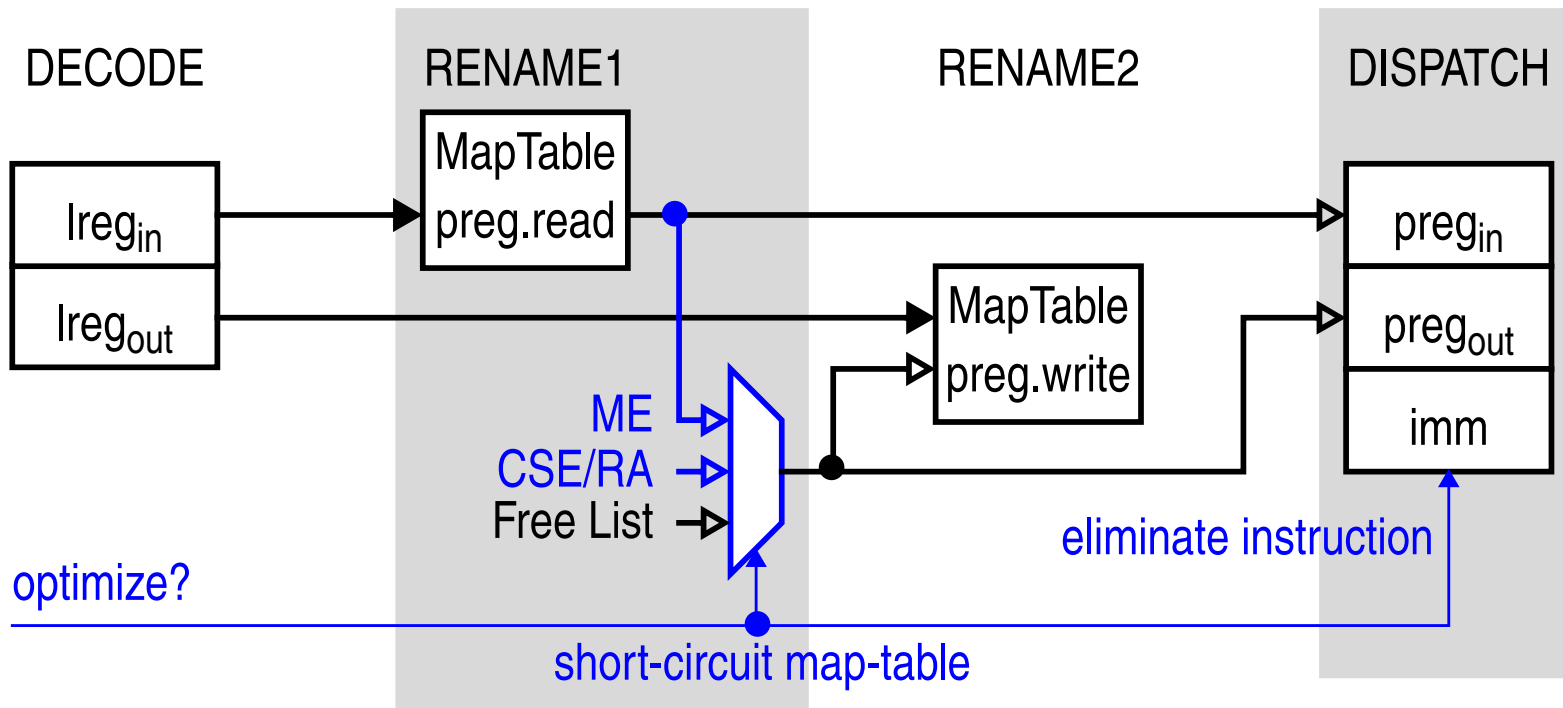


Example: scalar, one register input per instruction

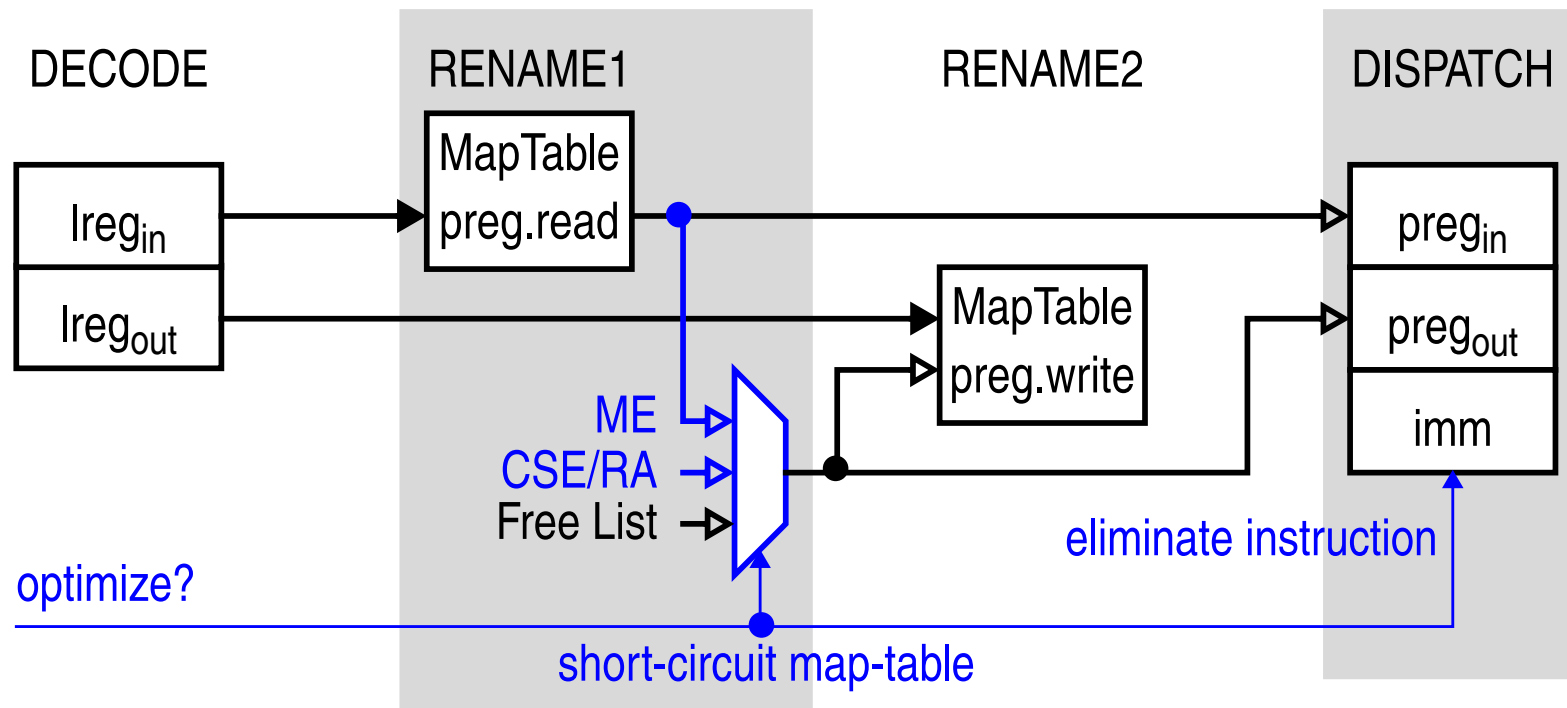
- See paper for superscalar



RENO_{ME/CSE/RA} Renaming



RENO_{ME/CSE/RA} Renaming

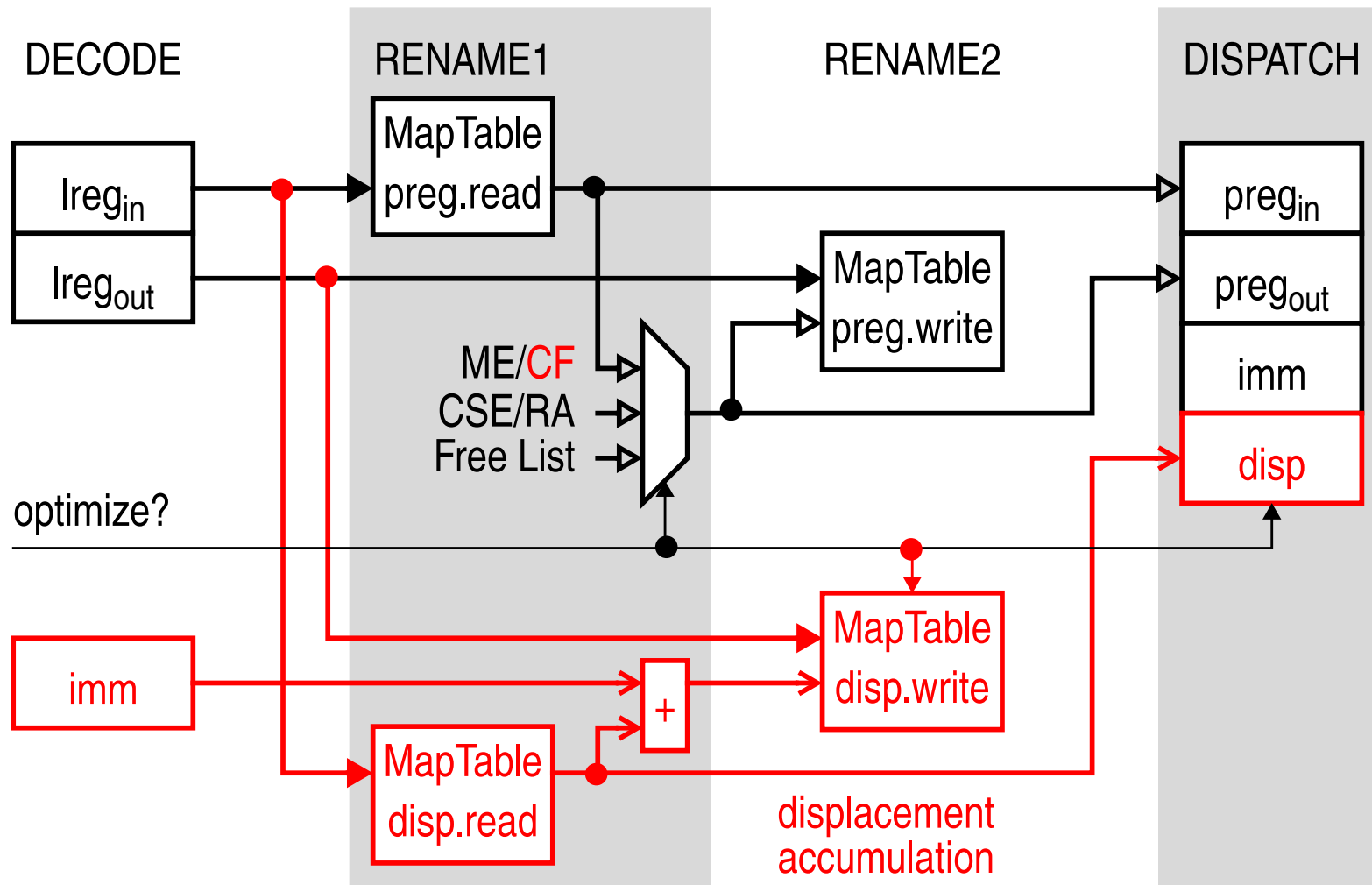


Superscalar? see paper

- $O(N^2)$: same as conventional renaming
- Don't optimize dependent insns in same cycle
- We *believe* additional renaming stages are unnecessary



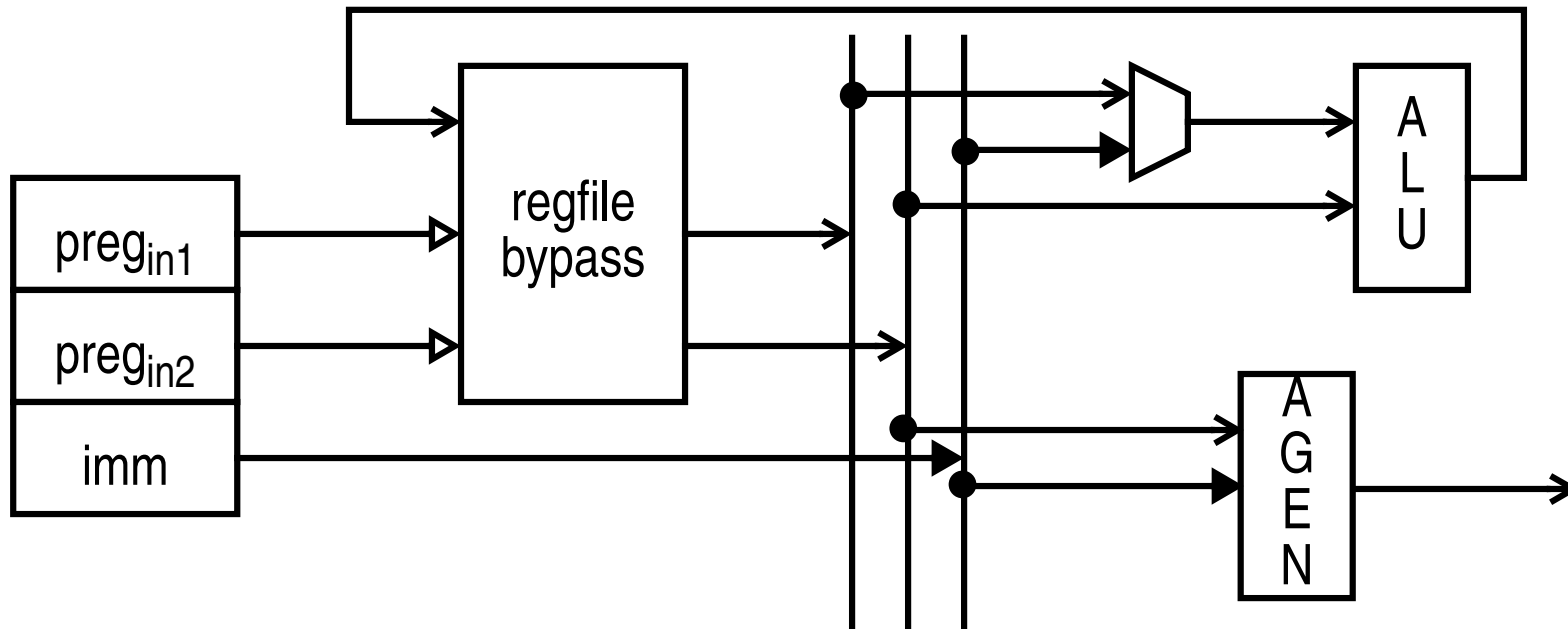
RENO_{CF} Renaming



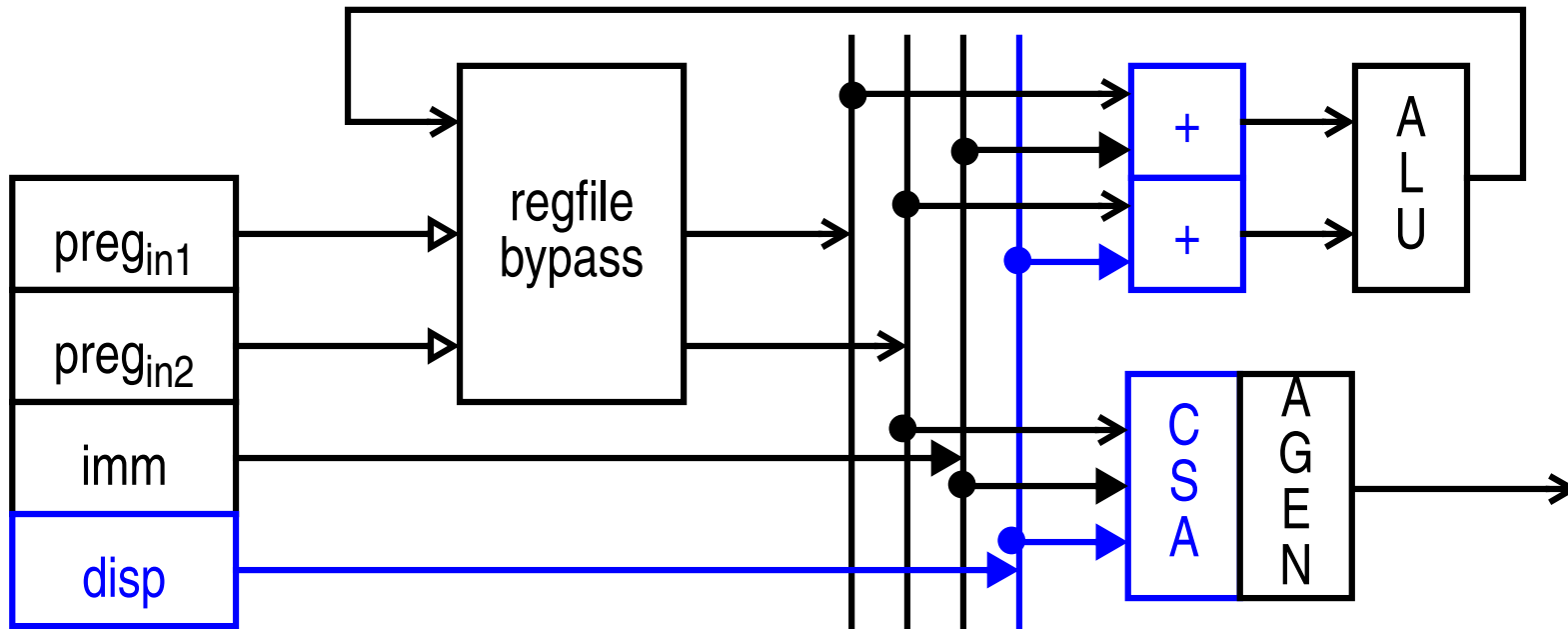
RENO_{ME} + parallel displacement accumulation circuit



Conventional Execution



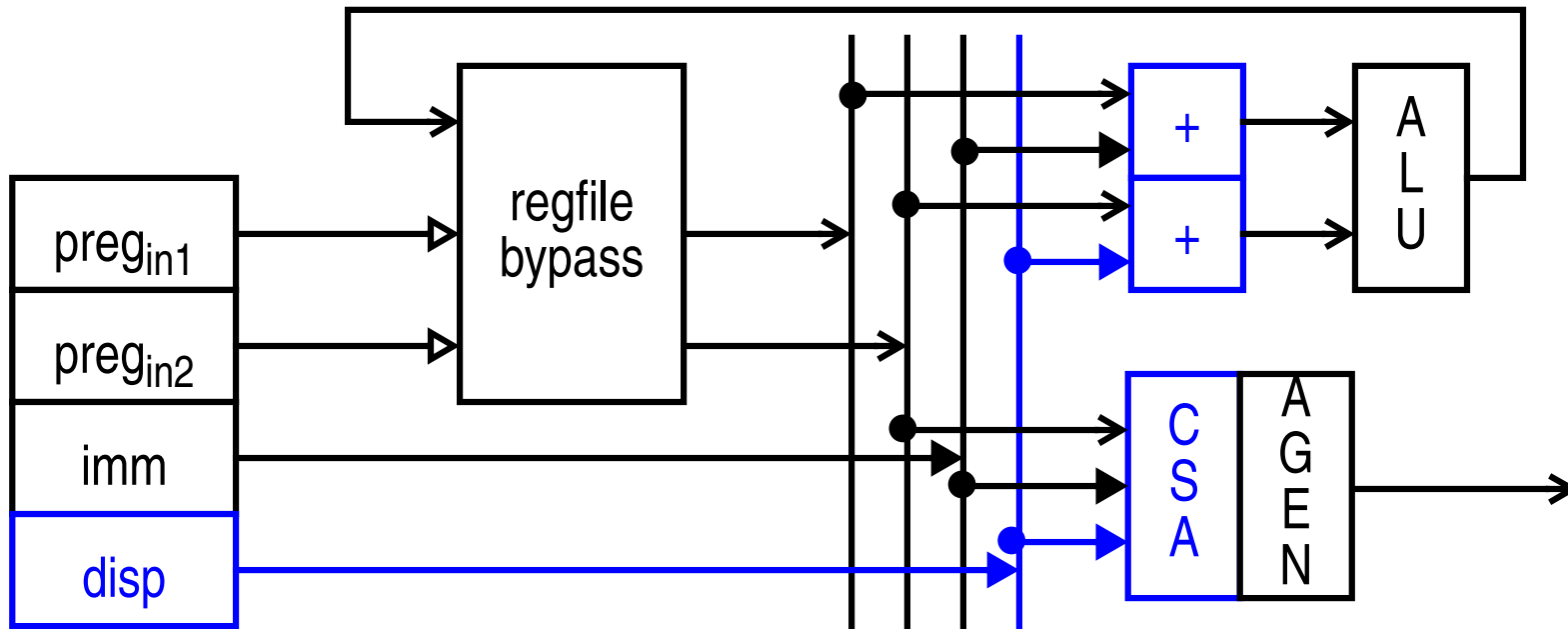
RENO_{CF} Execution



- All operations take an additional (immediate) input
- 1-cycle penalty for general “addi-X” fusion
 - + 0-cycle penalty “addi-addi” fusion (carry-save adder)
 - Most common and performance critical case (agen)



RENO_{CF} Execution



- All operations take an additional (immediate) input
- 1-cycle penalty for general “addi-X” fusion
 - + 0-cycle penalty “addi-addi” fusion (carry-save adder)
 - Most common and performance critical case (agen)

Regfile/bypass not complicated



Experimental Evaluation

What we care about

- Elimination rates (i.e., “coverage”) and performance
- RENO_{CF} synergy with $\text{RENO}_{CSE/RA}$
- See paper for sensitivity, bandwidth amplification, etc.

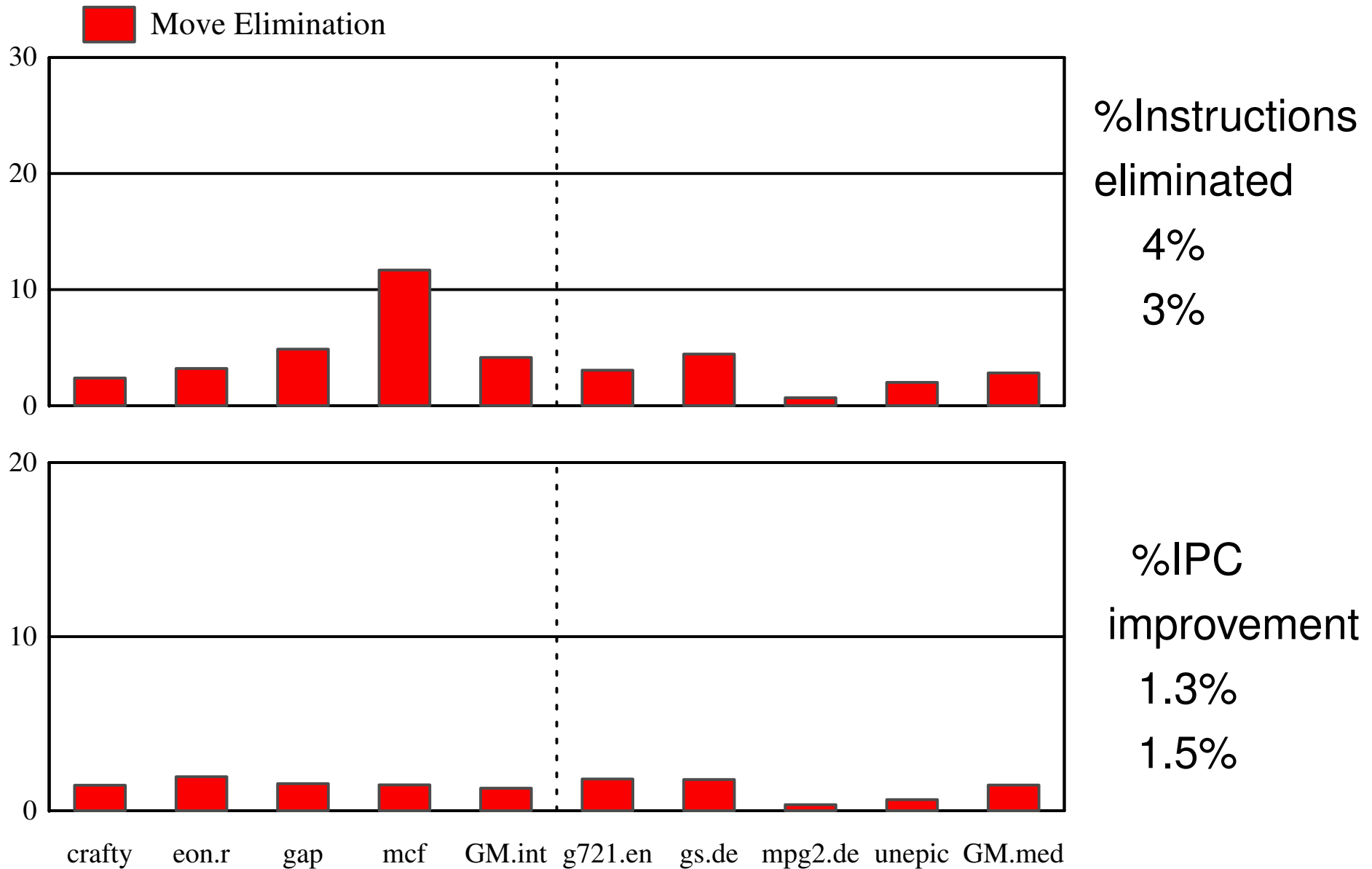
Benchmarks: SPECint2000, MediaBench

- DEC OSF **-O4**

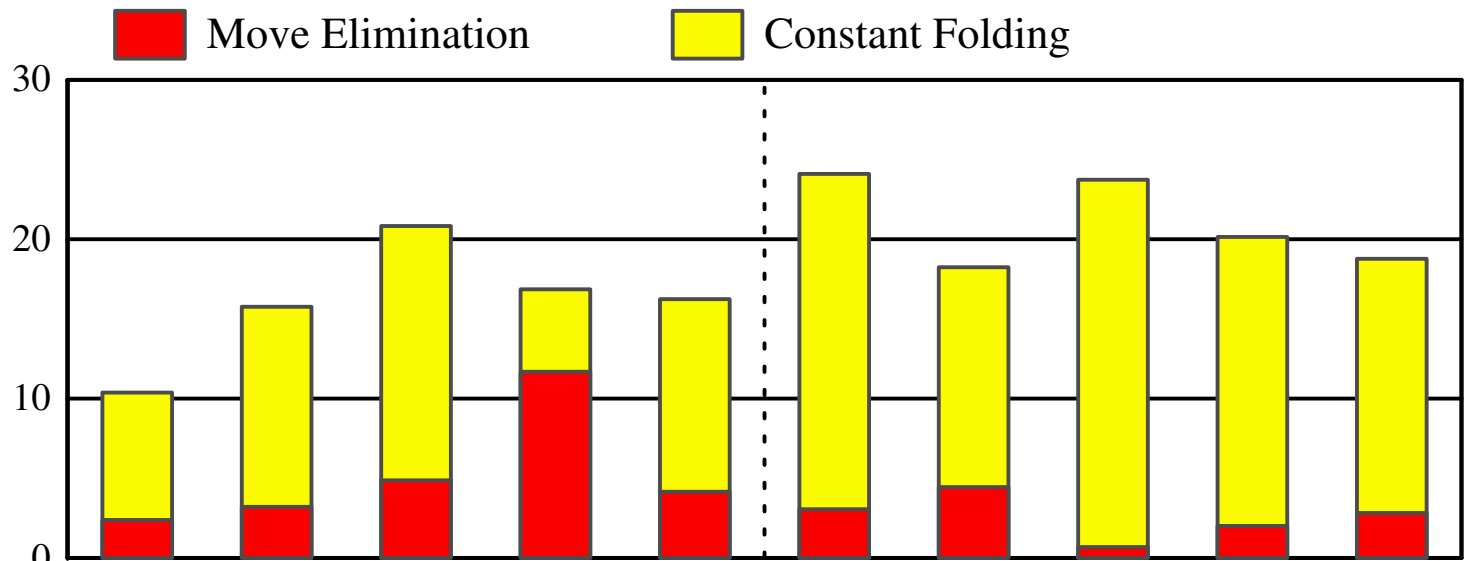
Simulator: Alpha SimpleScalar++

- 4-wide, out-of-order, 128-entry ROB, 50-entry issue queue
- 32KB D\$, 16KB I\$, 512KB L2
- RENO_{CF} : 16-bit displacements
- $\text{RENO}_{CSE/RA}$: 256-entry reuse table (**loads only**)

Coverage and Performance



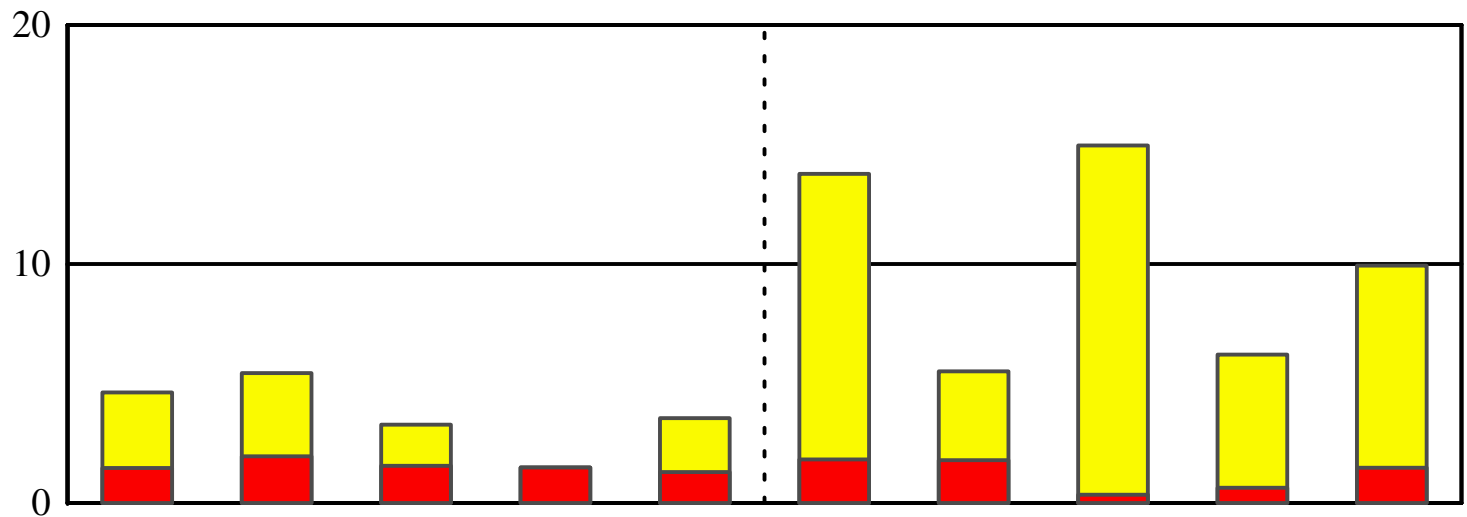
Coverage and Performance



%Instructions eliminated

16%

19%



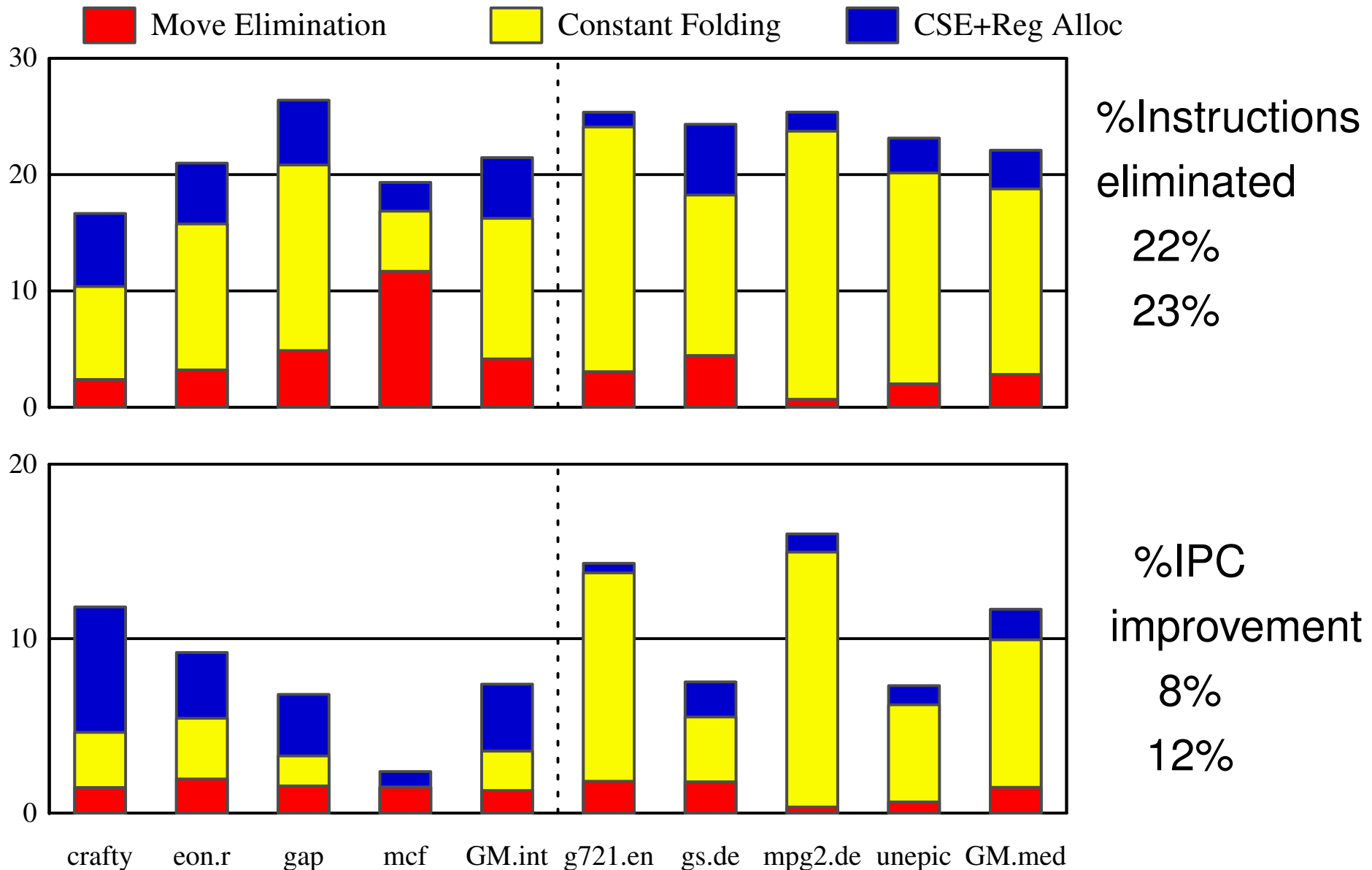
%IPC improvement

4%

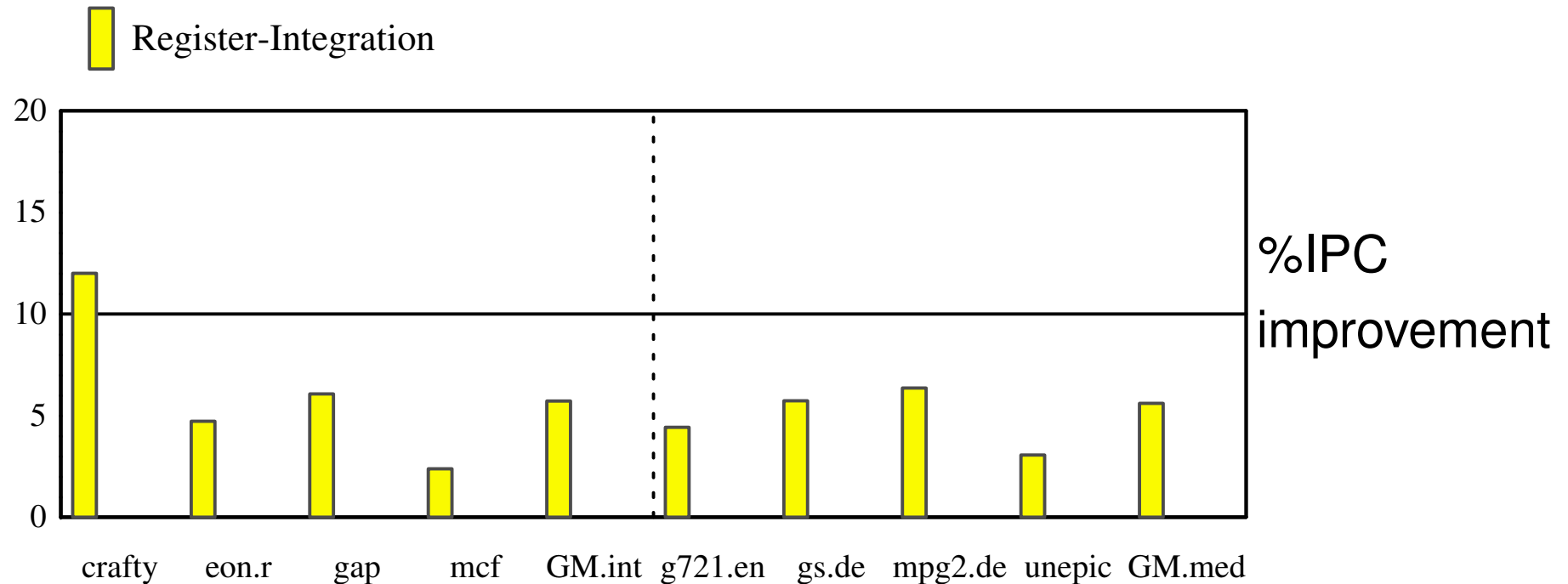
10%

crafty eon.r gap mcf GM.int g721.en gs.de mpg2.de unepic GM.med

Coverage and Performance



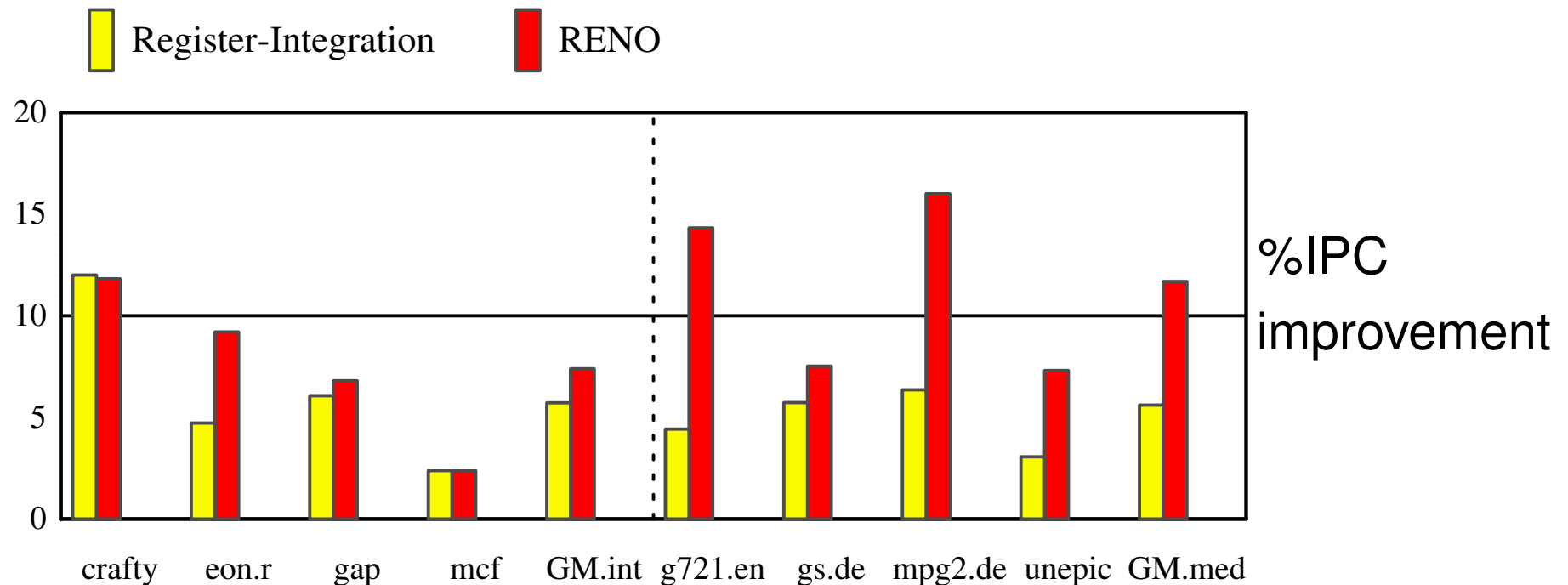
RENO_{CSE/RA} / RENO_{CF} Synergy



Register Integration: CSE for all, no folding

- Operations must be redundant
- Large tables, many lookups

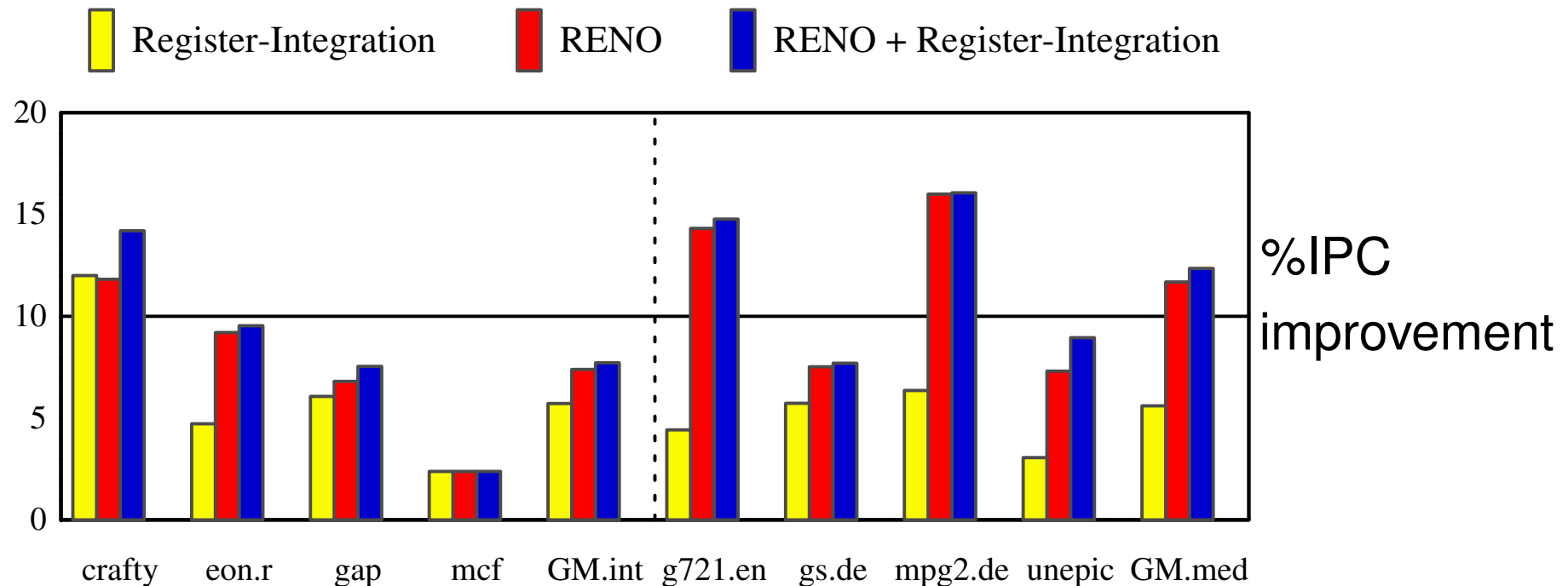
RENO_{CSE/RA} / RENO_{CF} Synergy



RENO: CSE for loads + folding

- + Better performance: addi's need not be redundant
- + Smaller tables, fewer lookups

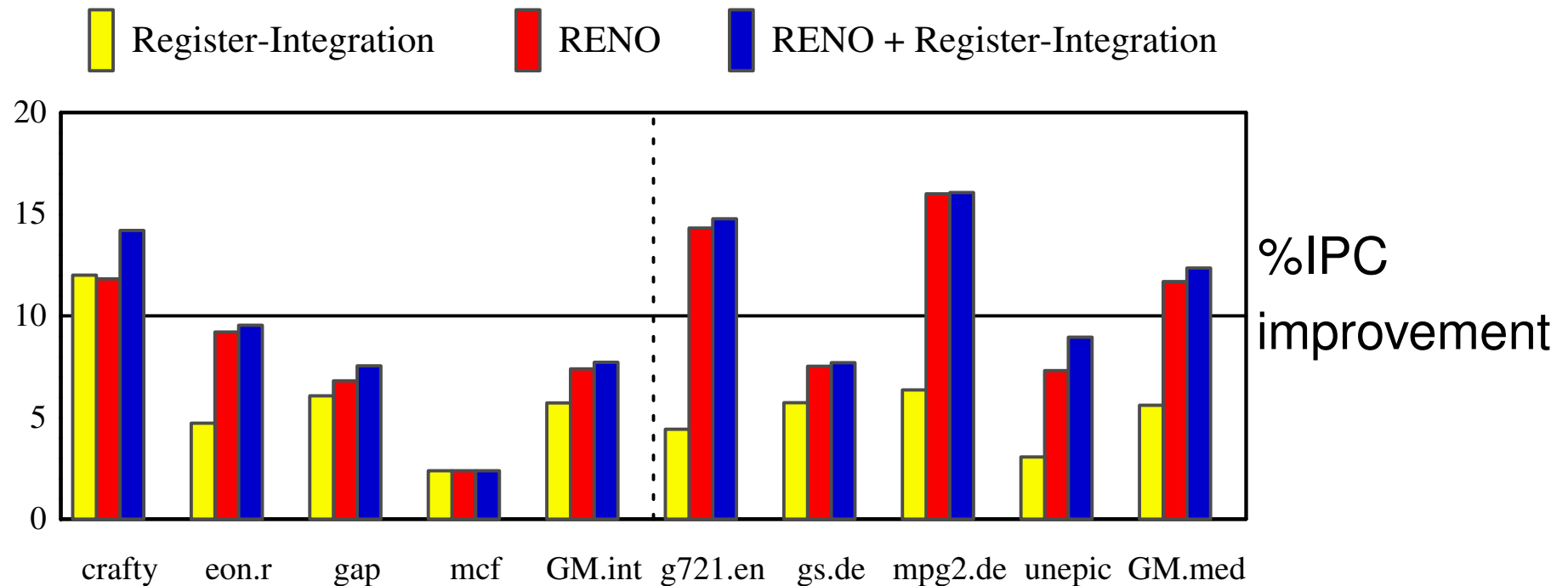
RENO_{CSE/RA} / RENO_{CF} Synergy



RENO+Register-Integration: CSE for all + folding

- Marginally better performance
- Not worth larger tables, more accesses

RENO_{CSE/RA} / RENO_{CF} Synergy



RENO sweetspot

folding much more important than non-load CSE

RENO Summary

Concept

- Dynamic analogs of compiler optimizations

Implementation

- Unifies several previously proposed techniques
 - Adds constant folding
- + Simple design

Effectiveness

- + Eliminates 23% of instructions from optimized programs

RENO and CO [Fahs+]

Similar ideas, different approaches

- Implementation complexity / optimization coverage tradeoff

RENO: purely register-name based

- Modified R10K renamer (adds immediates only)

- 3-input functional units

+ Simpler implementation

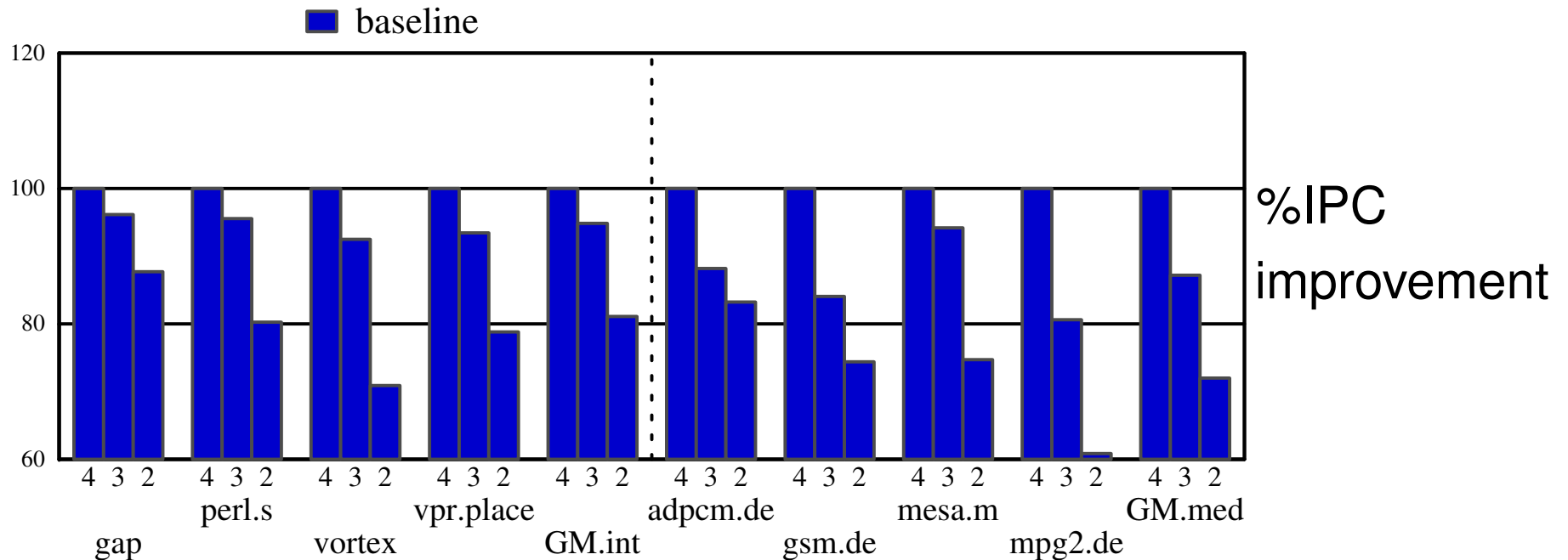
CO (our understanding): register-name+value based

- Short in-order execution pipeline inside renaming

- Feedback paths from main core register file

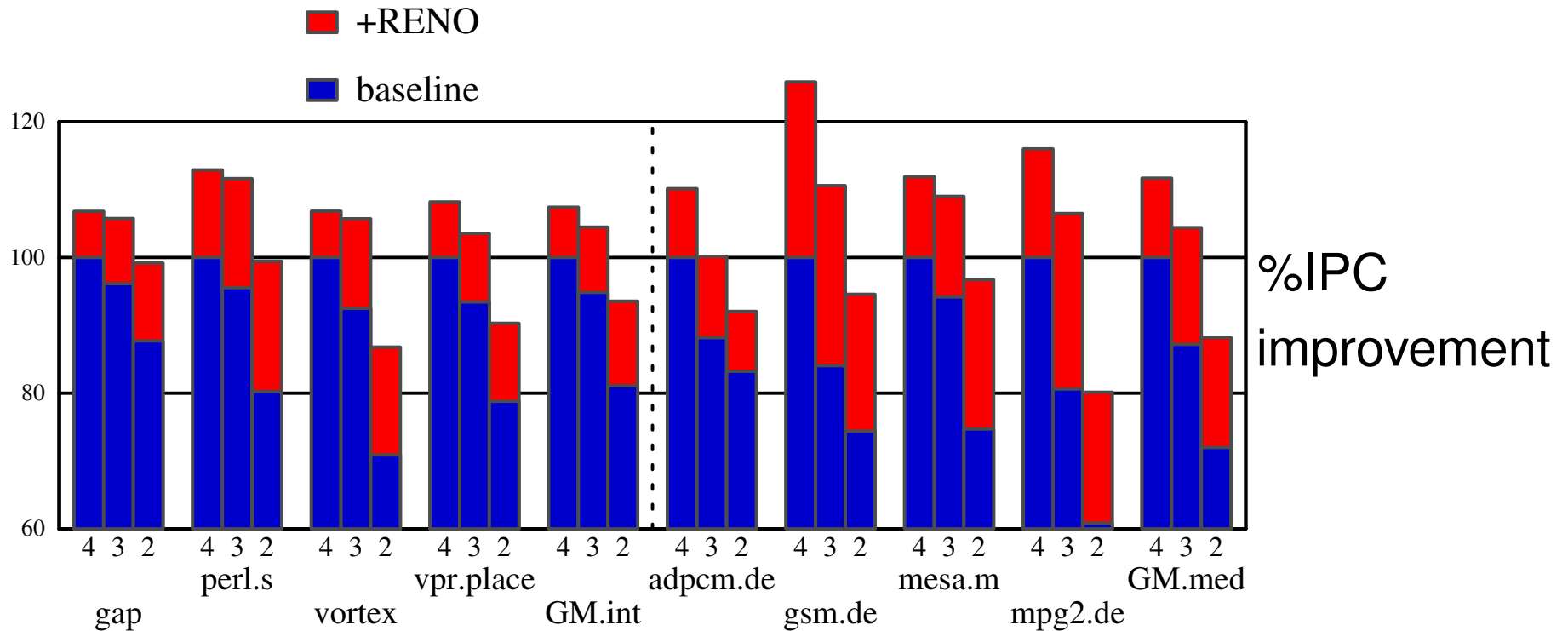
+ Broader optimization scope

Core Bandwidth Reduction



Scaled down execution core: 3-wide and 2-wide

Core Bandwidth Reduction



Scaled down execution core: 3-wide and 2-wide

3-wide: small speedup (4%) vs. 4-wide

2-wide: performance within 6% for SPEC and 12% for MediaBench vs 4-wide