

Just-In-Time Software Pipelining

Hongbo Rong Youfeng Wu Hyunchul Park Cheng Wang

Programming Systems Lab Intel Labs, Santa Clara



What is software pipelining? A loop optimization exposing instruction-level parallelism (ILP)

for
$$(i = 0; i < N; i++)$$
{
 $a: x = y + 1$
 $b: y = A[i] + x$
 $c: B[i+2] = B[i]*x$
 $d: A[i+2] = B[i+2]$
}



What is software pipelining? A loop optimization exposing instruction-level parallelism (ILP)

for
$$(i = 0; i < N; i + +)$$
{
 $a: x = y + 1$
 $b: y = A[i] + x$
 $c: B[i+2] = B[i] * x$
 $d: A[i+2] = B[i+2]$
}
Local dependence
Loop-carried dependence







































intel



intel



inte

Software pipelining has been static

- Extensively studied in 3 decades, and efficient for wide-issue architectures
 - VLIW [Lam 1988]
 - superscalar [Ruttenberg et al. 1996]



Software pipelining has been static

- Extensively studied in 3 decades, and efficient for wide-issue architectures
 - VLIW [Lam 1988]
 - superscalar [Ruttenberg et al. 1996]
- It is seen only in static compilers



Software pipelining has been static

- Extensively studied in 3 decades, and efficient for wide-issue architectures
 - VLIW [Lam 1988]
 - superscalar [Ruttenberg et al. 1996]
- It is seen only in static compilers
- Most works aim to minimize II but not compile overhead





- Dynamic languages are increasingly popular
 - JavaScript and PhP 88.9% and 81.5% in client and server websites (W3Techs)



- Dynamic languages are increasingly popular
 - JavaScript and PhP 88.9% and 81.5% in client and server websites (W3Techs)
- Huge amount of legacy code
 - Small optimization scope: a loop iteration
 - Software pipelining enlarges the scope to many iterations



- Dynamic languages are increasingly popular
 - JavaScript and PhP 88.9% and 81.5% in client and server websites (W3Techs)
- Huge amount of legacy code
 - Small optimization scope: a loop iteration
 - Software pipelining enlarges the scope to many iterations
- Minimizing compile overhead must be the 1st objective
 - Only simple/fast algorithms can be used



- Dynamic languages are increasingly popular
 - JavaScript and PhP 88.9% and 81.5% in client and server websites (W3Techs)
- Huge amount of legacy code
 - Small optimization scope: a loop iteration
 - Software pipelining enlarges the scope to many iterations
- Minimizing compile overhead must be the 1st objective
 - Only simple/fast algorithms can be used
 - linear-time algorithms are preferred



- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]



- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]



- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]





- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]



- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]



- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]
- Costly rollback
 - Software: Light-weight checkpointing



- Memory aliases kill parallelism
 - Hardware: Atomic region + rotating alias registers [MICRO-46]
- Costly rollback
 - Software: Light-weight checkpointing
- Scheduling is expensive for (j = 0; j < N; j + = M)for (i = j; i < j + M; i + +)for (i = 0; i < N; i + +)a Original b bAtomic optimizat С region С ion scope d 6

Framework (on Transmeta CMS)



x86 binary

Framework (on Transmeta CMS)





Framework (on Transmeta CMS)




















































Scheduling is expensive

• NP-complete problem to find an optimal schedule [Colland et al. 1996]



Scheduling is expensive

- NP-complete problem to find an optimal schedule [Colland et al. 1996]
- O(V³) at least, exponential at worst [Rau et al. 1992]

– V: number of operations



Scheduling is expensive

- NP-complete problem to find an optimal schedule [Colland et al. 1996]
- O(V³) at least, exponential at worst [Rau et al. 1992]

– V: number of operations

• Can we linearize software pipelining?





- Schedule in linear time
 - Use either simple or fast sub-algorithms
 - Avoid cubic or exponential complexity



- Schedule in linear time
 - Use either simple or fast sub-algorithms
 - Avoid cubic or exponential complexity
 - For iterative sub-algorithms, have a threshold: the maximum #iterations allowed



- Schedule in linear time
 - Use either simple or fast sub-algorithms
 - Avoid cubic or exponential complexity
 - For iterative sub-algorithms, have a threshold: the maximum #iterations allowed
 - The smaller the threshold, the less the compile overhead
 - Once exceeded, abort software pipelining



- Schedule in linear time
 - Use either simple or fast sub-algorithms
 - Avoid cubic or exponential complexity
 - For iterative sub-algorithms, have a threshold: the maximum #iterations allowed
 - The smaller the threshold, the less the compile overhead
 - Once exceeded, abort software pipelining
 - Key question: *How small can the threshold be?*



- Schedule in linear time
 - Use either simple or fast sub-algorithms
 - Avoid cubic or exponential complexity
 - For iterative sub-algorithms, have a threshold: the maximum #iterations allowed
 - The smaller the threshold, the less the compile overhead
 - Once exceeded, abort software pipelining
 - Key question: *How small can the threshold be?*
- Find a good enough schedule
 - No backtracking
 - Priority function: approximate and never update
 - Separate dependence and resource constraints
 - Separate local and loop-carried dependences



- Schedule in linear time
 - Use either simple or fast sub-algorithms
 - Avoid cubic or exponential complexity
 - For iterative sub-algorithms, have a threshold: the maximum #iterations allowed
 - The smaller the threshold, the less the compile overhead
 - Once exceeded, abort software pipelining
 - Key question: *How small can the threshold be?*
- Find a good enough schedule
 - No backtracking
 - Priority function: approximate and never update
 - Separate dependence and resource constraints
 - Separate local and loop-carried dependences
- Iteratively improve a schedule





Prepartition	
Н, В)

• Quickly creates an initial schedule to start with





- Quickly creates an initial schedule to start with
- Handles local dependences and resources





- Quickly creates an initial schedule to start with
- Handles local dependences and resources
- Adjusts the schedule for loop-carried dependences





- Quickly creates an initial schedule to start with
- Handles local dependences and resources
- Adjusts the schedule for loop-carried dependences





- Quickly creates an initial schedule to start with
- Handles local dependences and resources
- Adjusts the schedule for loop-carried dependences
- Iteratively improves the schedule





- Quickly creates an initial schedule to start with
- Handles local dependences and resources
- Adjusts the schedule for loop-carried dependences
- Iteratively improves the schedule
- Time complexity: O(V+E)
 - V: #operations E: #dependences









- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method



- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method

Conventional

O(V³)



- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method

Conventional

O(V³)

Turn the problem into a Markov decision process
– 1st time Howard algorithm is applied to pipelining



- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method

Conventional

Howard policy iteration algo.

O(V³) O(exponential*E)

- Turn the problem into a Markov decision process
 - $1^{\mbox{\scriptsize st}}$ time Howard algorithm is applied to pipelining



- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method



- Turn the problem into a Markov decision process
 1st time Howard algorithm is applied to pipelining
- Linearize Howard with a small constant ${\cal H}$



- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method



- Turn the problem into a Markov decision process
 1st time Howard algorithm is applied to pipelining
- Linearize Howard with a small constant ${\cal H}$


Calculating RecMII

- Recurrence Minimum II
 - II determined by the biggest dependence cycles
 - Needed by almost every software pipelining method



- Turn the problem into a Markov decision process
 1st time Howard algorithm is applied to pipelining
- Linearize Howard with a small constant ${\cal H}$
- A by-product: critical operations



Bellman-Ford

O(**∨***E)



Bellman-Ford

O(**V***E)

• Also an iterative sub-algorithm



Bellman-Ford

Linearized Bellman-Ford



O(<mark>B</mark>* E)

- Also an iterative sub-algorithm
- Linearize it with a constant B



Bellman-Ford

Linearized Bellman-Ford



O(<mark>B</mark>* E)

- Also an iterative sub-algorithm
- Linearize it with a constant B
- Specific order in visiting edges
 - Scan nodes in sequential order, and visit their incoming edges
 - Values are propagated along local edges in the $1^{\mbox{\scriptsize st}}$ itr.
 - Values are propagated along loop-carried edges in the 2nd itr.

Bellman-Ford Linearized Bellman-Ford



O(E)

- Also an iterative sub-algorithm
- Linearize it with a constant B
- Specific order in visiting edges
 - Scan nodes in sequential order, and visit their incoming edges
 - Values are propagated along local edges in the $1^{\mbox{\scriptsize st}}$ itr.
 - Values are propagated along loop-carried edges in the 2nd itr.









Local scheduling

- Any local scheduling algorithm can be used
 - E.g. list scheduling



Local scheduling

- Any local scheduling algorithm can be used
 E.g. list scheduling
- Weakness:
 - Loop-carried dependences may be violated



Local scheduling

- Any local scheduling algorithm can be used
 E.g. list scheduling
- Weakness:
 - Loop-carried dependences may be violated
- To reduce the chance of violation:
 - Before scheduling, priority function considers loopcarried dependences in advance
 - Prioritize critical operations















Experiments

• Transmeta CMS on SPEC2k traces



Experiments

- Transmeta CMS on SPEC2k traces
- Functional simulator to comprehensively
 - Explore thresholds H and B
 - Evaluate compile overhead and schedules' quality



Experiments

- Transmeta CMS on SPEC2k traces
- Functional simulator to comprehensively
 - Explore thresholds H and B
 - Evaluate compile overhead and schedules' quality
- Cycle-accurate simulator
 - Simulates cache misses, latencies, ...
 - Initial performance study





























• $B \le 3$ for 98.8% of 11,992 loops





• $B \le 3$ for 98.8% of 11,992 loops





- $B \le 3$ for 98.8% of 11,992 loops
- $B \le 5$ for all the loops





- $B \le 3$ for 98.8% of 11,992 loops
- $B \le 5$ for all the loops
- From now on, we set H=10, $B=5 \rightarrow 11,910$ loops scheduled



Scheduling overhead & schedules' quality





Scheduling overhead & schedules' quality



• JITSP achieves optimal schedules for 95% loops



Scheduling overhead & schedules' quality



• JITSP achieves optimal schedules for 95% loops


Scheduling overhead & schedules' quality



• JITSP achieves optimal schedules for 95% loops



Scheduling overhead & schedules' quality



• JITSP achieves optimal schedules for 95% loops



Scheduling overhead & schedules' quality



• JITSP achieves optimal schedules for 95% loops



Compile overhead distribution



Note: acyclic scheduler handles acyclic code, or loops NOT selected for software pipelining



• 40 hot loops













- 40 hot loops
- The architecture has bottleneck in registers
 - 2/7/24/32 predicate/static alias/integer/floating point available for pipelining 26





- 40 hot loops
- The architecture has bottleneck in registers
 - 2/7/24/32 predicate/static alias/integer/floating point available for pipelining 26





- 40 hot loops
- The architecture has bottleneck in registers
 - 2/7/24/32 predicate/static alias/integer/floating point available for pipelining 26



II/MII (The lower, the better)





II/MII (The lower, the better)



• JITSP achieves optimal schedules for all but 1 loops



Speedup





Speedup



• JITSP 10~36% speedup. Better than the others

 Exception: loop 2. Optimal schedule but slowdown due to memory aliases → retranslation needed



Speedup



• JITSP 10~36% speedup. Better than the others

- Exception: loop 2. Optimal schedule but slowdown due to memory aliases → retranslation needed
- Speedup swim(5.3%), ammp(4.4%), mcf(3.6%)



• 1st linear software pipelining algorithm implemented for dynamic compilers



- 1st linear software pipelining algorithm implemented for dynamic compilers
- Turns a traditionally-expensive optimization into linear time O(V+E)



- 1st linear software pipelining algorithm implemented for dynamic compilers
- Turns a traditionally-expensive optimization into linear time O(V+E)
 - Taking advantages of results from various domains:
 - hardware circuit design (Retiming)
 - stochastic control (Howard algorithm)
 - Graph (Bellman-Ford)
 - software pipelining (Rotation scheduling and DESP)



- 1st linear software pipelining algorithm implemented for dynamic compilers
- Turns a traditionally-expensive optimization into linear time O(V+E)
 - Taking advantages of results from various domains:
 - hardware circuit design (Retiming)
 - stochastic control (Howard algorithm)
 - Graph (Bellman-Ford)
 - software pipelining (Rotation scheduling and DESP)
- Generates optimal or near-optimal schedules with reasonable compile overhead



Future work

- Register availability
 - Add more architecture registers
 - Algorithm: Register pressure-aware
- Implementation
 - Loop selection
 - Re-translation
- Evaluation
 - Benchmarks



Backup slides



- Each operation is rewarded to reach a cycle via 1 policy edge
 - The bigger the cycle, the more the reward



- Each operation is rewarded to reach a cycle via 1 policy edge
 - The bigger the cycle, the more the reward





- Each operation is rewarded to reach a cycle via 1 policy edge
 - The bigger the cycle, the more the reward





- Each operation is rewarded to reach a cycle via 1 policy edge
 - The bigger the cycle, the more the reward





Distribution statistics of JITSP

21	min	median	mean	max
# operations	4	10	13.82	96
# dependences	12	38	50.48	353
# local dependences	1	13	22.36	283
# loop-carried deps	5	22	28.12	156
MII	1	3	4.78	55
II - MII	0	0	0.05	5
II / MII	1	1	1.01	1.5
# local scheduling	1	1	1.26	4

