Clairvoyance: Look-Ahead Compile-Time Scheduling

Kim-Anh Tran, Trevor E. Carlson, Konstantinos Koukos, Magnus Själander, Vasileios Spiliopoulos, Stefanos Kaxiras, Alexandra Jimborean

Presenters: Lee Moore, Nicholas Wilson, Jiong Zhu, Do June Min

Memory latency in memory-bound applications

```
for (i=0; i<N; i++) {
   t1 = load x[i]
   if(t1) {
     t2 = load node[i]->val
     t3 = load y[i]
      store t^2 + t^3,
            node[i]->val
```

- Chains of dependent loads and stores
- Last-level cache
 misses (LLC) can cost
 hundreds of cycles

How can we hide memory latency?

- Hardware-level Out-of-Order (OoO) Execution Engine
 - Aggressive OoO Engine needs aggressive energy





Fast and power-hungry OoO

Simple but energy-efficient OoO

Could compile-time scheduling help in hiding memory latency?

How can we hide memory latency?

- Compile-time scheduling

Challenges to handle:

- Insufficient independent instructions
- Chains of dependent loads
- Statically unknown dependencies
- Register pressure

```
for (i=0; i<N; i++) {
  t1 = load x[i]
  if(t1) {
     t2 = load node[i]->val
     t3 = load y[i]
     store t^2 + t^3,
            node[i]->val
```

Introducing Clairvoyance

- Goal: hide memory latency for simple OoO engine in compile time.
- How to approach this goal?



Basic Clairvoyance Algorithm

- (2) Unroll Loop
- (3) Initialize Access
- (4) Identify load instructions
- (5-9) Identify instructions required by load instructions
- (10-11) Create access and execute phases
- (12-13) Join access and execute phases

Input: Loop L, Unroll Count $count_{unroll}$ Output: Clairvoyance Loop $L_{Clairvoyance}$

	0	
2	$L_{unrolled} \leftarrow \texttt{Unroll}(L, \operatorname{count}_{unroll})$	
3	$L_{access} \leftarrow Copy(L_{unrolled})$	
4	$hoist_list \leftarrow FindLoads(L_{access})$	
5	$to_keep \gets \emptyset$	
6	for load <i>in</i> hoist_list do	
7	$requirements \leftarrow FindRequirements(load)$	
8	to_keep \leftarrow Union(to_keep, requirements)	
9	end	
10	$L_{access} \leftarrow \texttt{RemoveUnlisted}(L_{access}, \texttt{to}_{\texttt{keep}})$	
11	$L_{execute} \leftarrow \texttt{ReplaceListed}(L_{access}, L_{unrolled})$	
12	$L_{Clairvoyance} \leftarrow \texttt{Combine}(L_{access}, L_{unrolled})$	
13	return L _{Clairvoyance}	
14 end		

Example Transformation



Improvement to Basic Clairvoyance

- (2) Unroll Loop
- (3) Initialize Access
- (4) Identify load instructions
- (5-9) Identify instructions required by load instructions
- (10-11) Create access and execute phases
- (12-13) Join access and execute phases

Input: Loop L, Unroll Count $count_{unroll}$ Output: Clairvoyance Loop $L_{Clairvoyance}$

2	$L_{unrolled} \leftarrow \texttt{Unroll}(L, \operatorname{count_{unroll}})$	
3	$L_{access} \leftarrow \texttt{Copy}(L_{unrolled})$	
4	$hoist \ list \leftarrow FindLoads(L_{access})$	
5	$to_keep \gets \emptyset$	
6	for load in hoist_list do	
7	$requirements \leftarrow FindRequirements(load)$	
8	to_keep \leftarrow Union(to_keep, requirements)	
9	end	
10	$L_{access} \gets \texttt{RemoveUnlisted}(L_{access}, \texttt{to}_\texttt{keep})$	
11	$L_{execute} \gets \texttt{ReplaceListed}(L_{access}, L_{unrolled})$	
12	$L_{Clairvoyance} \gets \texttt{Combine}(L_{access}, L_{unrolled})$	
13	return L _{Clairvoyance}	
14 end		

Identifying Critical Loads

- "critical loads" are load instructions selected for the access phase
- Indirection Count:
 - number of loads a load instruction is dependent on
 - ex: x[y[z[i]]] = 2
- Load instructions with indirection count <= count_{indir} are selected as critical loads

$$count_{Indir} = 0$$

$$count_{Indir} = 1$$

$$coun$$

for t₁

t₁

Improvement to Basic Clairvoyance

- (2) Unroll Loop
- (3) Initialize Access
- (4) Identify load instructions
- (5-9) Identify instructions required by load instructions
- (10-11) Create access and execute phases
- (12-13) Join access and execute phases

Input: Loop L, Unroll Count $count_{unroll}$ Output: Clairvoyance Loop $L_{Clairvoyance}$

2	$L_{unrolled} \leftarrow \texttt{Unroll}(L, \operatorname{count_{unroll}})$	
3	$L_{access} \leftarrow \mathtt{Copy}(L_{unrolled})$	
4	hoist list \leftarrow FindLoads(L _{access})	
5	$to_keep \gets \emptyset$	
6	for load in hoist_list do	
7	$requirements \leftarrow FindRequirements(load)$	
8	$to_keep \leftarrow Union(to_keep, requirements)$	
9	end	
10	$L_{access} \gets \texttt{RemoveUnlisted}(L_{access}, \texttt{to}_\texttt{keep})$	
11	$L_{execute} \gets \texttt{ReplaceListed}(L_{access}, L_{unrolled})$	
12	$L_{Clairvoyance} \leftarrow \texttt{Combine}(L_{access}, L_{unrolled})$	
13	return L _{Clairvoyance}	
14 end		

Handling Unknown Dependencies

• Load should be hoisted only if data dependency is respected.

- Aliases have 3 cases
 - a. No-Alias \Rightarrow Hoist!
 - b. Must-Alias \Rightarrow Don't hoist!
 - c. May-Alias \Rightarrow Prefetch!



Improvement to Basic Clairvoyance

- (2) Unroll Loop
- (3) Initialize Access
- (4) Identify load instructions
- (5-9) Identify instructions required by load instructions
- (10-11) Create access and execute phases
- (12-13) Join access and execute phases

Input: Loop L, Unroll Count $count_{unroll}$ Output: Clairvoyance Loop $L_{Clairvoyance}$

2	$L_{unrolled} \leftarrow \texttt{Unroll}(L, \operatorname{count_{unroll}})$	
3	$L_{access} \leftarrow \texttt{Copy}(L_{unrolled})$	
4	$hoist_list \leftarrow FindLoads(L_{access})$	
5	$to_keep \gets \emptyset$	
6	for load <i>in</i> hoist_list do	
7	requirements \leftarrow FindRequirements (load)	
8	to_keep ← Union(to_keep, requirements)	
9	end	
10	$L_{access} \leftarrow \texttt{RemoveUnlisted}(L_{access}, \texttt{to}_\texttt{keep})$	
11	$L_{execute} \leftarrow \texttt{ReplaceListed}(L_{access}, L_{unrolled})$	
12	$L_{Clairvoyance} \leftarrow \texttt{Combine}(L_{access}, L_{unrolled})$	
13	return L _{Clairvoyance}	
14 end		

Handling Chains of Dependent Loads

Input: Set of *loads* Output: List of sets *phase_loads*



Parameters of the Clairvoyance Pass

- Turn on Clairvoyance Optimization:
 - Heuristics: Disable if loads/branches < 0.7

- Loop Unroll Count, Indirection Count:
 - Choose by using state-of-the-art runtime version selectors

Benchmark Results

Memory-bound benchmarks

- Yield highest performance improvements
- Expose most opportunities for MLP

DAF

Clairvoyance-Best shows potential improvements using speculative heuristics, and given better pointer analysis.

Compute-bound benchmarks expose less opportunities for Clairvoyance improvements



Conclusion

Idea is increasing instruction-level, memory-level parallelism to make better use of core resources.

Key technique is reordering dependent-loads across loop iterations, controlling register pressure.

Most improvement in memory-bound scenarios: up to 43%, 7% geomean for conservative, 13% geomean for speculative

Thoughts

- + Implementation software only
- + Effective use of software prefetching
- + Open-source, possible to replicate
- Not enough detail on register pressure
- Effect of indirection count parameter on performance of tests not clear
- Test hardware not aligned with paper motivation

Questions?