Constraint-Based Register Allocation and Instruction Scheduling

Candidate:Roberto Castañeda LozanoOpponent:Prof. Laurent MichelMain Supervisor:Prof. Christian Schulte



Doctoral defense, September 3, 2018

This research has been partially funded by Ericsson AB and the Swedish Research Council (VR 621-2011-6229)



- The integration of register allocation and instruction scheduling using constraint programming is practical and effective.
 - practical: complete, fairly scalable, executableeffective: better code than heuristic approaches
 - Enables compiler back-ends that can:
 - ✓ trade compilation time for speed or code size
 - ✓ adapt to new processor features and criteria

Earlier Approaches

- Incomplete model
 - mostly local (single basic blocks)
 - limited coverage of standard subproblems
- Do not scale beyond small problems (up to ~ 100 ins)
- No executable code
- Ours:
 - complete model
 - fair scalability (up to 647 ins)
 - executable code
 - key: constraint programming

Constraint Programming



Key feature: global constraints

Our Approach



```
    processor instructions
    temps (program variables)
    unscheduled
    t<sub>2</sub> ← mul t<sub>1</sub>, 2
t<sub>3</sub> ← sub t<sub>1</sub>, 1
t<sub>4</sub> ← add t<sub>2</sub>, t<sub>3</sub>
```

ret t_4

processor instructions
registers
scheduled
0: R0 ← mul R0,2 || R1 ← sub R0,1
1: R0 ← add R0,R1
2: ret R0

Contributions

C1 Survey of combinatorial approaches



Survey on Combinatorial Register Allocation and Instruction Scheduling. To appear in ACM CSUR, 2018

Contributions

- C2-3 Model and program representation
 - C4 Model extensions
 - C5 Solving method for scalability
 - C6 Experimental evaluation
 - C7 Study of model accuracy

Controlmband Register Advention and Instruction Workship
Annual Contract Contract (Section 2014) and Contract Cont
Where $M_{\rm eff}=M_{\rm eff}=M_{\rm eff}$ is the second secon
4 Evaluation Department of the start of the factor of the start of the Particle start of the start of the start of the start of the particle start of the start of the start of the start of the department of the start of the

Constraint-based Register Allocation and Instruction Scheduling. In CP, 2012 Combinatorial Spill Code Optimization and Ultimate Coalescing. In LCTES, 2014



Combinatorial Register Allocation and Instruction Scheduling. Tech. report, 2018 (submitted for publication)

Contributions

C8 Unison: a mature, open-source tool



 Image: Control of the state of the

Register Allocation and Instruction Scheduling in Unison. In CC, 2016 unison-code.github.io

Outline











Register Assignment

- Spilling, Splitting, and Coalescing
- Instruction Scheduling
- Global Scope
- Putting It All Together

2 Solver



4 Future Work

Register Assignment

Map temporaries to individual registers

$$\begin{array}{c} t_{2} \leftarrow \text{mul } t_{1}, 2 \\ t_{3} \leftarrow \text{sub } t_{1}, 1 \\ t_{4} \leftarrow \text{add } t_{2}, t_{3} \\ \text{ret } t_{4} \end{array} \xrightarrow[t_{1} \mapsto \text{R1}]{} \begin{array}{c} \text{R2} \leftarrow \text{mul } \text{R1}, 2 \\ \text{R1} \leftarrow \text{sub } \text{R1}, 1 \\ \text{R1} \leftarrow \text{add } \text{R2}, \text{R1} \\ \text{ret } \text{R1} \\ \text{ret } \text{R1} \end{array}$$

- Problem: limited number of registers
- Reuse registers as much as possible
- When can two temps reuse a register?

Register Assignment: Liveness

Live range of a temp t: points at which t might be needed

between definition and last use

$$t_2 \leftarrow \text{mul} \cdots$$

$$\cdots \leftarrow \text{sub} \cdots$$

$$\cdots \leftarrow \text{add} t_2, \cdots$$

 $ret \cdots$

Register Assignment: Register Reuse

When can two temps reuse a register?

if their live ranges do not interfere



t_3 and t_4 can use the same register

t₂ and t₃ cannot

Register Assignment: Model

to which register is each temp assigned?



Based on abstraction by Pereira and Palsberg (PLDI 2008)



Register Assignment

Spilling, Splitting, and Coalescing

- Instruction Scheduling
- Global Scope
- Putting It All Together

2 Solver



4 Future Work

Spilling

Allocate temporaries to memory

applied when no more registers are available



Live Range Splitting

Split a temporary live range into shorter ones can improve register reuse



Coalescing

Merge split temporaries

to reduce overhead of move instructions



Spilling, Splitting, and Coalescing: Model

Common mechanism: insert or remove copies

- Spilling: insert reg-to-mem and mem-to-reg copies
- Splitting: insert reg-to-reg copies
- Coalescing: remove reg-to-reg copies
- Model memory locations as registers (columns)
- 2 Extend input program with optional copies:

$$t_2 \leftarrow \{\texttt{store}, \texttt{move}, \bot\} \ t_1$$





Other Register Allocation Subproblems

- Register packing
- Load-store optimization
- Rematerialization
- Multi-allocation

1 Model

- Register Assignment
- Spilling, Splitting, and Coalescing
- Instruction Scheduling
- Global Scope
- Putting It All Together

2 Solver



4 Future Work

Instruction Scheduling: Model

- Map instructions to time points (issue cycles)
- Classic constraint-based scheduling model:

in which cycle is each instruction issued?

- Subject to:
 - precedences
 - resource constraints
- Connection to register allocation: live ranges



1 Model

- Register Assignment
- Spilling, Splitting, and Coalescing
- Instruction Scheduling

Global Scope

Putting It All Together

2 Solver



4 Future Work

Global Scope

- Extend the rectangle packing model to entire functions
- Temps live in multiple basic blocks?



- Decompose them into local temps
- Define decomposed temps as congruent (=)

Global Scope: Model

- Only relation between basic blocks: congruences
- Congruent temps are assigned to the same register



1 Model

- Register Assignment
- Spilling, Splitting, and Coalescing
- Instruction Scheduling
- Global Scope
- Putting It All Together

2 Solver



4 Future Work

Model

```
minimize \sum weight(b) \times cost(b) subject to
            no-overlap ({(reg(t), reg(t) + width(t), start(t), end(t)): t \in T_b \land live(t)}) \forall b \in B
                                                                                     reg[temp(p)] = r \quad \forall p \in P : p \triangleright r
                                                   reg[temp(p)] \in class[p, ins(operation(p))] \forall p \in P : active(operation(p))
                                                                                               active(o) \forall o \in O : \neg copv(o)
                                          live(t) \iff active(operation(definer(t))) \iff
                                          \exists p \in users(t) : active(operation(p)) \land temp(p) = t \quad \forall t \in T
                                                                    \operatorname{reg[temp(p)]} = \operatorname{reg[temp(q)]} \quad \forall p, q \in P : p \equiv q
                  issue(operation(q)) \ge issue(operation(p)) + lat[ins(operation(p))] \quad \forall t \in T, \forall p \in \{definer(t)\}, \forall q \in users(t) :
                                                                                                              active(operation(q)) \land temp(q) = t
cumulative ({(issue(o), dur[ins(o), s], con[ins(o), s]) : o \in O_b \land active(o)}, cap(s))
                                                                                                              \forall b \in B, \forall s \in S
                                                        start(t) = issue(operation(definer(t)))
                                                                                                              \forall t \in T : live(t)
                                                              end(t) = max issue(operation(p)) \forall t \in T: live(t)
                                                                   p \in users(t) : temp(p) = t
```

- Generic objective function: speed, code size, ...
- Compact: size independent of registers and cycles
- Global constraints:
 - rectangle packing, scheduling with resources





3 Results

4 Future Work

Solver: Decomposition

- Split the problem into simpler subproblems
- Solve each subproblem, recombine, and iterate
- Improves scalability



Solver: Decomposition

- Split the problem into simpler subproblems
- Solve each subproblem, recombine, and iterate
- Improves scalability



Solver: Other Improvements



"Many a little makes a mickle." (English proverb)









Experimental Setup

- How practical is our approach?
 - evaluation of scalability
- How effective?
 - evaluation of code quality (speed and code size)
- 100 functions from MediaBench and SPEC CPU2006
 - medium size: up to 1000 instructions
- Baseline: LLVM 3.8 (state-of-the-art compiler)
 - for comparison
 - as a quality lower bound
- Processors:
 - Hexagon V4: multiple-issue, signal processing
 - ARMv6: single-issue, general-purpose
 - MIPS32: single-issue, embedded
- Solver time limit: 15 minutes per function







4 Future Work

Solving Time to Optimality by Function Size



- Scales up to medium-size functions
 - optimally up to 647 ins, improves up to 874 ins
- Scalability depends on processor







4 Future Work

Estimated Speedup over LLVM: ARM



Slight speedup: 1% mean

Better spilling and coalescing

Estimated Speedup over LLVM: MIPS



- Moderate speedup: 5.1% mean
- Better spill code placement

Estimated Speedup over LLVM: Hexagon



Significant speedup: 9.3% mean

- Better exploitation of multiple-issue capabilities
 - more instructions per cycle
- Better spill code placement

Actual Application Speedup over LLVM

- Processor: Hexagon (highest estimated speedup)
- Benchmarks: MediaBench apps (top 5 functions)
- Execution: Qualcomm's Hexagon Simulator (high accuracy)



Moderate speedup: 5.4% mean, accuracy?

Speedup Accuracy: Estimated vs. Actual Speedup



- Strong monotonic relationship (Spearman correlation ρ = 0.77)
 Inaccuracies (Wilcoxon signed-rank test)
 - mostly due to dynamic processor behavior

Code Size Minimization

- Slight to moderate code size reduction over LLVM (0.8%-3.9% mean depending on the processor)
- Aggressive coalescing
 - scheduling and spilling subordinated
- Exploitation of processor-specific features
 - example: ARM's 16-bits Thumb-2 instructions



2 Solver

3 Results

4 Future Work

Future Work

- Extend instruction scheduling
- Increase model accuracy
- Improve scalability
- Integrate with instruction selection





3 Results

4 Future Work

Conclusion

The integration of register allocation and instruction scheduling using constraint programming is practical and effective.

practical: complete, fairly scalable, executable

- finally catching up with heuristic approaches
- used in industry today
- effective: better code than heuristic approaches
 - works today: actual speedup and code size reduction
 - works tomorrow:
 - exploits ever-evolving processor features
 - adapts to different optimization criteria

```
unison-code.github.io
```