

JIT-Compiled Frameworks: Compiler/Hardware Co-optimization

Wonsun Ahn

Assistant Professor

Department of Computer Science

University of Pittsburgh

<http://cs.pitt.edu/~wahn/>



Pitt

CS2001, October 2016

What is **Just-In-Time (JIT)** Compilation?

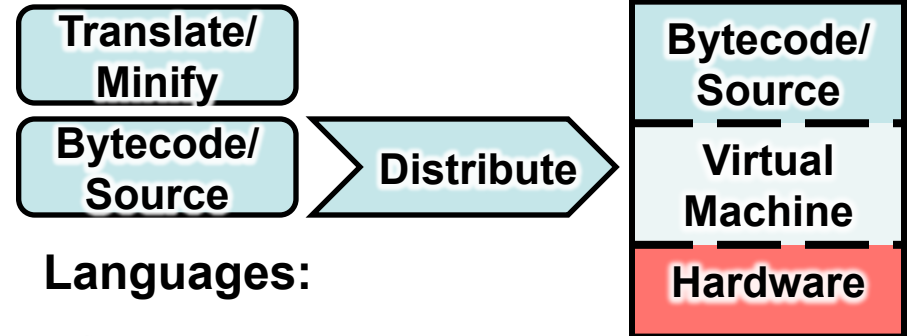
Ahead-Of-Time (AOT) Compilation



Languages:



Just-In-Time (JIT) Compilation



Languages:



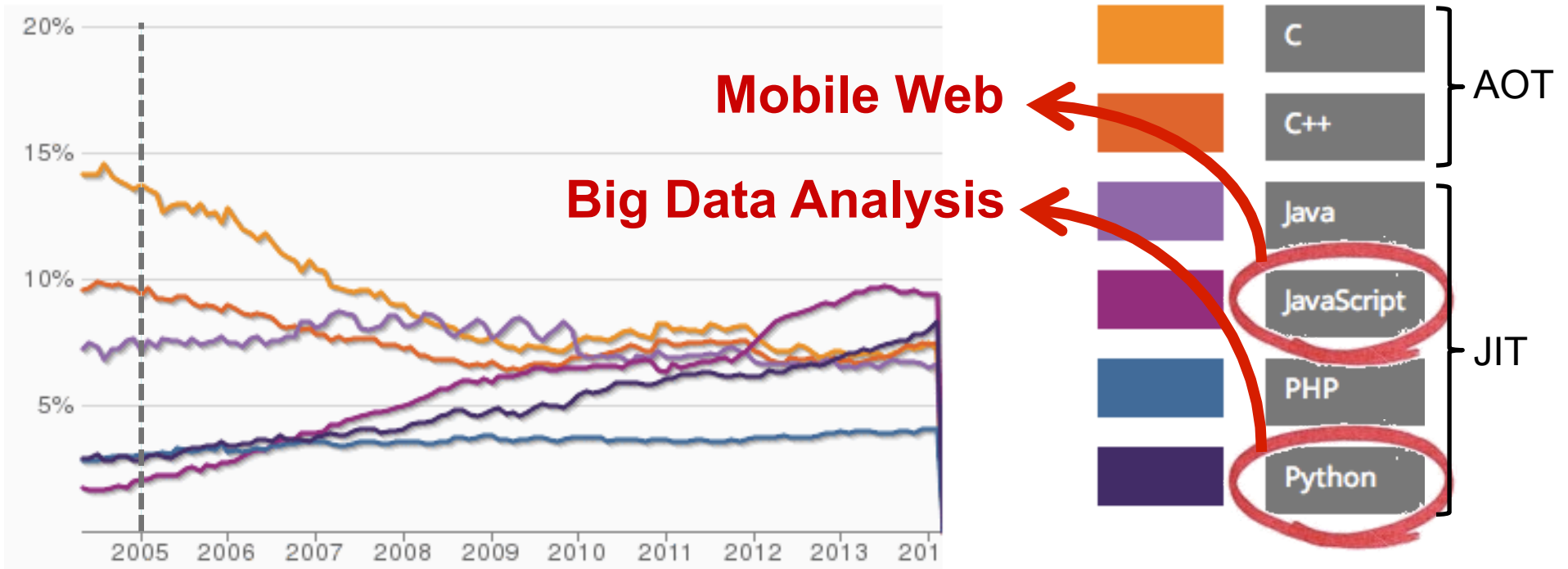
OpenCL



Pitt

Growth of JIT-Compiled / Scripting Languages

* Number of programmers coding in given language over time
(Reference: www.ohloh.net)



- C/C++ dominated in 2005
- Now, JavaScript and Python are the most used languages



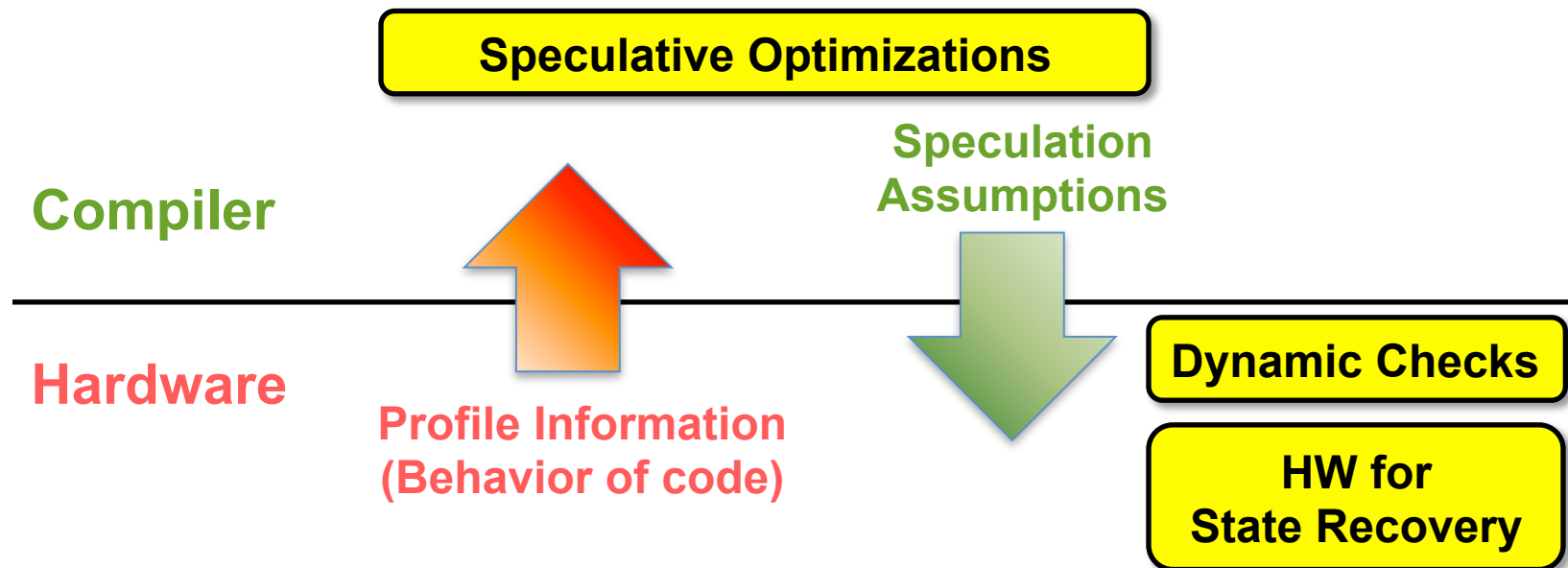
Will JIT / Scripting be Popular with IoT?

- Reasons IoT devices will want to use JIT:
 - JIT languages are **easier to use**
 - Big data scientists use Python, R, Matlab
 - JIT languages are **portable**
 - JavaScript is the only language that runs on all devices
 - JIT languages are more **secure**
 - Executes in a safe managed environment
- However many IoT devices still use C
 - Many IoT devices are battery-powered or energy-harvesting
 - C is more energy / memory efficient than JIT languages
 - Challenge: how to make JIT / Scripting efficient

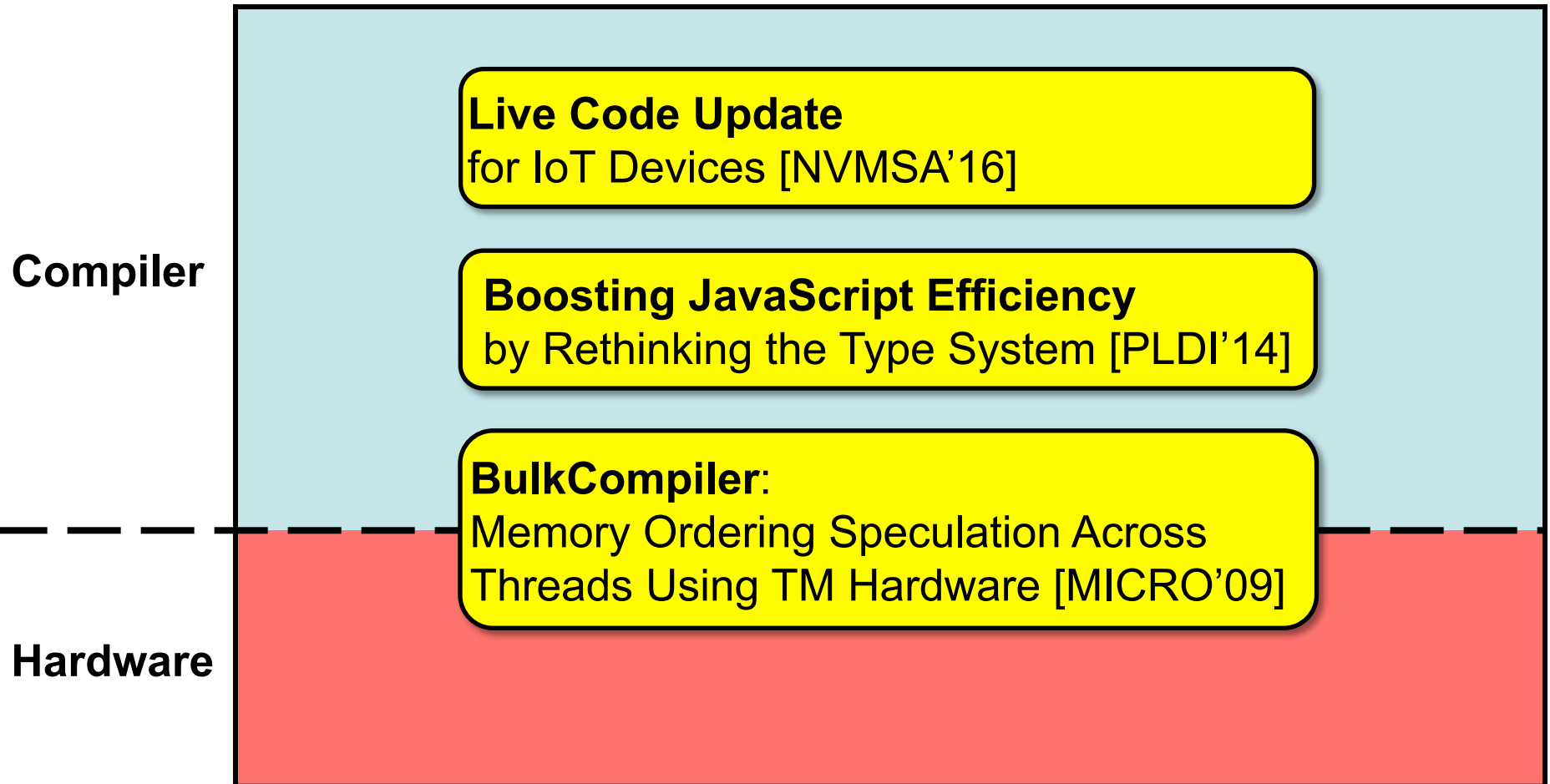


My Answer: Hardware/Compiler Co-Optimization

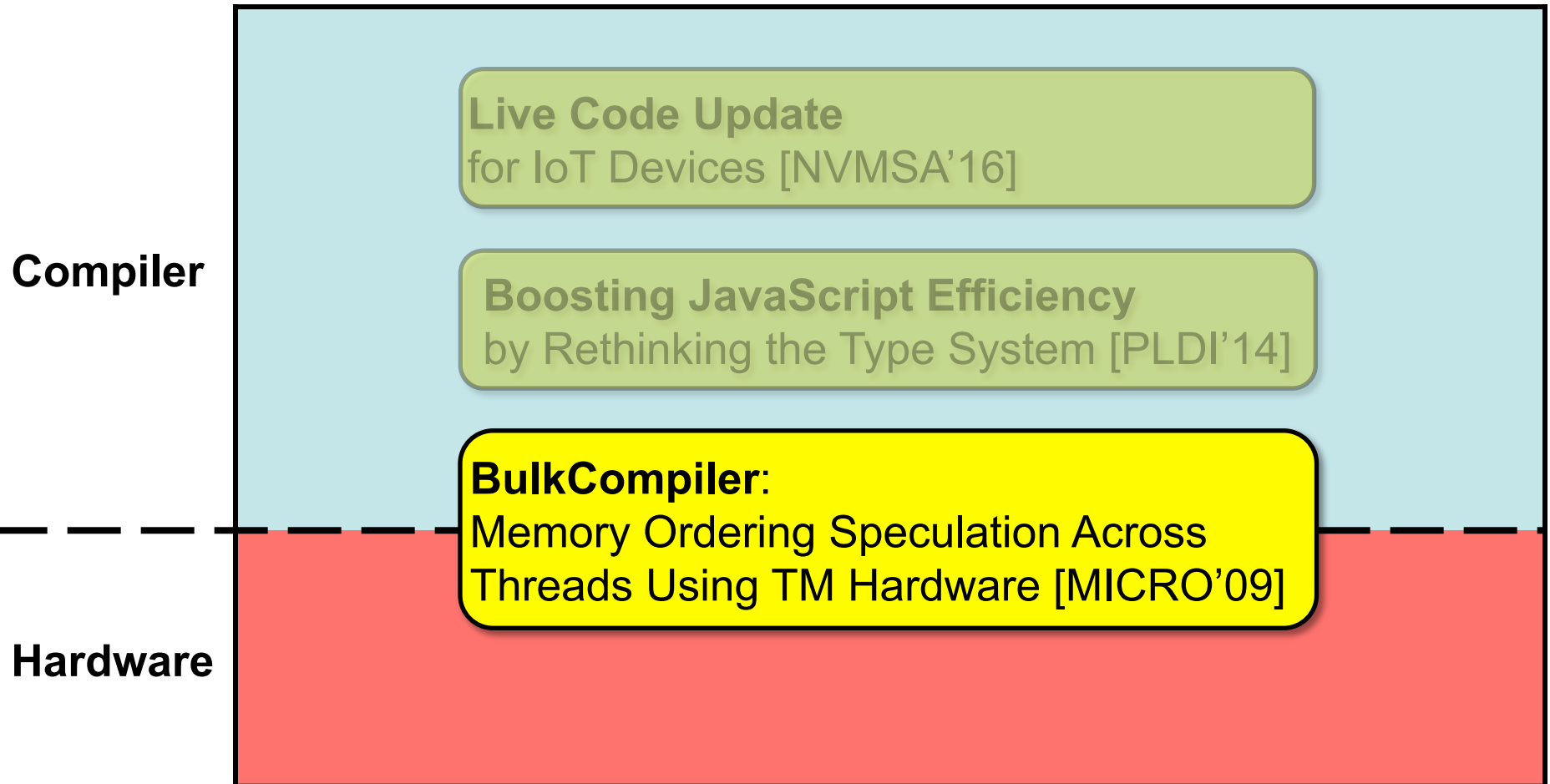
- One advantage JIT compilers have: **profile information**
- But often limited by:
 - Cost of collecting profile information
 - Cost of checking if profile-based prediction is correct
 - Cost of recovering state when prediction is incorrect



My Contributions

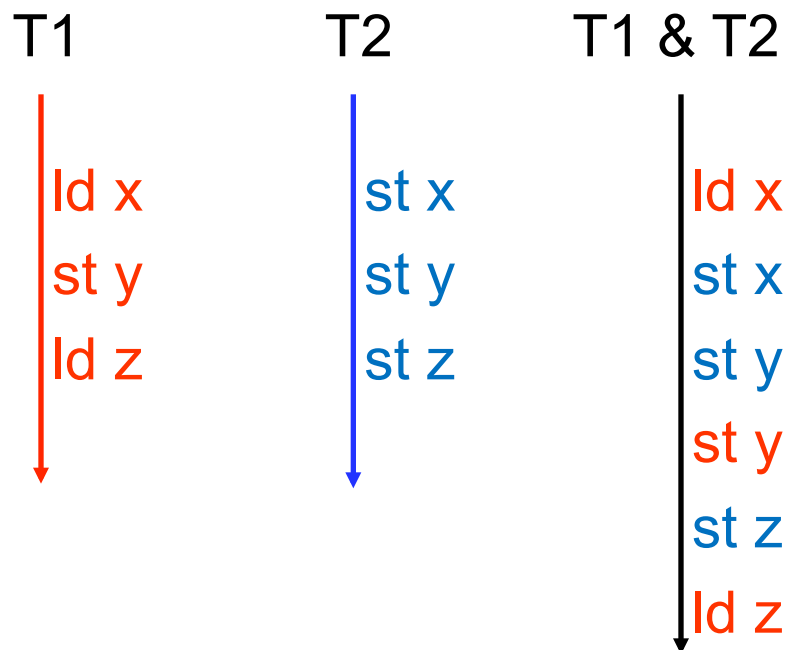


My Contributions



Memory Consistency (Ordering) Model

- Memory Consistency Model:
 - Defines how memory accesses can interleave across threads
 - Sequential Consistency (SC): only model that is intuitive to humans



- Accesses must appear to:
 - Interleave in a single total order
 - Follow per-thread program order



SC: Ideal Language Memory Consistency Model

- Much easier to **debug**:
 - Only memory model programmers can reason about
 - Debuggers can reproduce buggy interleaving
- **Software verification** tools assume SC:
 - More relaxed memory models result in too many interleavings
- Provides fool-proof **secure** language memory model:
 - Java Memory Model (JMM) took years to design / verify
 - Complexity of JMM is the cause of insecure implementations

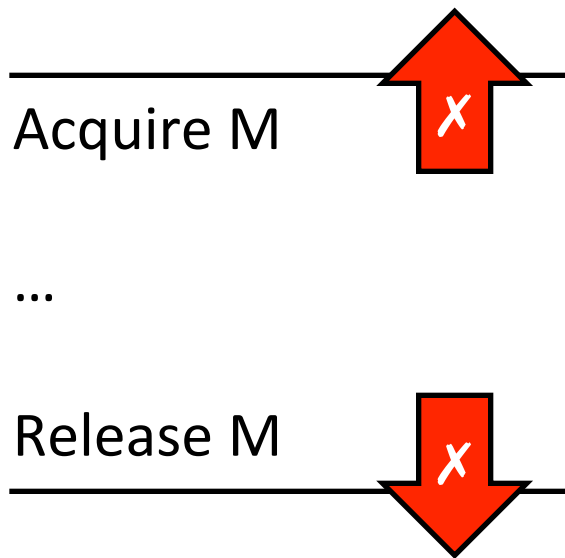
- Current memory consistency models are much more relaxed
 - To allow Compiler / HW reordering of memory accesses
 - For the purposes of performance optimization



Current Memory Models

Sometimes Prevent Optimizations

- All current memory models impose strict ordering around synchronization points (e.g. acquire / release)



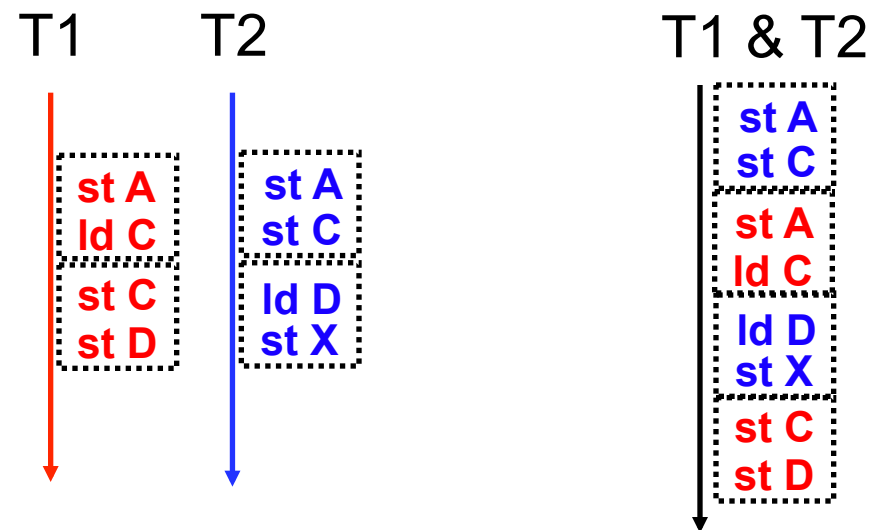
```
X Illegal  
for (...) {  
  Acquire M  
  a = b + c;  
  A[i] = a;  
  Release M  
}
```

- Current memory consistency models have two problems
 - Too relaxed: hard to program and debug
 - Too restrictive: prevents optimization across sync points



Idea: Use Transactions for Memory Ordering Speculation

- Group chunks of dynamic instructions into transactions
- Commit transactions in a total order → **SC is achieved**



- **Higher performance** through memory ordering speculation
 - Within a transaction, Compiler / HW can freely reorder instructions
 - Transaction guarantees isolation (rollback if violated)
- Transactional Memory (TM) hardware is already commercialized
 - Intel Haswell, IBM Power 8, IBM Bluegene/Q



BulkCompiler: Transformations on Synchronization

- **Low-contention** critical sections:
 - Group many of them in same transaction and **optimize**

beginAtomic

*i*₁

acquire M1

*i*₂

release M1

*i*₃

acquire M2

*i*₄

release M2

*i*₅

endAtomic

beginAtomic

*i*₁

while (M1 == taken) {}

while (M2 == taken) {}

while (M2 == taken) {}

*i*₅

*i*₂

*i*₃

while (M2 == taken) {}

*i*₄

*i*₁

*i*₃

*i*₁

endAtomic



Speedups relative to current compilers

- Remove acquire / release
- **Optimize** and reorder the code

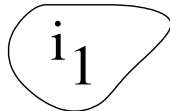


Pitt

BulkCompiler: Transformations on Synchronization

- **High-contention** critical sections:
 - Tight-fit a transaction
 - Reduce chance of rollbacks

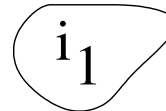
beginAtomic
acquire M1



release M1
endAtomic



beginAtomic
while (M1 == taken) { }

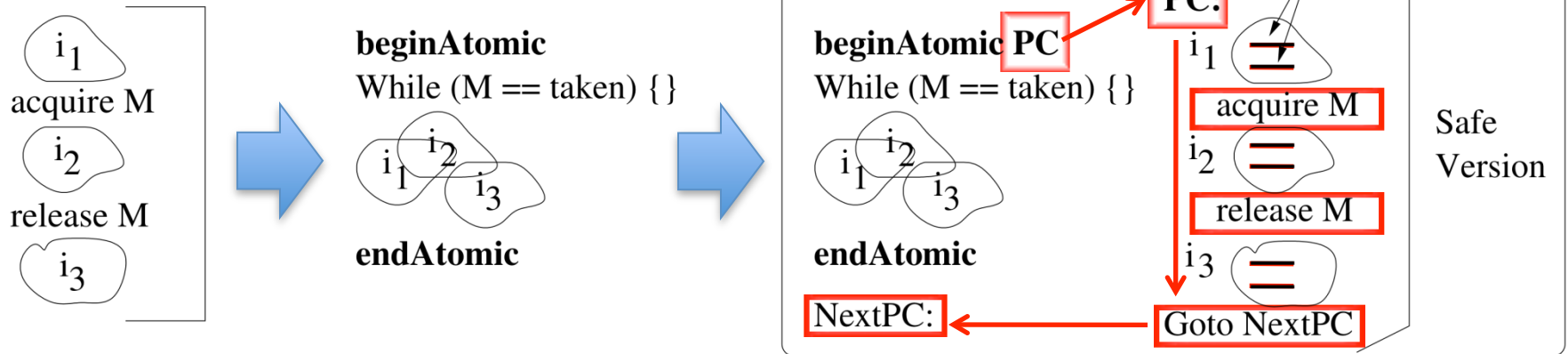


endAtomic

→ Prevent slowdowns due to chunk rollbacks



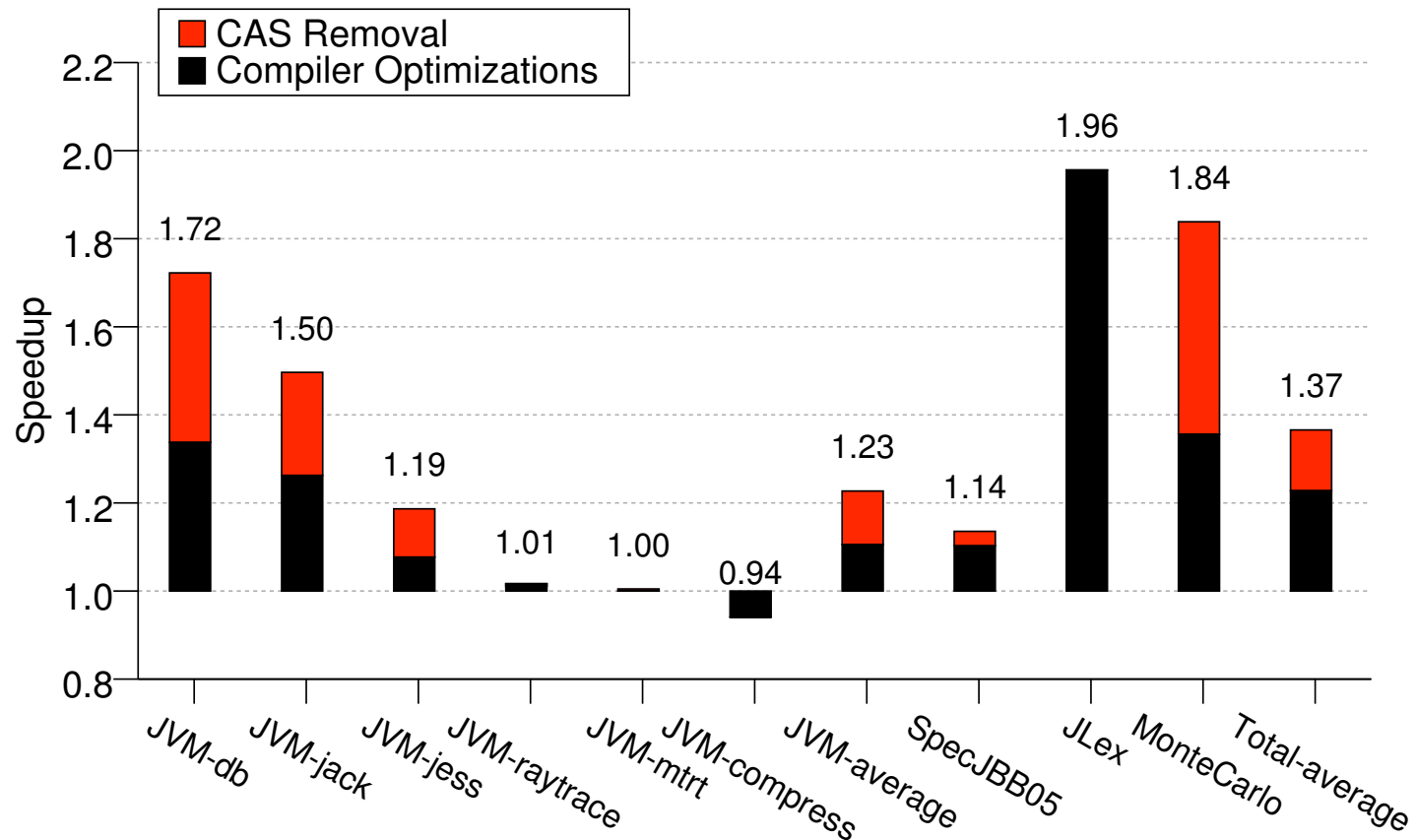
Guaranteeing Forward Progress: Safe Versions



- Semantically a duplicate of the transaction. But..
 - Does not use transactions to ensure SC (instead use fences)
 - Includes original synchronization accesses
- Transactions fall back on Safe Versions when forward progress problem
 - Safe Version is guaranteed forward progress (no rollbacks)



BulkCompiler: Speedup over HotSpot Java Compiler



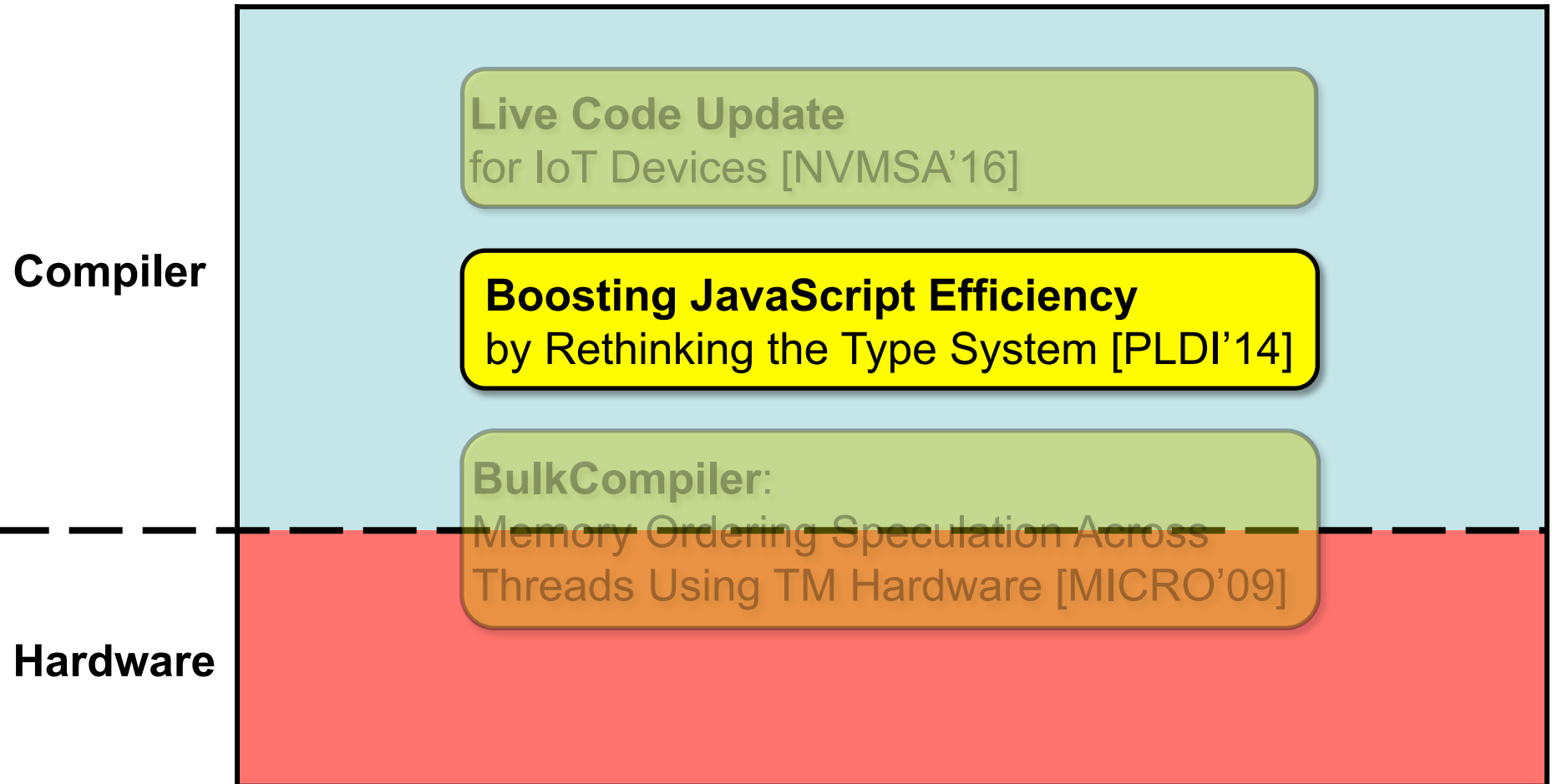
- Average speedup is 37% (23% from compiler opts)
- While still supporting SC



BulkCompiler Summary

- Built BulkCompiler on top of Sun Hotspot Java Compiler
 - Optimizing inside transactions by leveraging atomicity
- **37% speedup** over original Sun Hotspot Java Compiler
 - While upgrading Java Memory Model to **Sequential Consistency**
 - Improving **programmability** and **security**

My Contributions



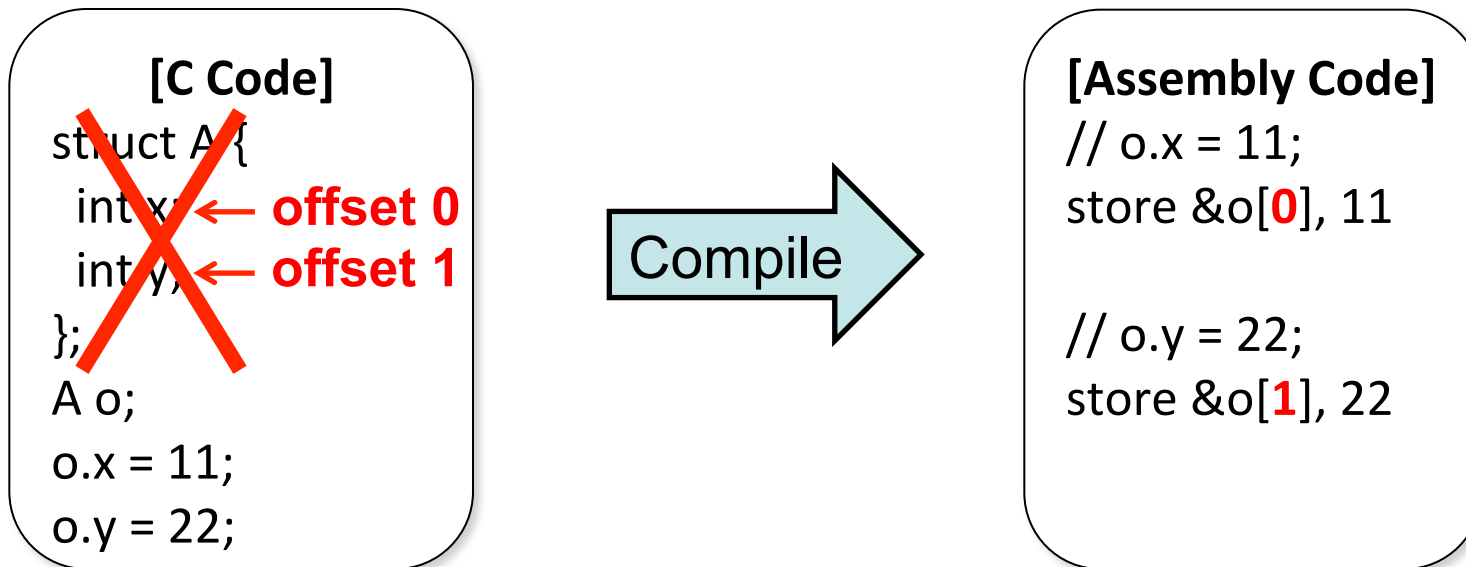
Scripting Languages and Efficiency

- JavaScript, PHP, Python, Perl, Ruby, Lua, R
- Have common features that make programming easier but efficient code generation much much harder
 - **Dynamic Typing**
 - Dynamic Dispatch
 - First Class Functions
 - Runtime Evaluation (`eval('string')`)



Why Types Are Important for Performance

- Types are crucial to generating efficient code

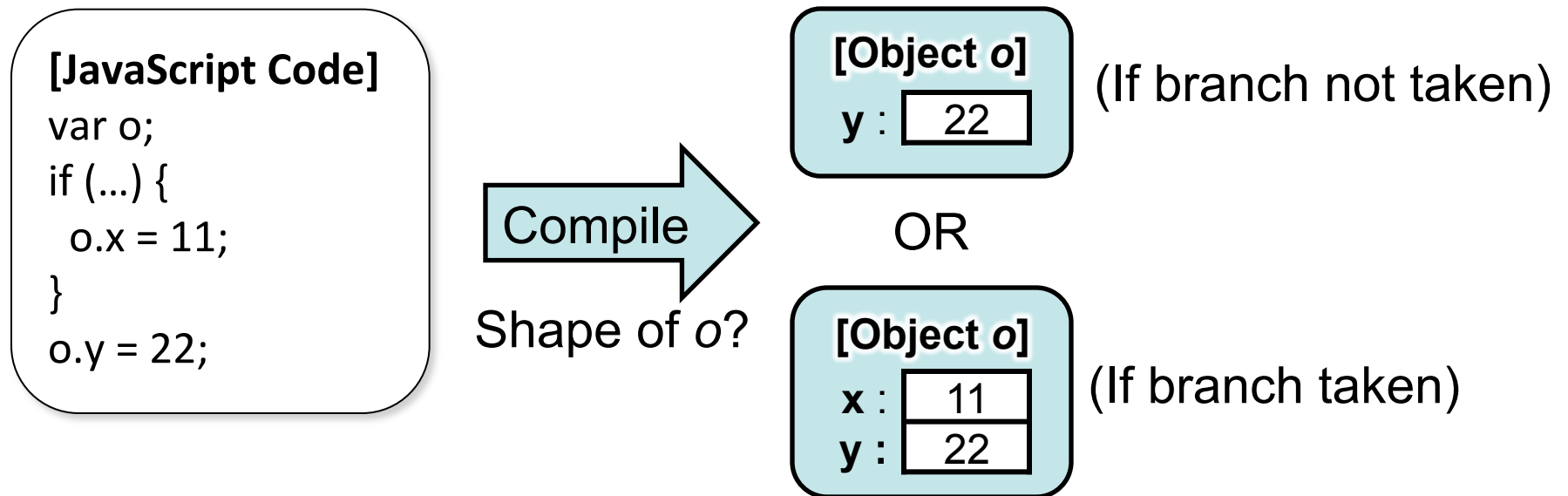


Types tell compiler the shape of `o` (fields and their offsets)



Scripting Languages Have No Types → Slow

- Objects are simply dictionaries from properties to values
- Properties can be added and removed at any time



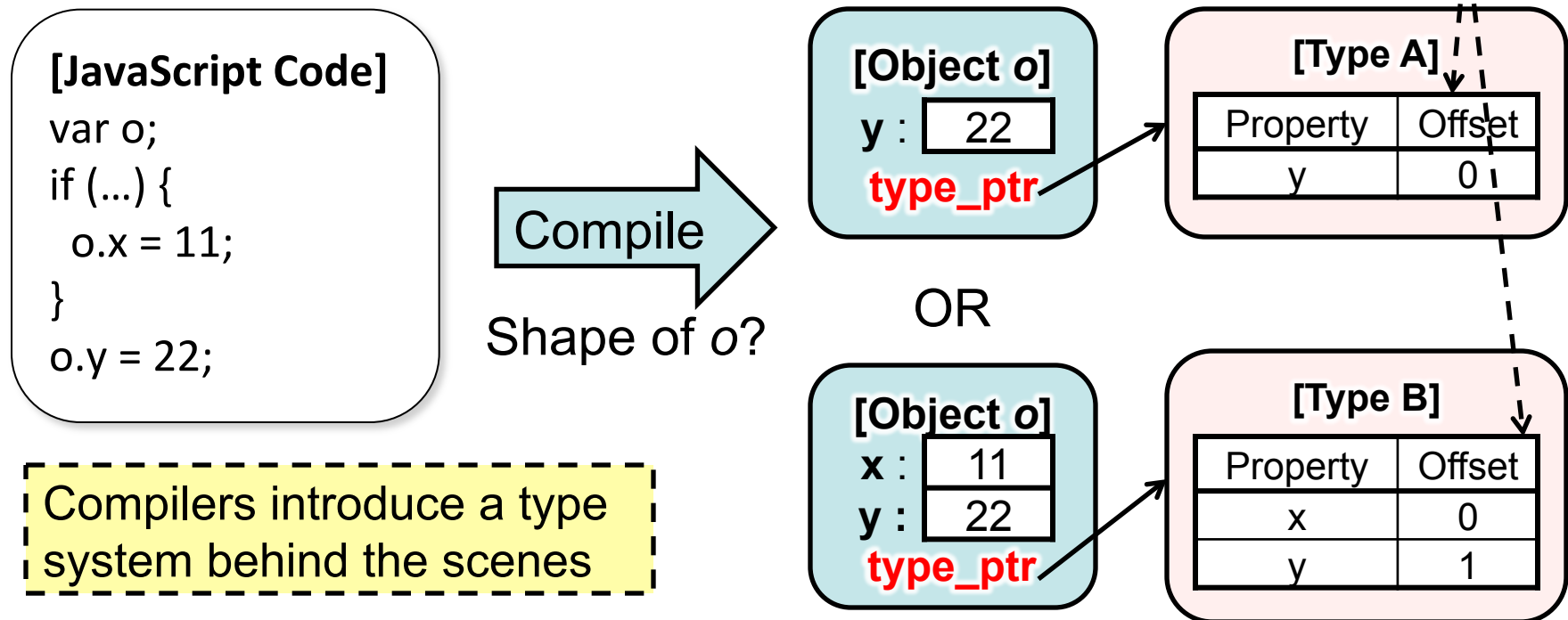
How to generate code when shape of *o* is unknown?



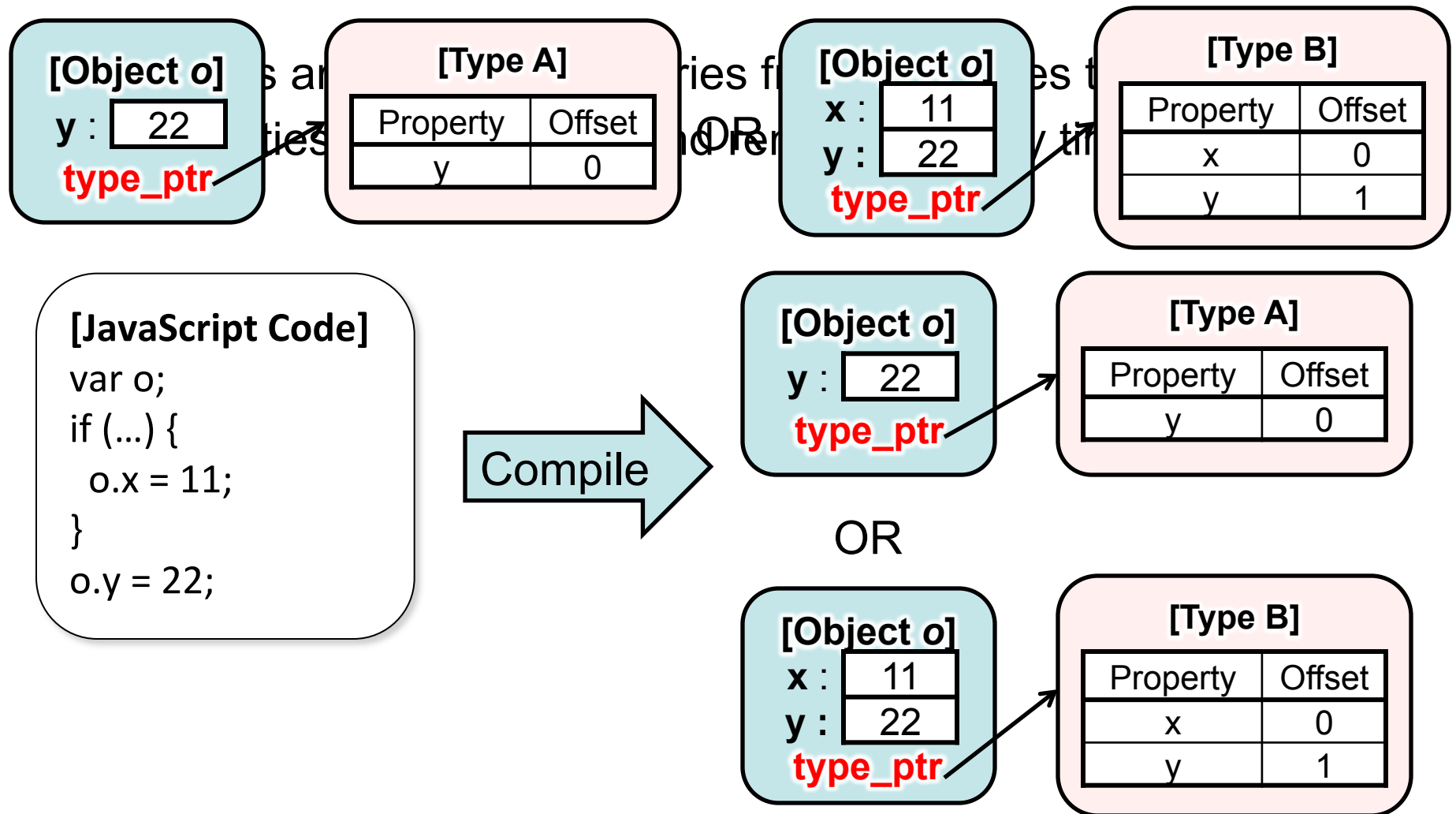
Scripting Languages Have No Types → Slow

- Objects are simply dictionaries from properties to values
- Properties can be added and removed at any time

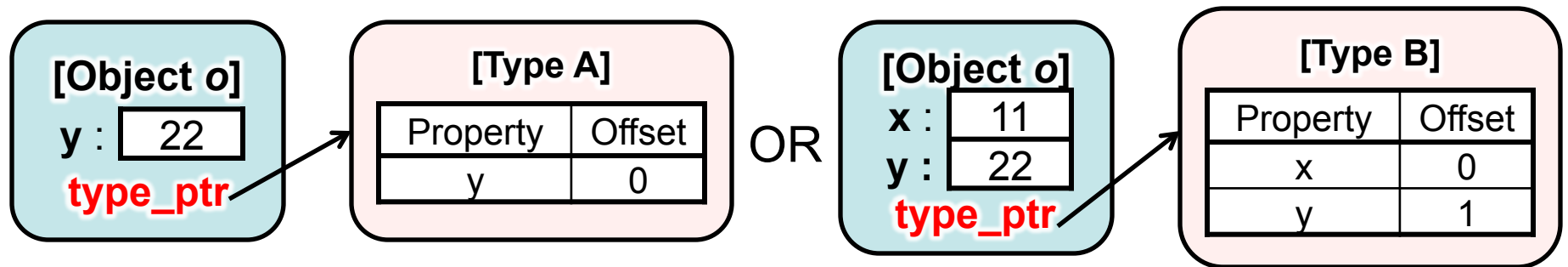
Hash tables



Scripting Languages Have No Types → Slow



Scripting Languages Have No Types → Slow



[JavaScript Code]

```
var o;  
if (...) {  
  o.x = 11;  
}  
o.y = 22;
```

Compile

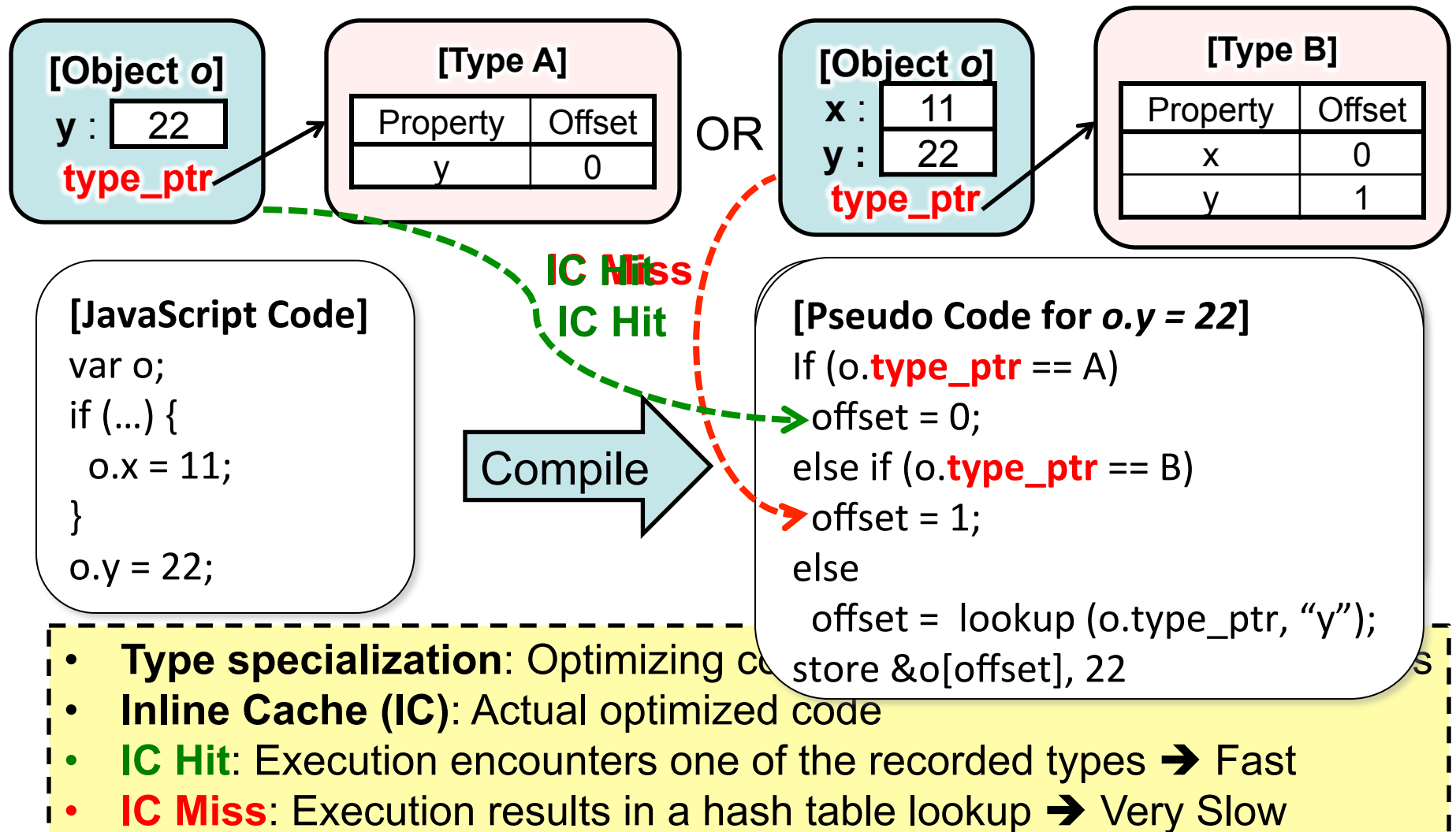
[Pseudo Code for o.y = 22]

```
offset = lookup (o.type_ptr, "y");  
store &o[offset], 22;
```

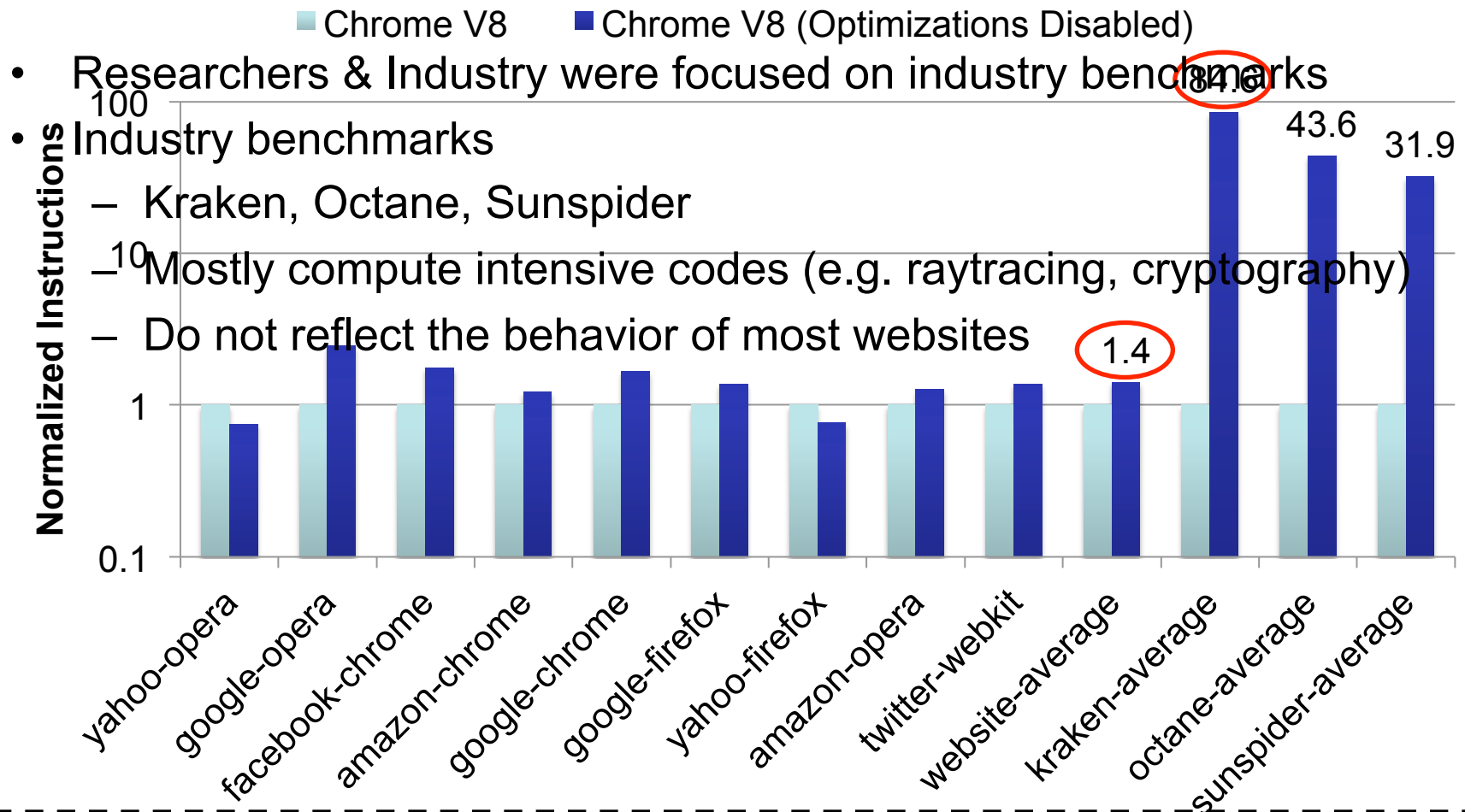
A field access always entails a **hash table lookup** to get the offset



State-of-the-art Compilers do Better: Type Specialization Optimization



Google Chrome V8 JavaScript Compiler: Ineffective for Real Websites



Chrome V8 compiler not optimized for dynamism in real websites



Problem: V8 Type System Too Brittle [PLDI 2014]

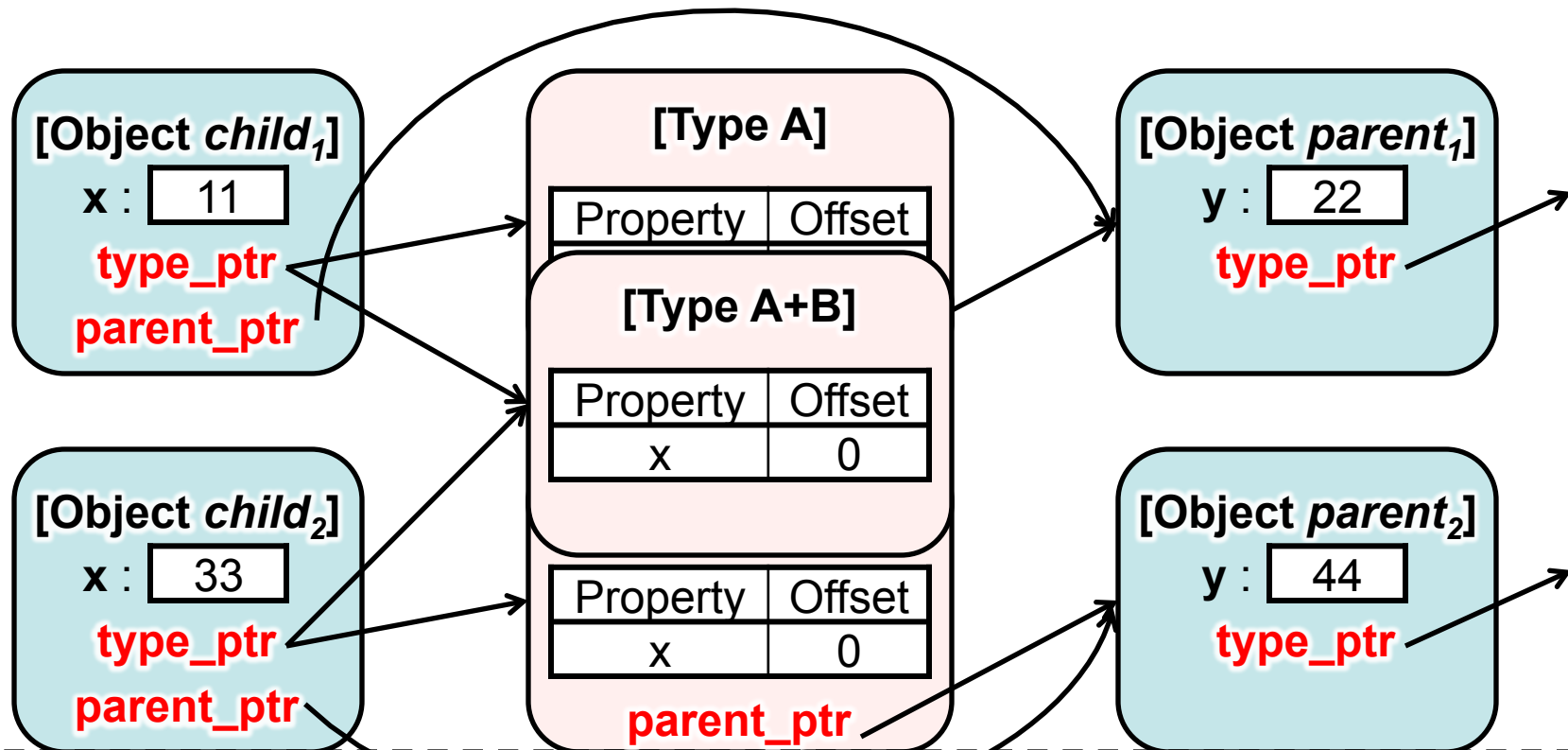
- Chrome V8 type system encodes (other than property offsets):
 - **Inheritance** (i.e. address of parent object)
 - **Method bindings** (i.e. addresses of method properties)
- Helps in generating efficient code during type specialization
 - Inheritance: helps when accessing parent properties
 - Method bindings: helps resolve targets for method calls
- V8 Assumption: inheritance and method bindings rarely change
 - Reasonable since always true for statically typed languages

Reality: Assumption **NOT TRUE** for dynamic website code
→ Leads to type dynamism and frequent inline cache misses
(and terrible performance)



Solution: Decouple Inheritance from Type System

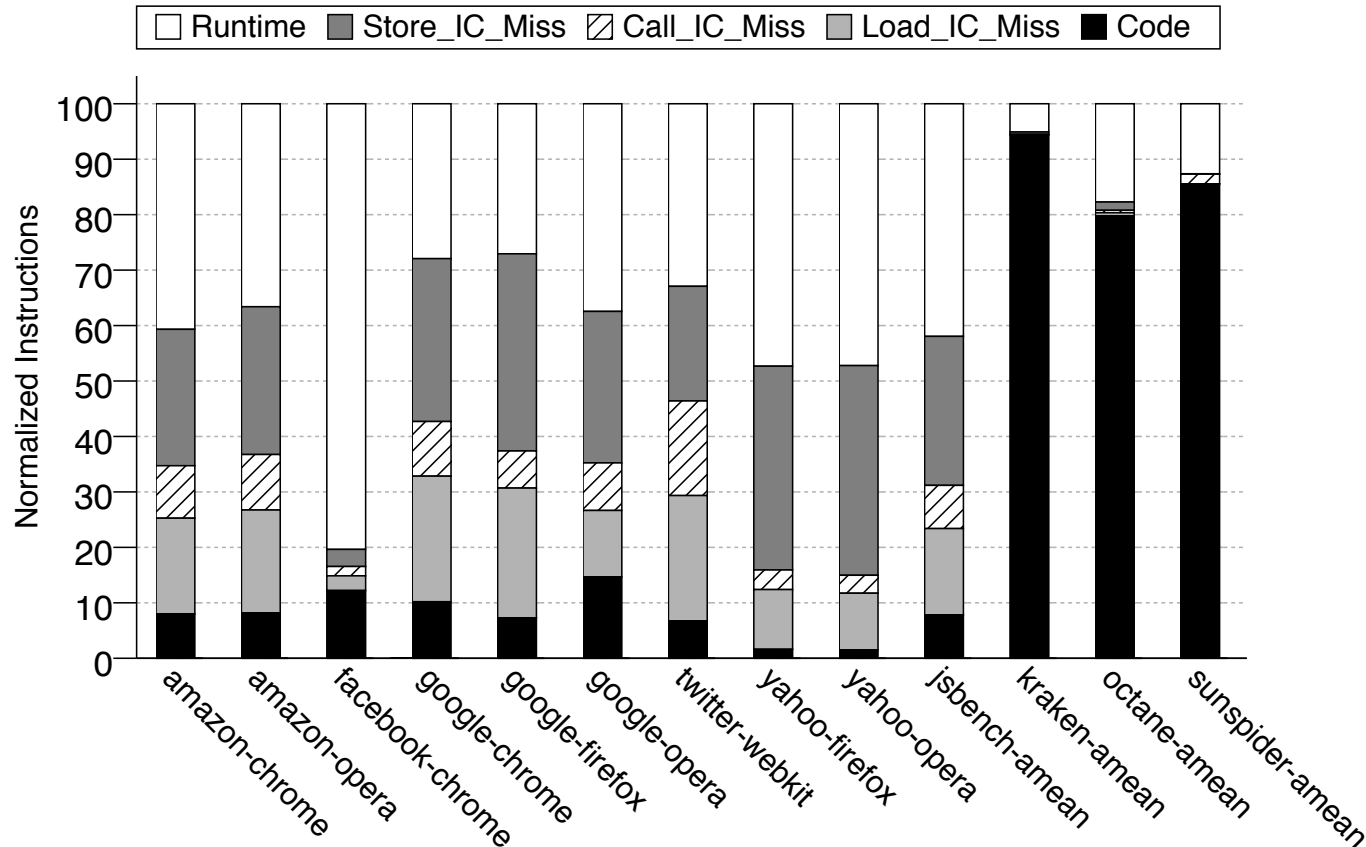
- Decoupling Inheritance from the Type System



- Types A & B have exactly same shapes but different types
- Leads to needless type dynamism and inline cache misses



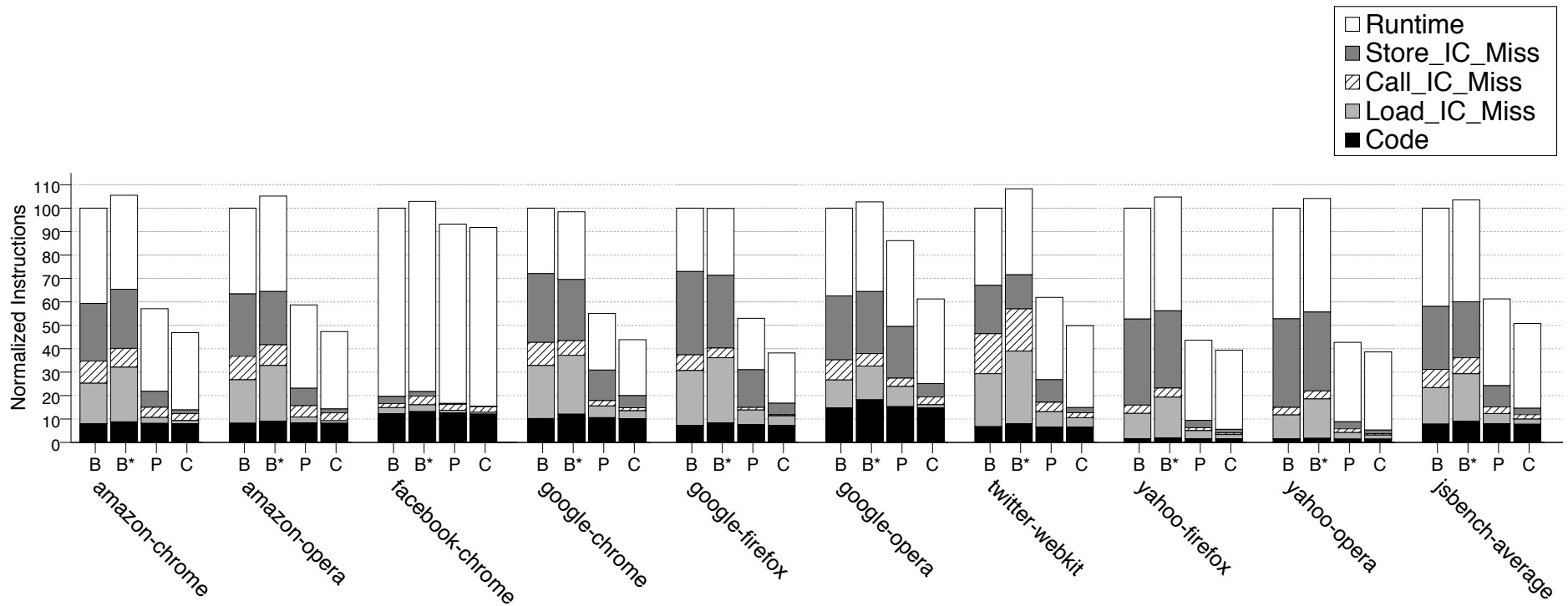
Instruction Breakdown for Chrome V8



- JSBench: 50% of instructions spent on inline cache miss handling
- Kraken, Octane, Sunspider: Most instructions spent on generated code



Instruction Count Reduction over Chrome V8



- Almost all instructions due to inline cache misses removed
- Total instruction count decreased by 49% compared to Chrome V8
- Negligible change for Kraken, Octane, Sunspider (not shown)



Summary

- Made the type system of Chrome V8 JS Compiler more flexible:
 - By decoupling inheritance and method binding from types

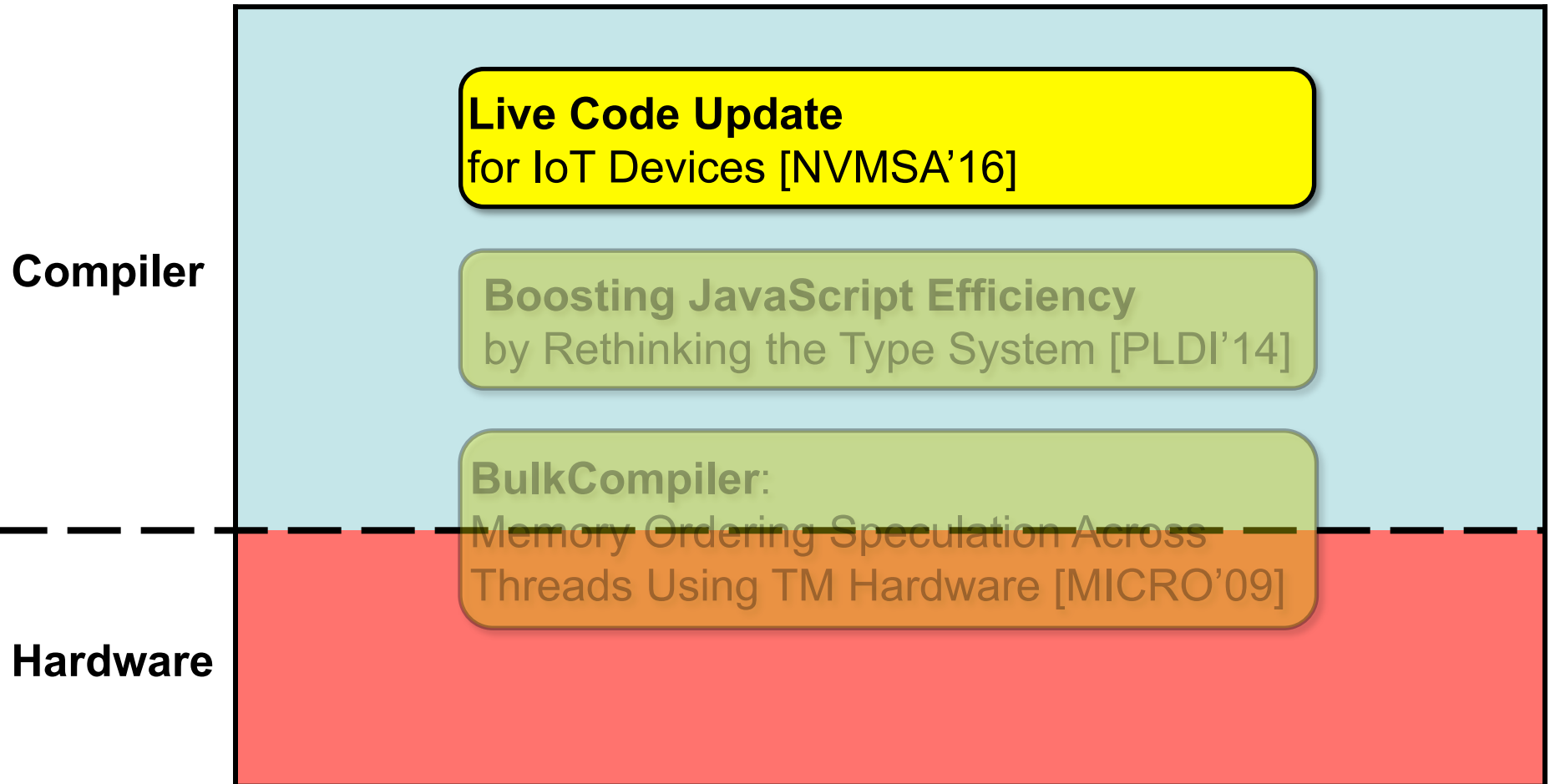
Results in:

49% in dynamic instruction count reduction

20% in heap memory usage reduction



My Contributions



Code Update for IoT Devices Design Goals

- Minimize **device down time**
 - Consists of reboot time after code update
 - Prolonged down time can impact QoS
 - For energy harvesting devices, it may take days
- Minimize **reprogramming time**
 - Consists of code transfer time + code write time
 - Important in iterative debug / tuning scenarios
 - For energy harvesting devices, it may take weeks
- Minimize **extra code memory**
 - Consists of memory to store new updated image
 - Larger memory means more costly devices



Previous Code Update Approaches

Reprogramming time

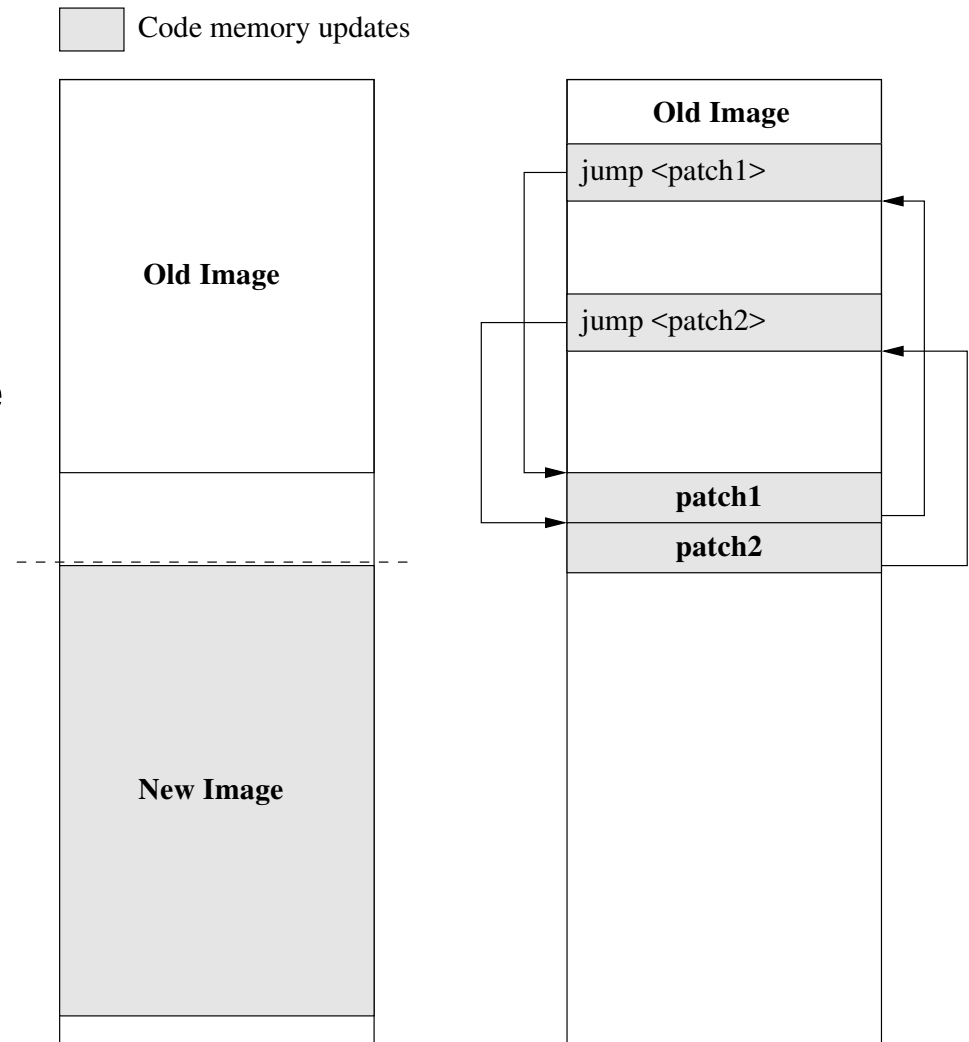
- Naïve:
 1. Rebuild code image at base station
 2. Transfer image wirelessly to each IoT device
 3. Write image in separate area in code memory in device
 4. Reboot device from image
- Extra memory
- Device down time

- Incremental: Improves upon naïve by transferring only diff
 - Decreases transfer time
 - But must still rewrite image and reboot



Our Code Update Approach [NVMSA 2016]

- Live update:
 - Apply patches to **live** image in device **in-place**
 - No down time
 - Minimal reprogramming time
 - Minimal extra memory



Previous Approaches **Our Approach**



Atomic Code Patching

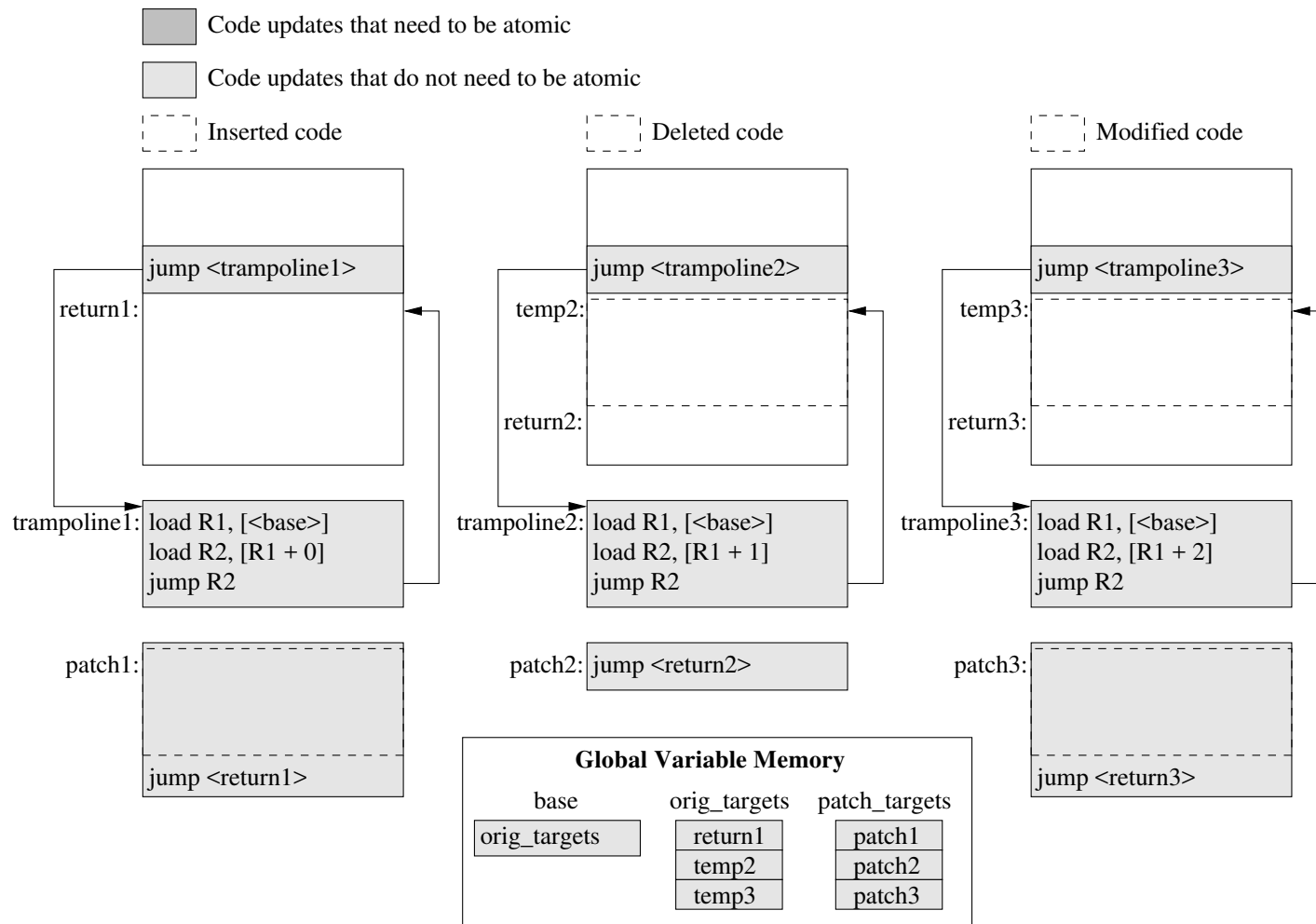
- Problem: Need to perform all patches **atomically**
 - Cannot expose live code in an inconsistent state
- Solution: Use **indirect jump vectors**
- Indirect jump vectors
 - Table holding addresses pointing to pieces of code
 - Used to implement interrupt service routines in OSes

→ Idea: change functionality instantaneously by switching over from one jump vector table to another jump vector table

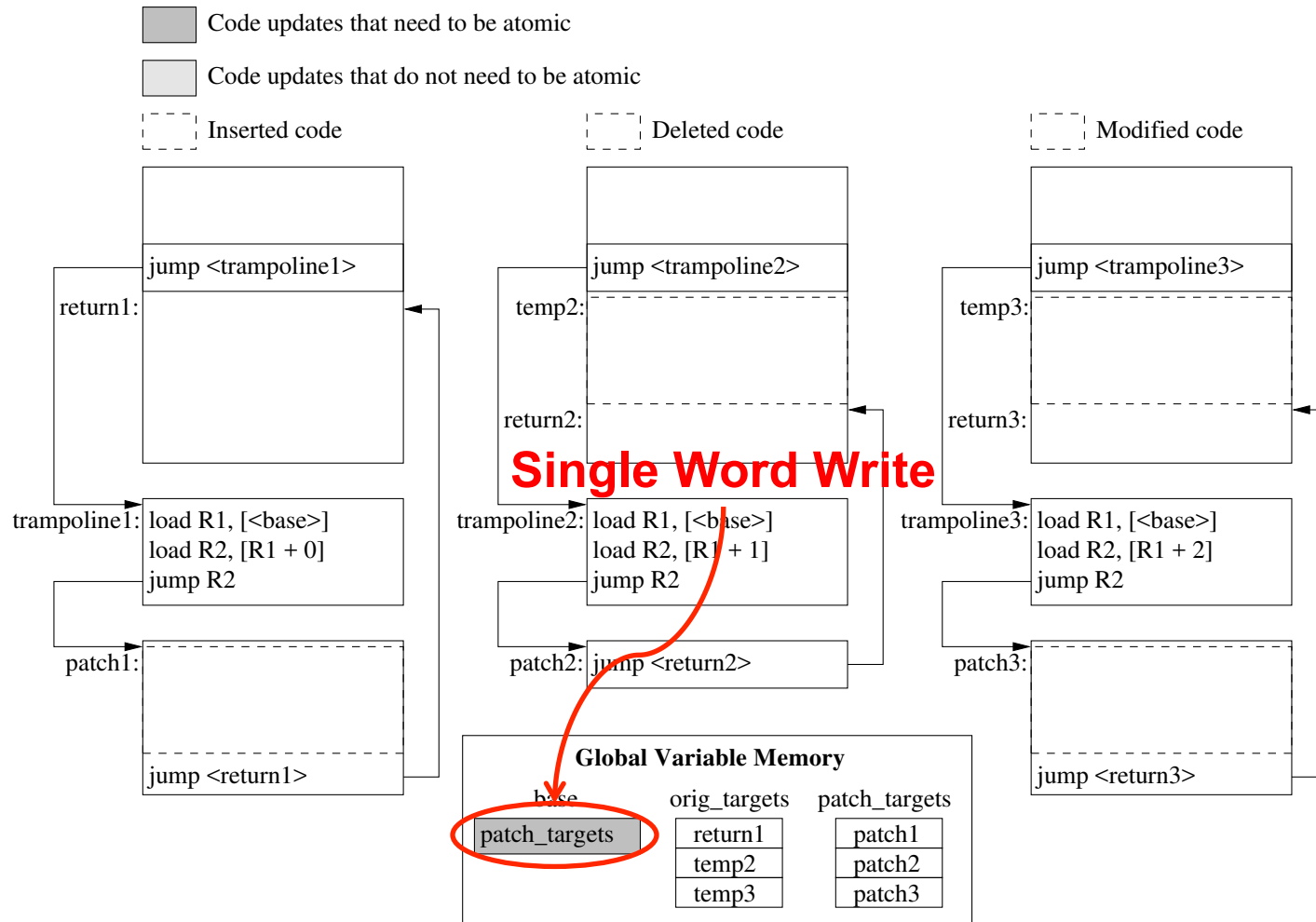
→ Perform the switch through an (atomic) **single word write**



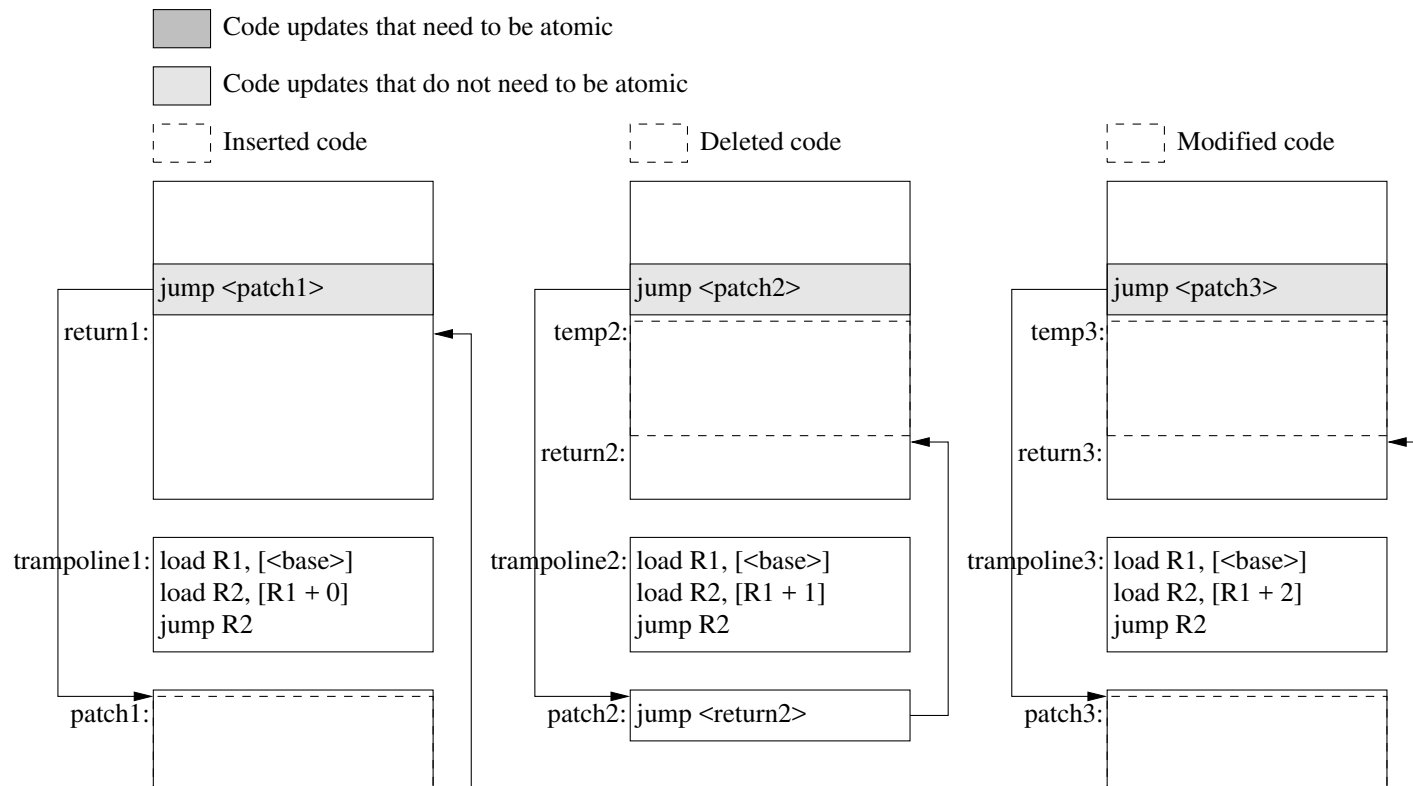
Phase 1: Trampoline Insertion



Phase 2: Trampoline Switching



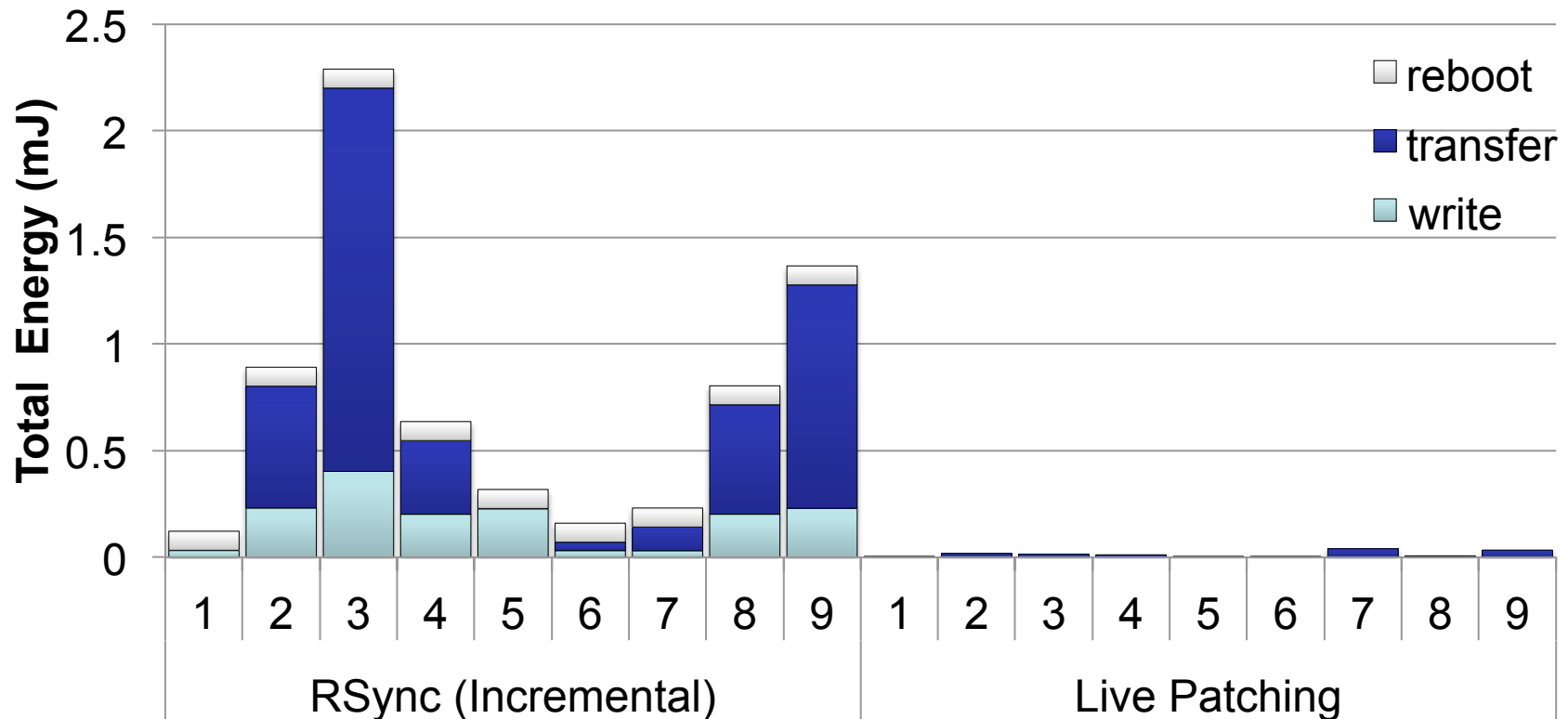
Phase 3: Trampoline Removal



- Control flow is streamlined for better performance
- Trampolines are recycled to implement next patch



Reprogramming Time Reduction



- Reboot eliminated for live patching (no device down time)
- Transfer & write reduced due to no code shifts for live patching



Summary

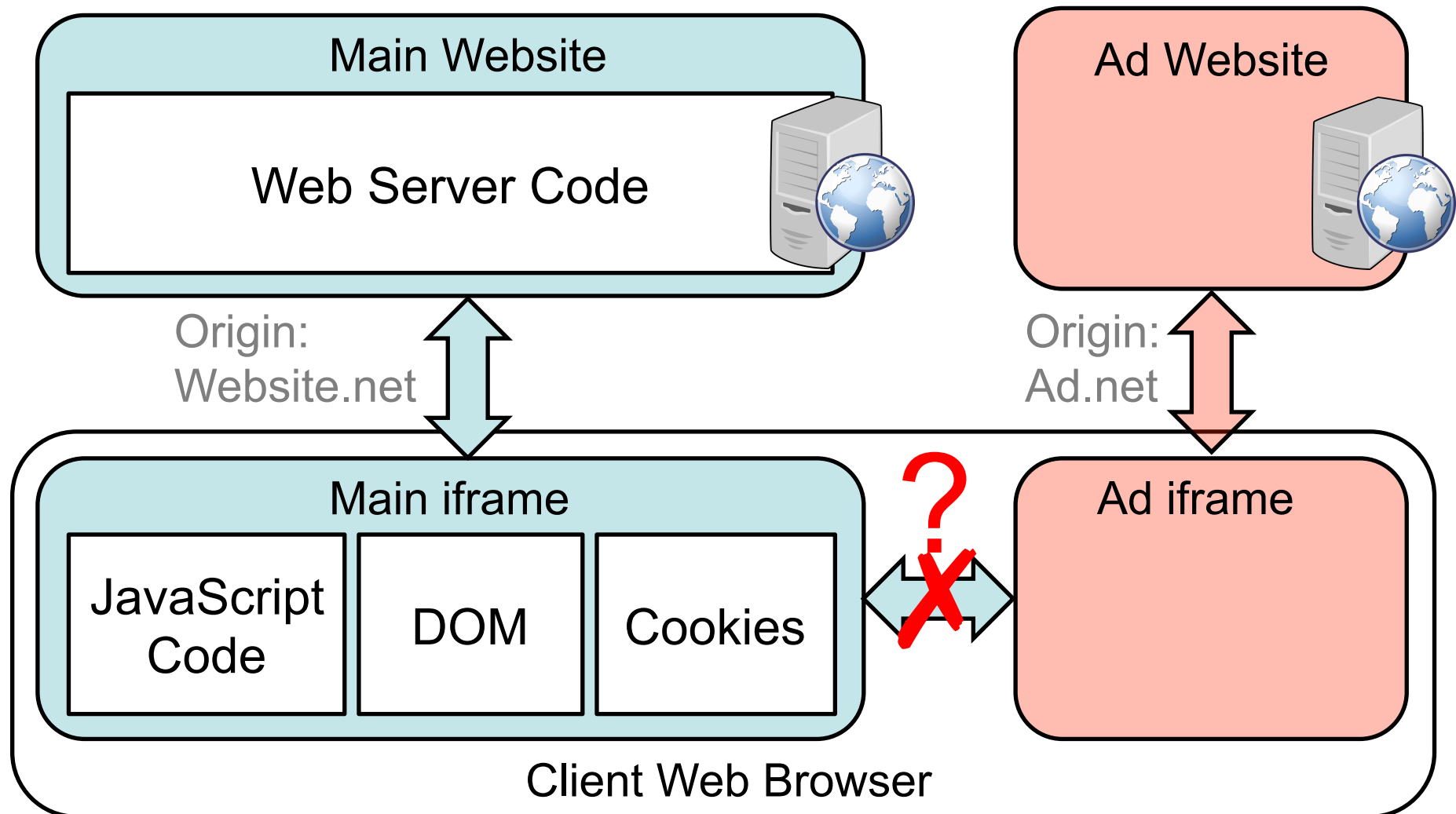
- Proposed live code update scheme for IoT devices that can:
 - Improve QoS by **eliminating down time**
 - Improve usability by **minimizing reprogramming time**
 - With minimal runtime overhead ($< 2.7\%$)

Future Work: Hardware/Compiler co-optimization

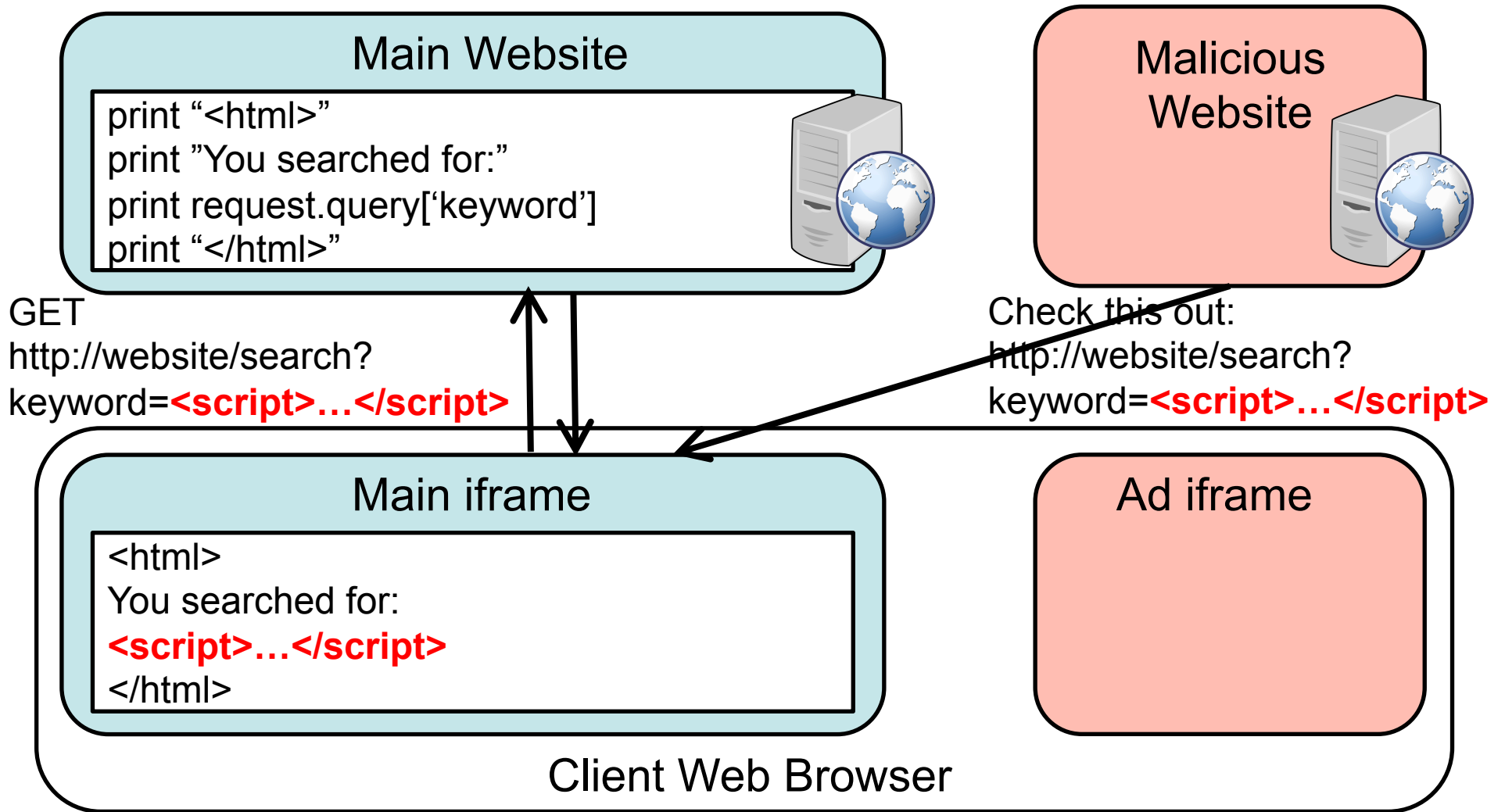
- **Scripting on IoT** for improved **security** / **programmability**
 - Re-engineer VM to fit in limited device memory
 - Offload compilation to base station
 - Device VM needs only runtime (garbage collector etc.)
 - More efficiency improvements through SW / HW co-design
 - Decrease heap usage through HW-assisted compiler optimizations such as stack promotion
 - Energy efficient GC through co-design
 - Security improvements through SW / HW co-design



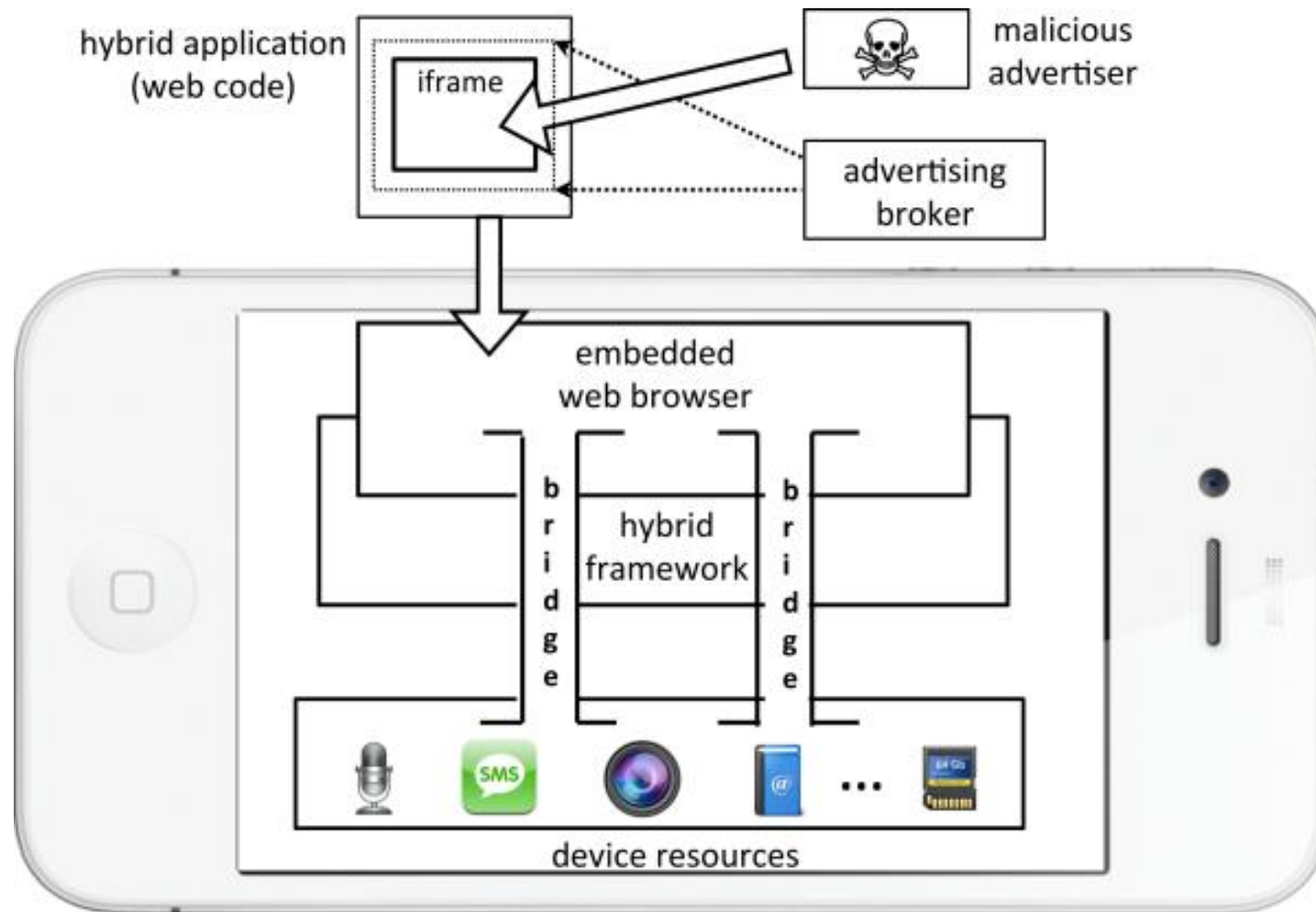
Web Browser Sandboxing – Same Origin Policy



Conventional Vulnerability – Cross Site Scripting



New Vulnerability with IoT – Holes in the Sandbox



Future Work: Hardware/Compiler co-optimization

- **Bytecode as ISA**
 - Design rich bytecode that expresses
 - More semantic / behavioral information
 - Security expectations
 - Fault tolerance expectations
 - Quality of Service expectations
 - Approximate computing expectations
 - JIT compile bytecode into underlying hardware
 - Use combination of compiler / HW to meet expectations
 - Use HW support for profiling / runtime checking / speculation