

Cimple: Instruction and Memory Level Parallelism DSL

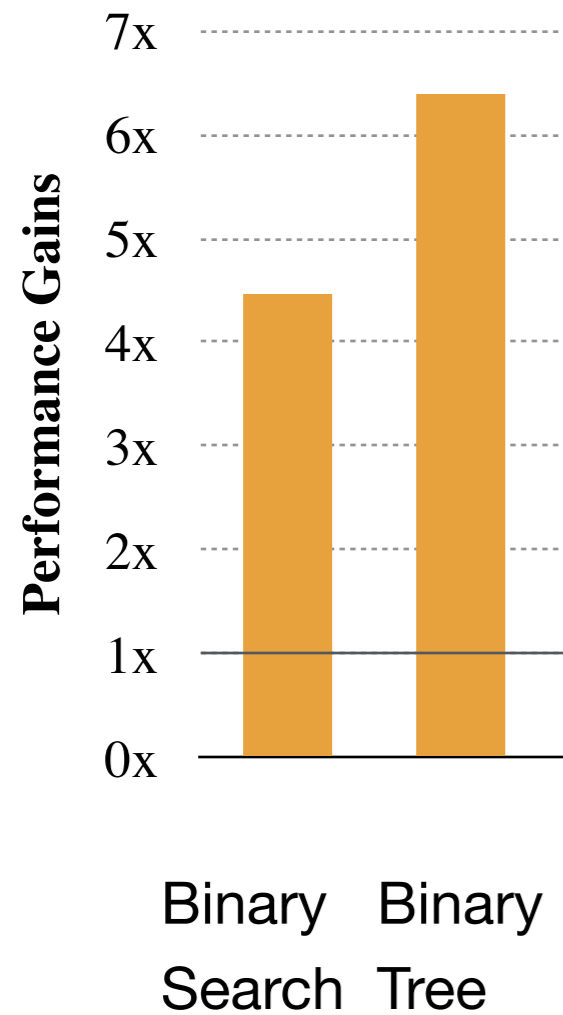
Vladimir Kiriansky, Haoran Xu,
Martin Rinard, Saman Amarasinghe

PACT'18

November 3, 2018
Limassol, Cyprus

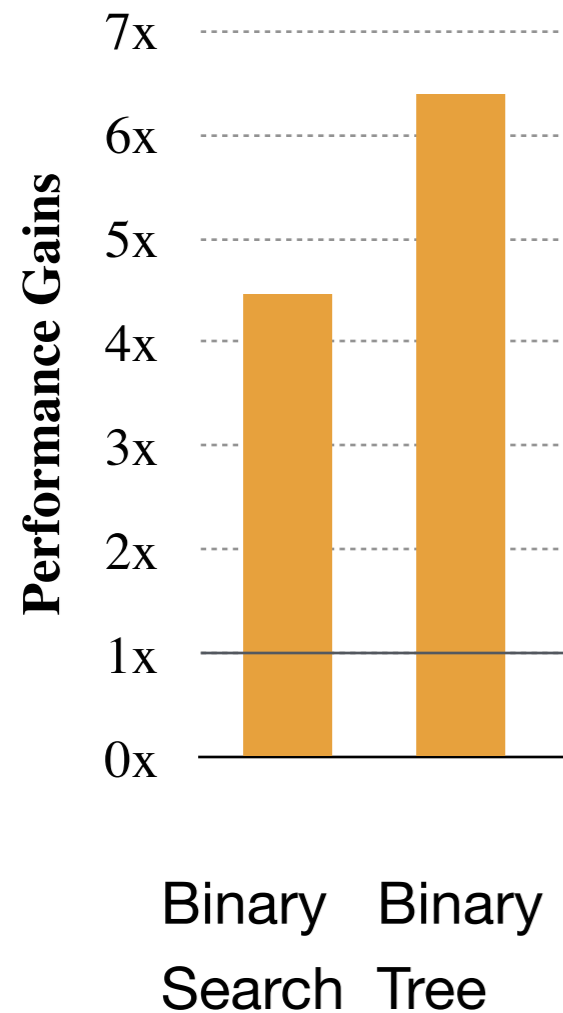


Cimple Performance Gains



[Haswell]

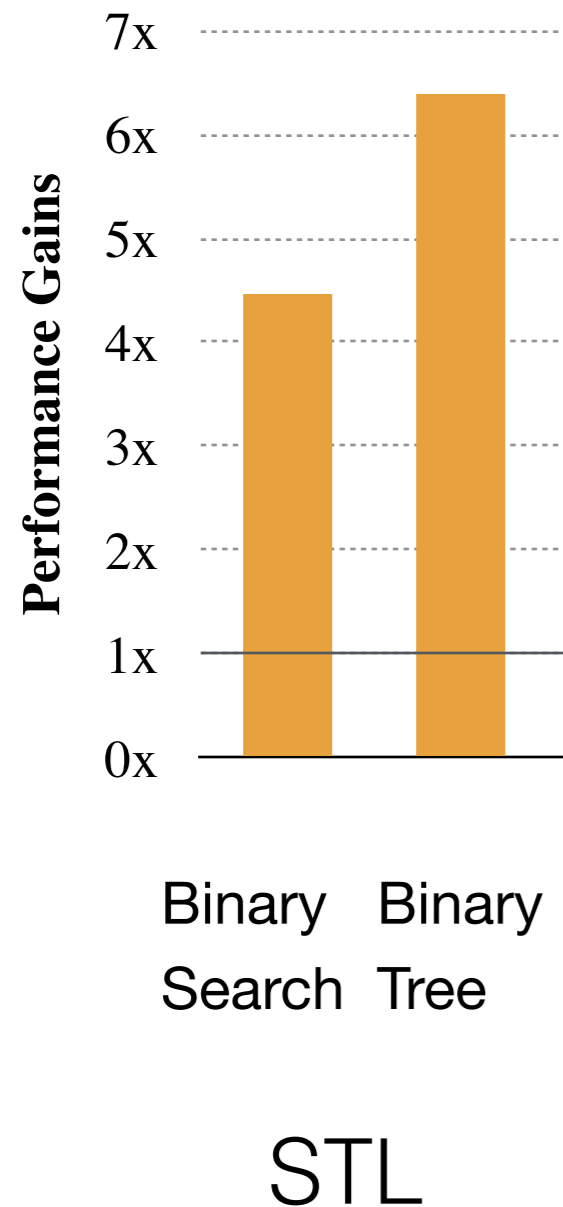
Cimple Performance Gains



- Are we teaching the wrong algorithms?

[Haswell]

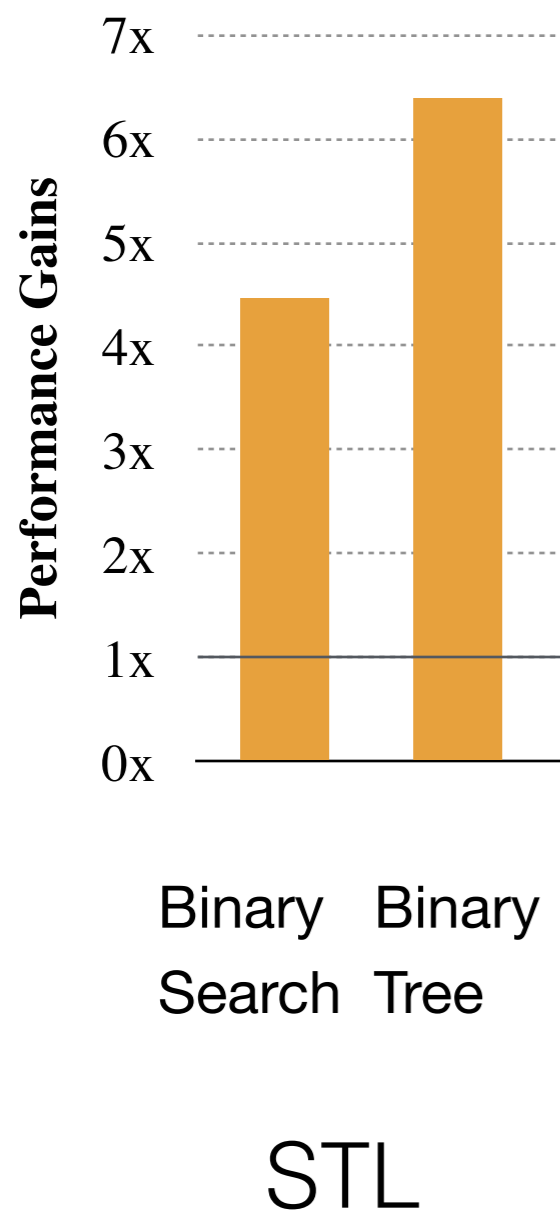
Cimple Performance Gains



- Are we teaching the wrong algorithms?

[Haswell]

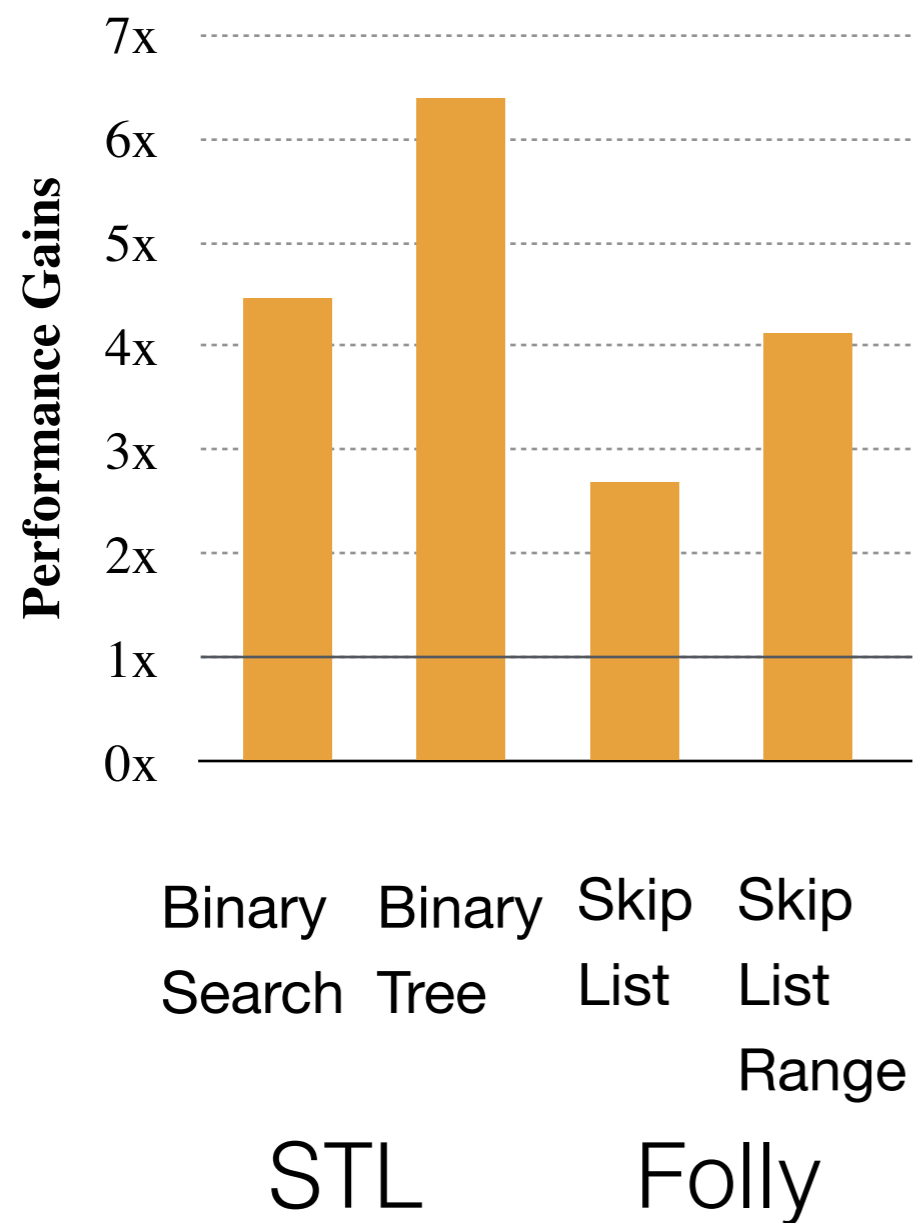
Cimple Performance Gains



- Are we teaching the wrong algorithms?
- Is STL using bad implementations?

[Haswell]

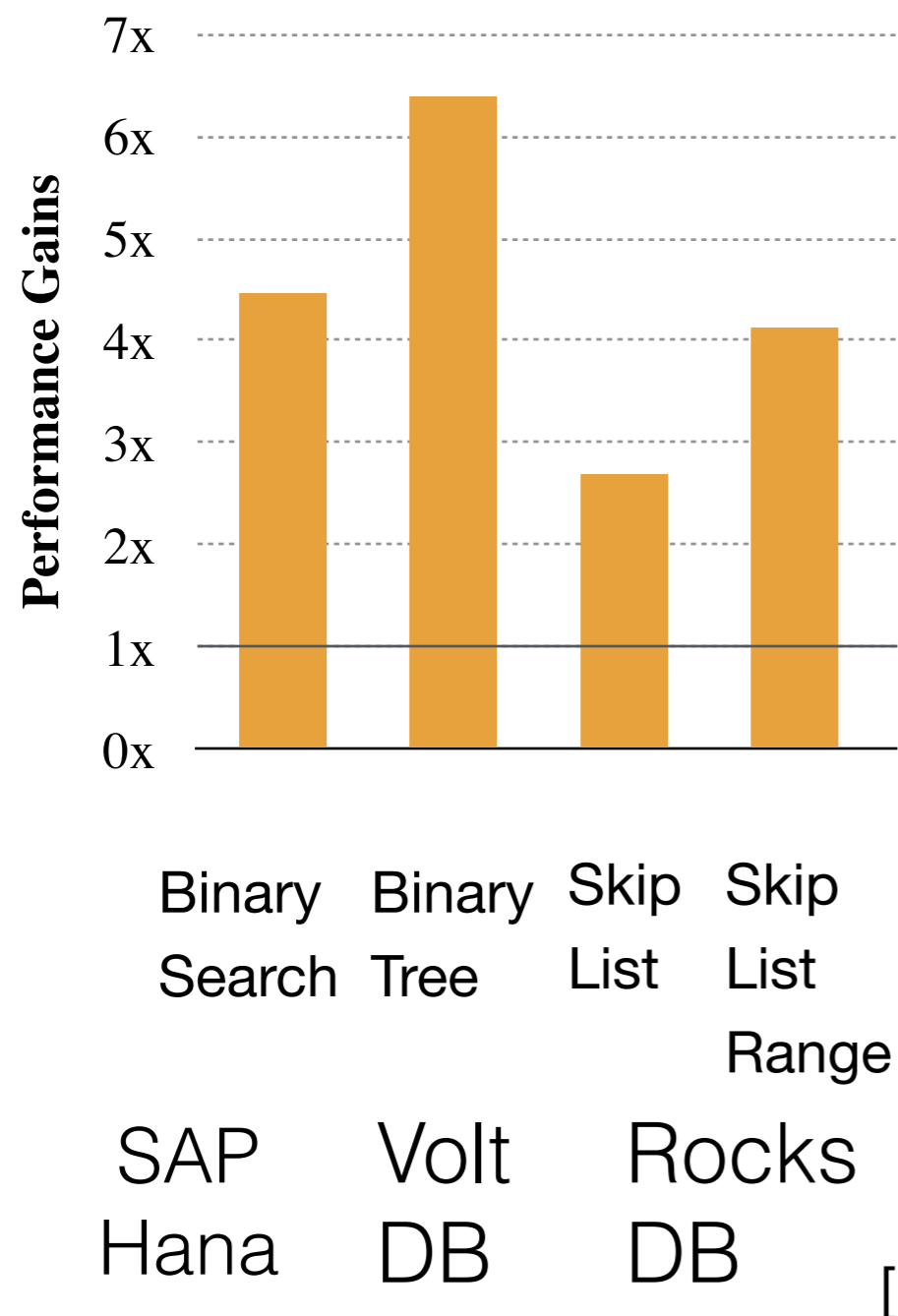
Cimple Performance Gains



- Are we teaching the wrong algorithms?
- Is STL using bad implementations?

[Haswell]

Cimple Performance Gains

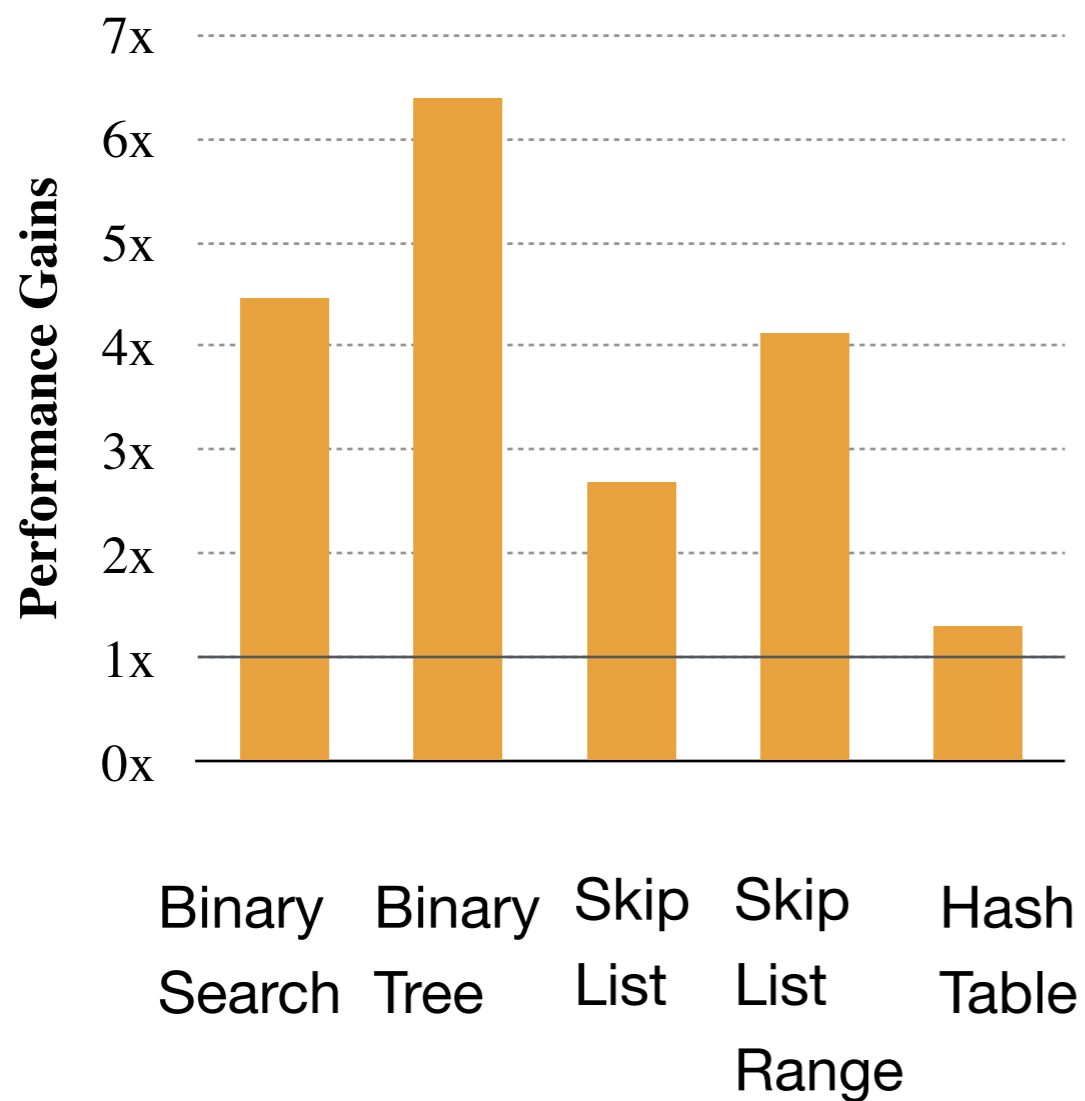


- Are we teaching the wrong algorithms?
- Is STL using bad implementations?

In-Memory Databases

- Terabyte Working Sets
 - AWS 12 TB VM
- Binary Search
 - SAP Hana
- Binary Tree (partitioned)
 - VoltDB
- Skip List (shared)
 - RocksDB, MemSQL

Cimple Performance Gains



- Are we teaching the wrong algorithms?
- Is STL using bad implementations?
- Does our programming model match hardware?

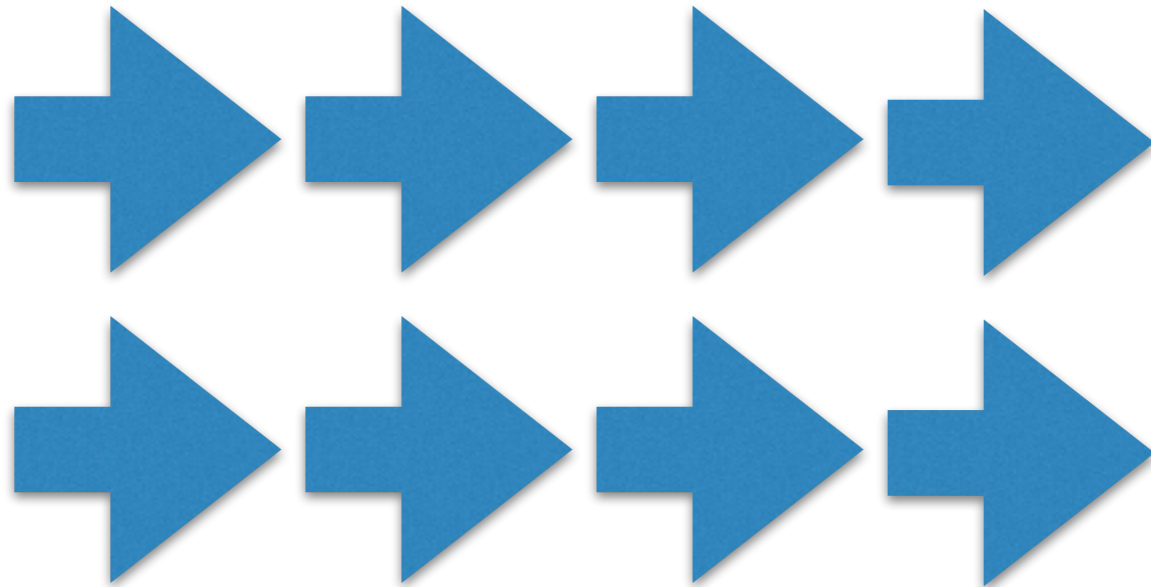
[Haswell]

Little's Law

$$L = \lambda W$$

Arrival
Rate

λ



Concurrency

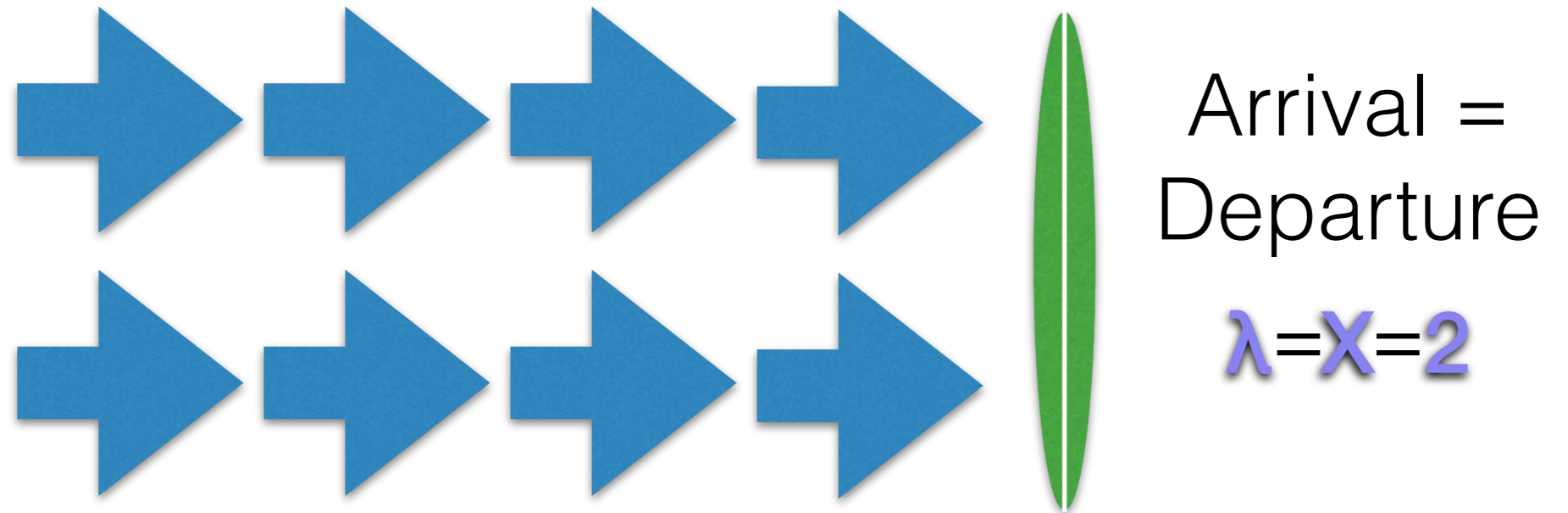
L

Latency

W

Little's Law

$$L = \lambda W$$



Bandwidth

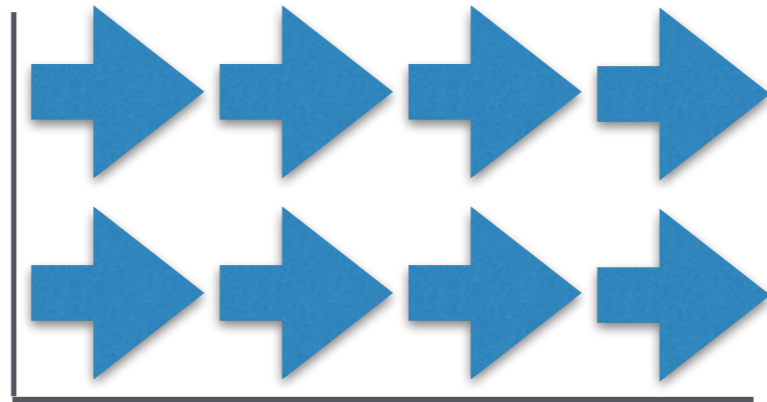
Concurrency

$$L = 8$$

Latency

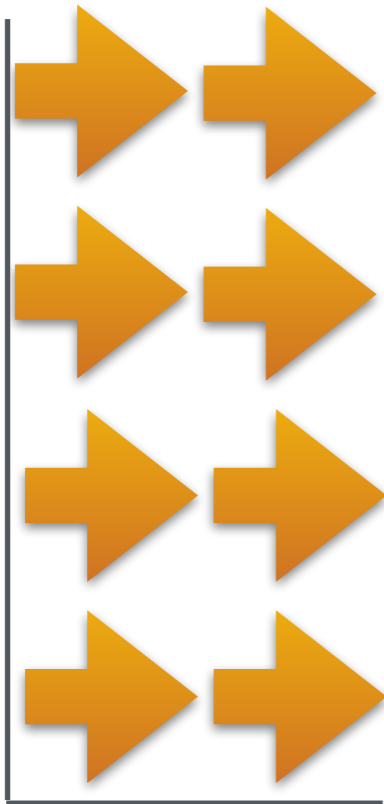
$$W = 4$$

Three Improvement Paths



A) Latency

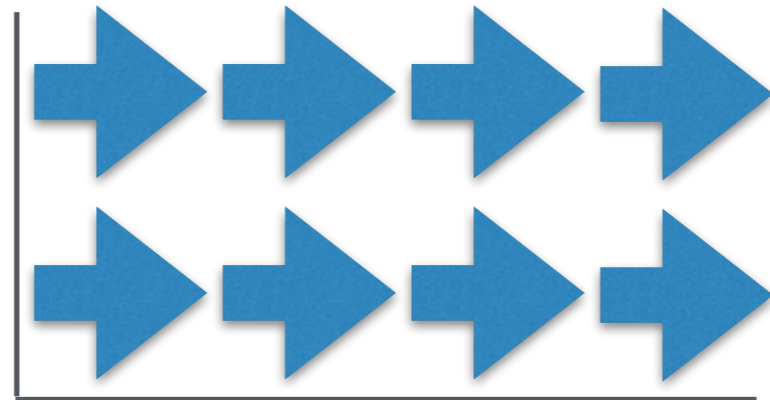
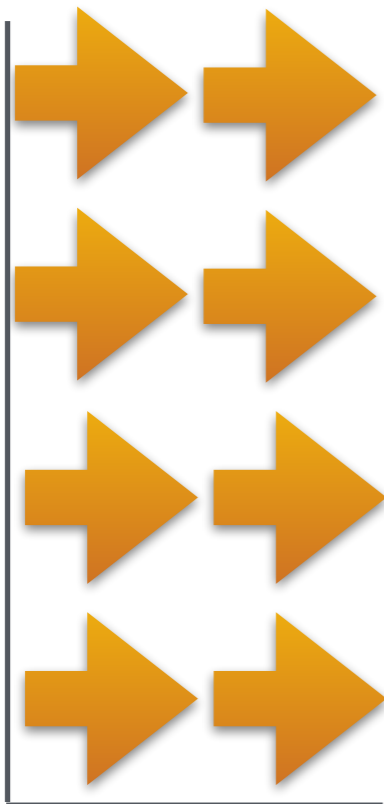
W ↓



Three Improvement Paths

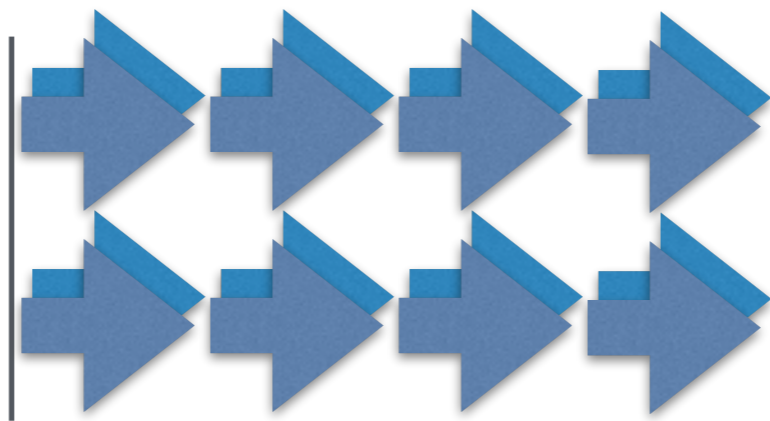
A) Latency

$W \downarrow$



B) Spatial locality

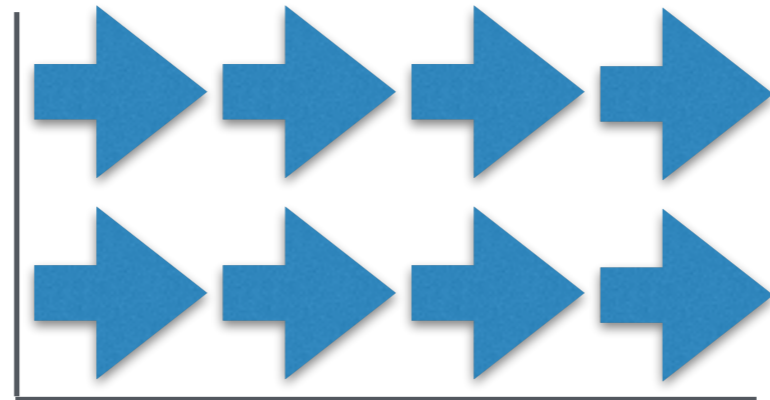
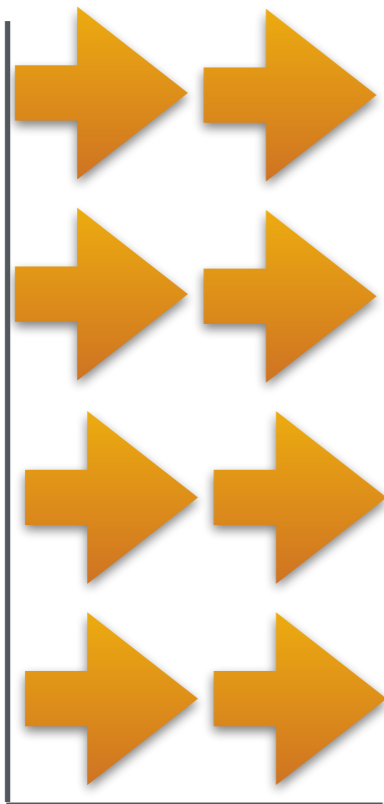
$L \uparrow$



Three Improvement Paths

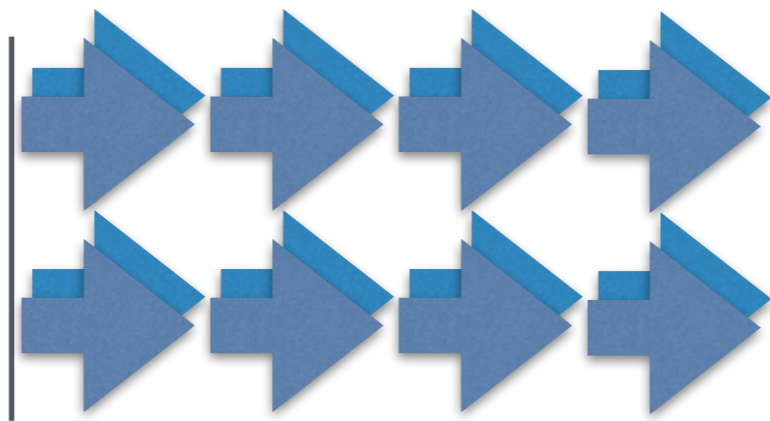
A) Latency

$W \downarrow$



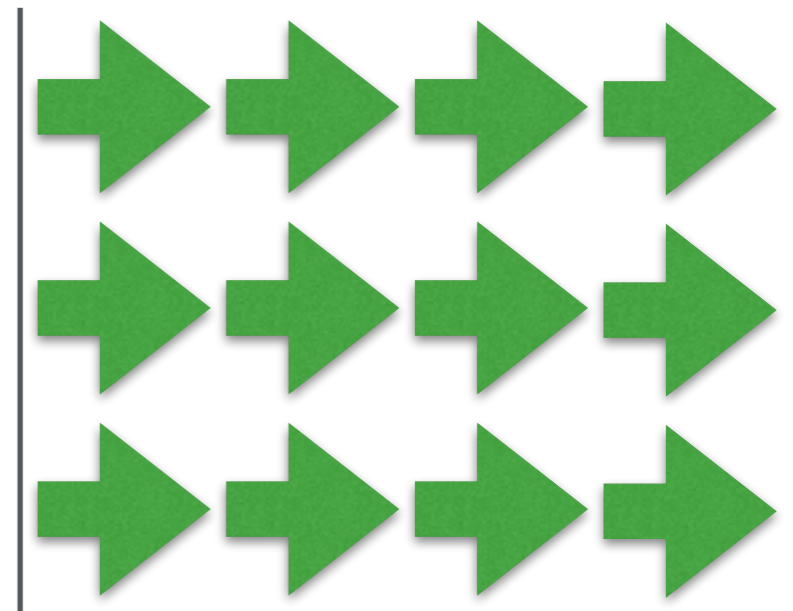
B) Spatial locality

$L \uparrow$



C) Memory Level Parallelism (MLP)

$L \uparrow$



Memory Wall

- Speculative out-of-order processors:
automatically discover MLP



Memory Wall

- Speculative out-of-order processors: automatically discover MLP
- Non-blocking caches - Miss Status Handling Registers (MSHRs)
- Large Instruction Windows



Memory Wall Conquered?

- Speculative out-of-order processors:
automatically discover MLP
- Non-blocking caches -
Miss Status Handling Registers
(MSHRs) = 10 misses
- Large Instruction Windows
~200 instructions



Branch Mispredictions



Memory Wall Conquered?

✘ Branch Misprediction

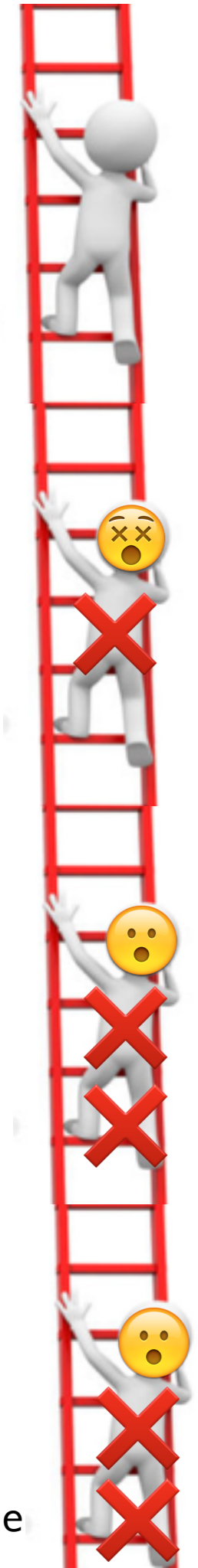


Memory Wall Conquered?

✗ Branch Misprediction



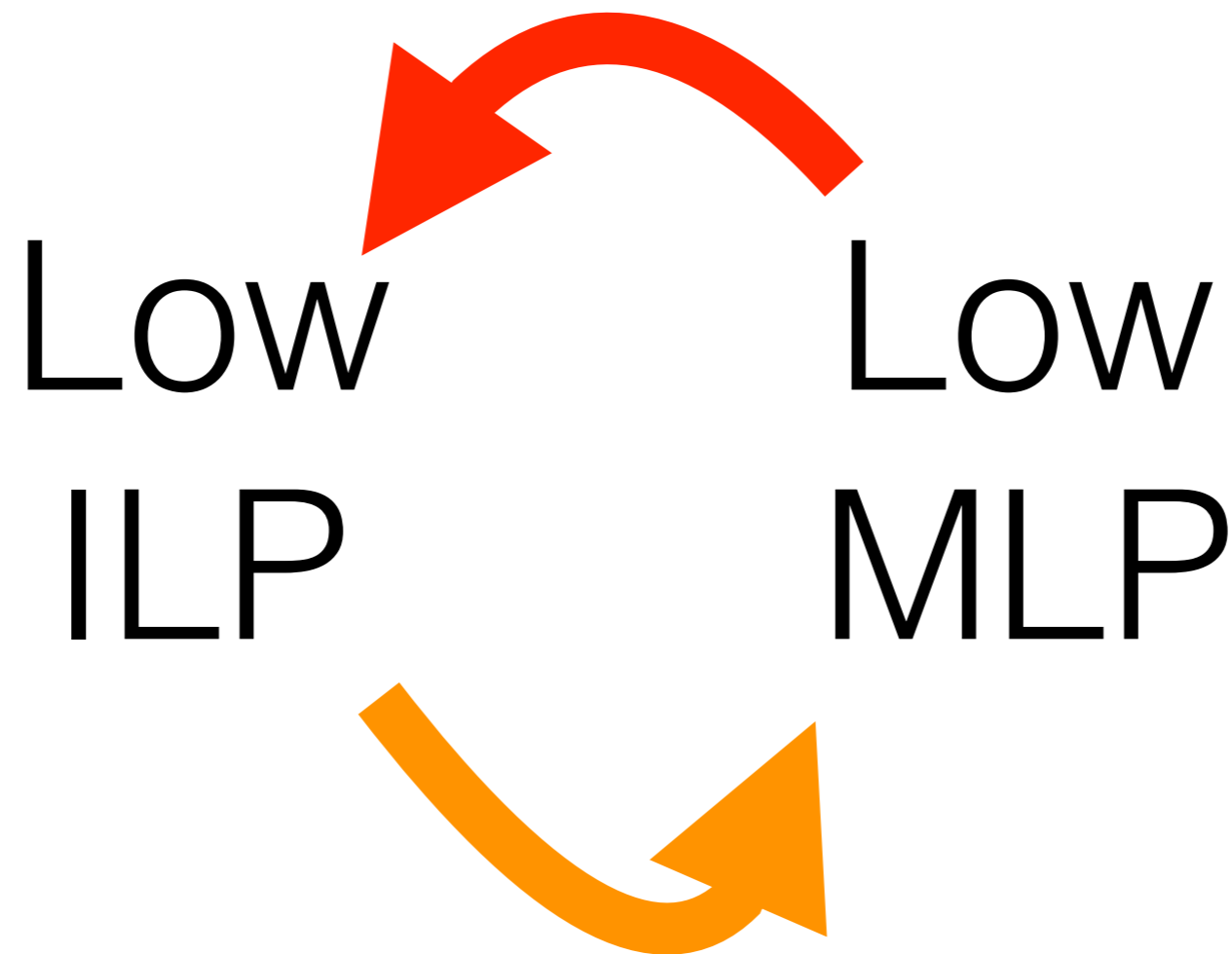
Memory Wall



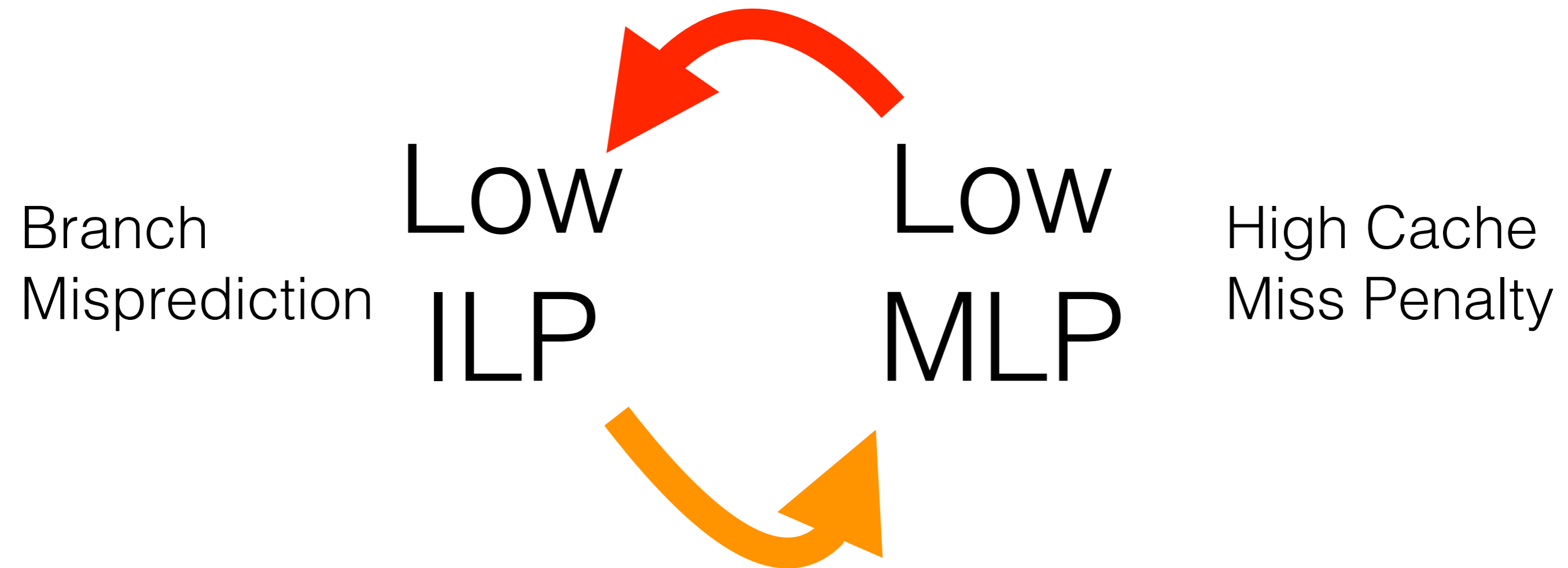
✗ Branch Misprediction

✗ Re-execution of correct path of
✗ independent tasks

ILP / MLP Vicious cycle

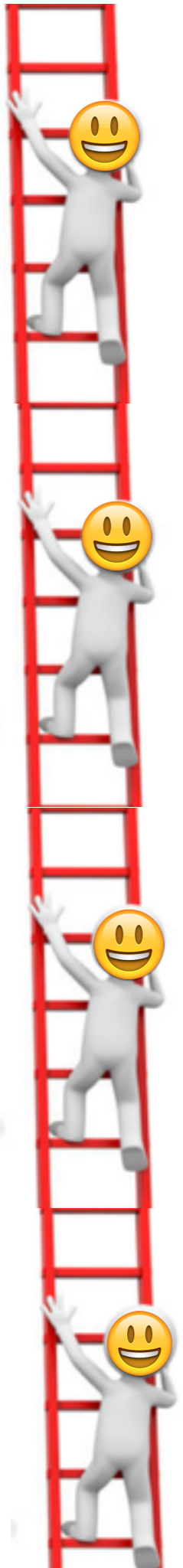


ILP / MLP Vicious cycle



Simple Alternative

- Avoid speculation for MLP - harness Request Level Parallelism (RLP)
- Tasks pipelined on one thread
- Cooperatively context switch on *likely* cache miss



Outline

- Cimple Co-Routines Overview
- Static and Dynamic Schedulers
- Related Work
- Cimple DSL and Code Generation
- Performance Evaluation
- Conclusion & Q/A

Cimple Overview

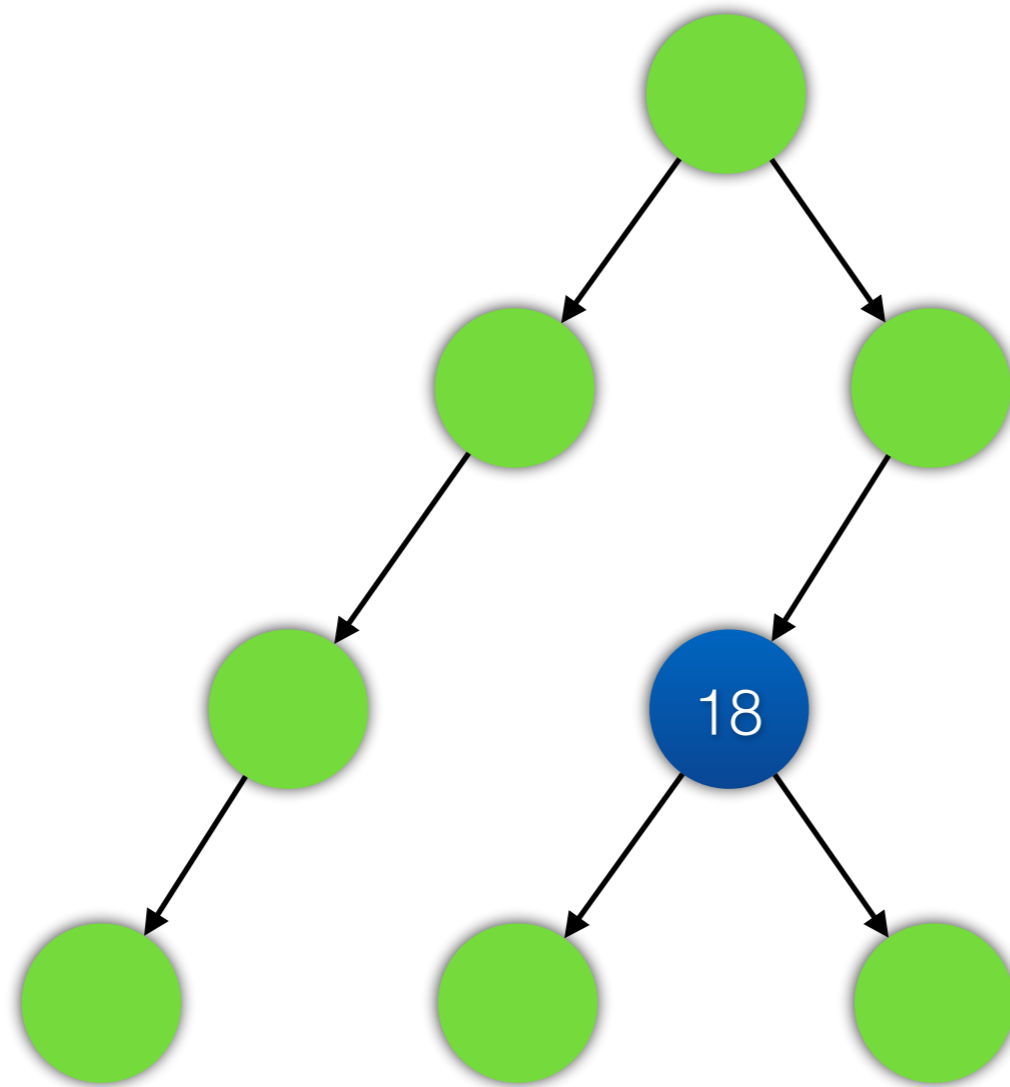
Traditional

One task per thread

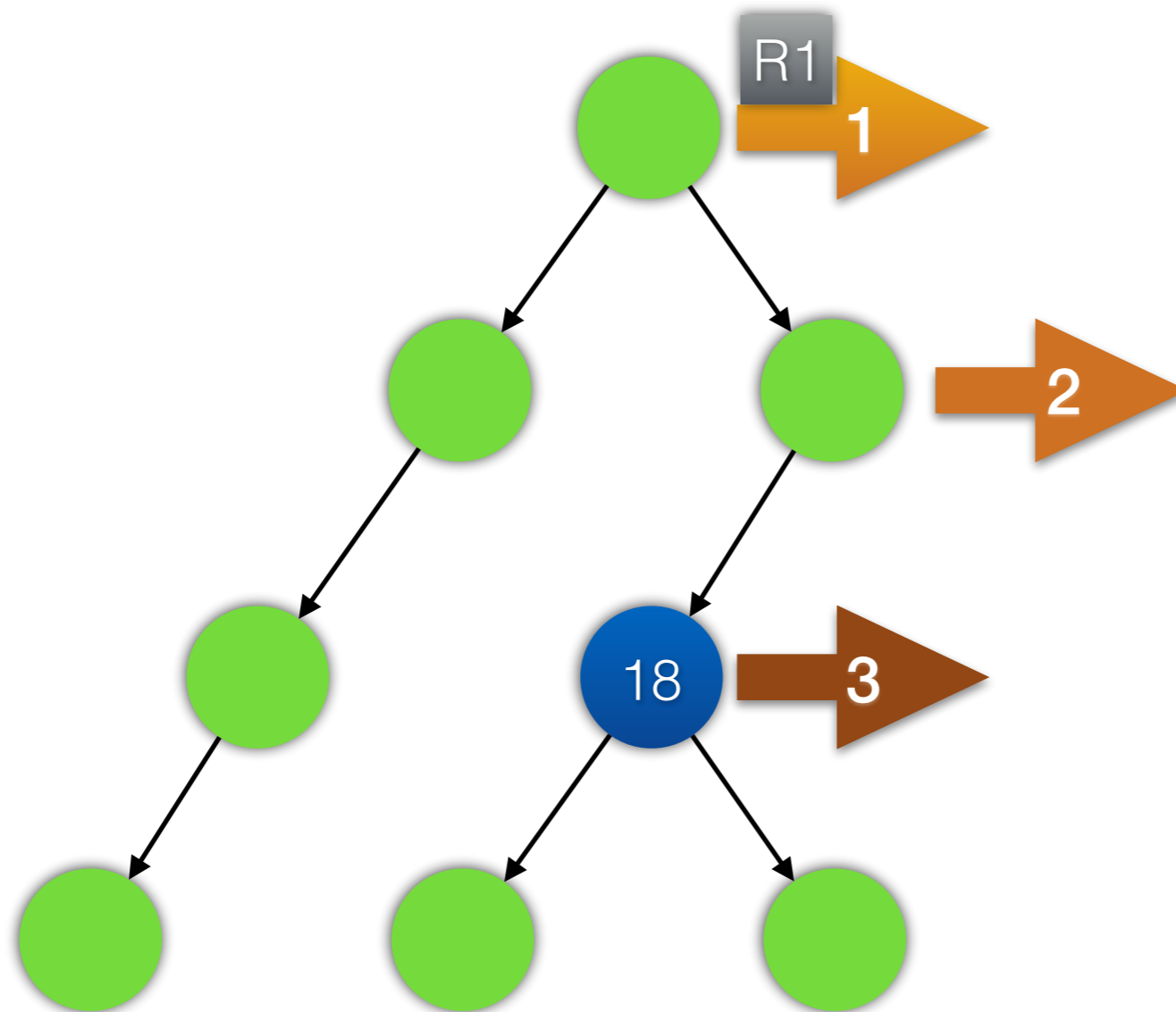
- Traditional dependence chain
- Request 1 executed to completion



Binary Tree



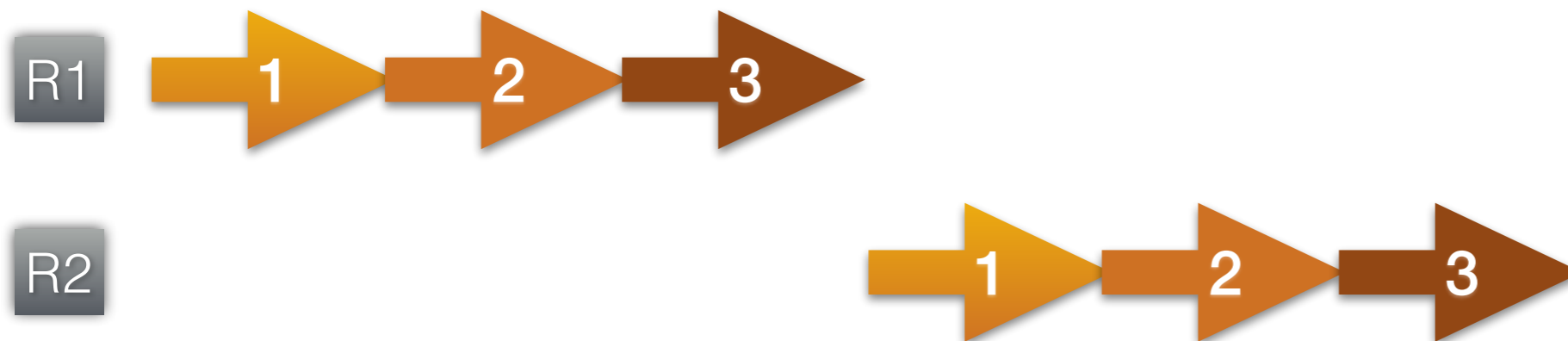
Binary Tree



Traditional

One task per thread

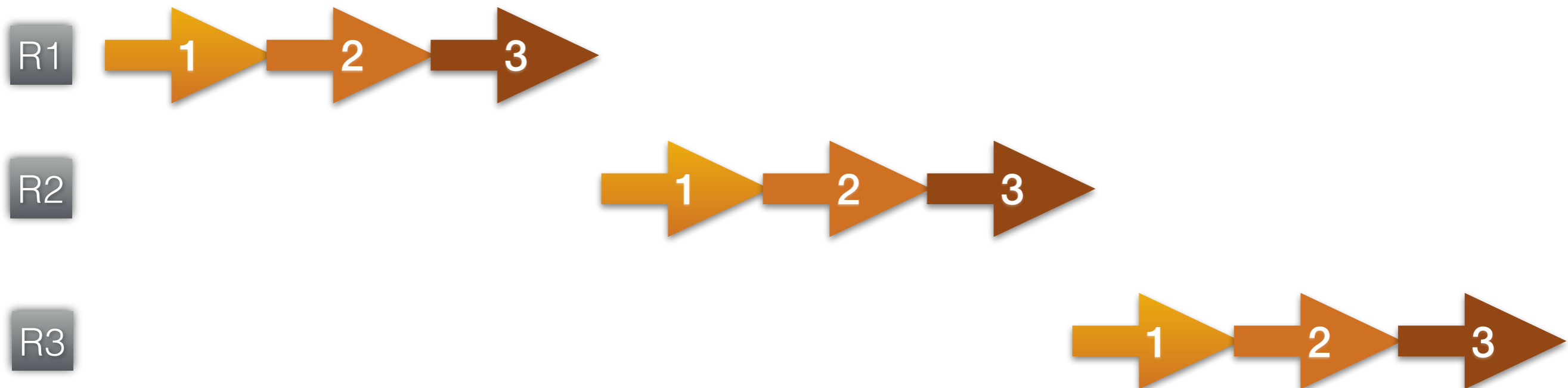
- Traditional dependence chain
- When Request 1 is complete, start Request 2



Traditional

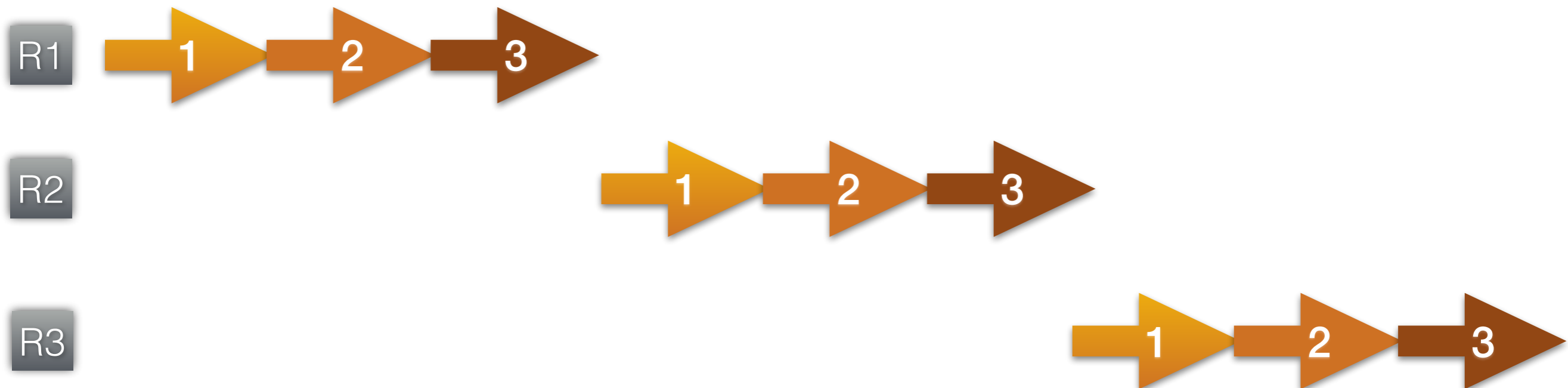
One task per thread

- Traditional dependence chain
- When Request 2 is complete, start Request 3



Traditional Limited HW Reordering

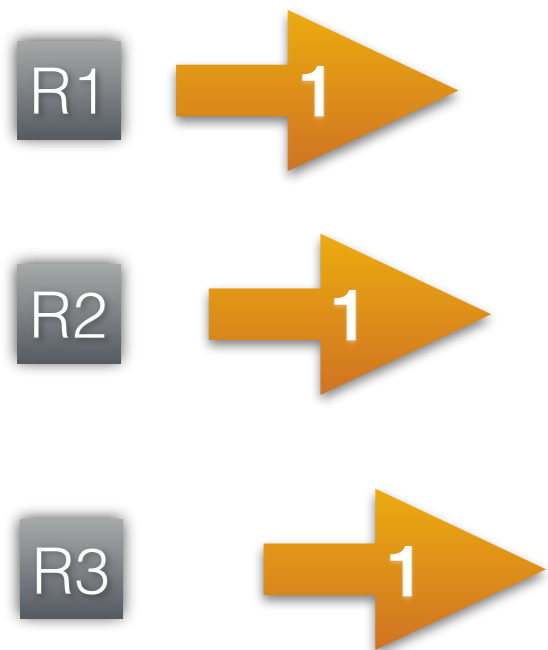
- Traditional dependence chain
- HW out-of-order execution only if predictable&short



Cimple Co-routines

Co-operative Scheduling

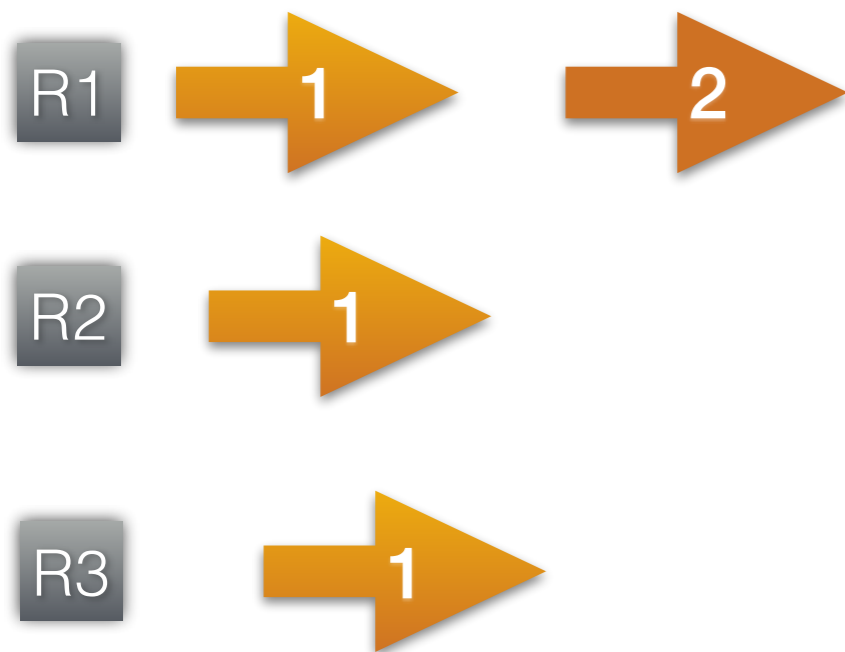
- Voluntary context switches after memory access



Cimple Co-routines

Co-operative Scheduling

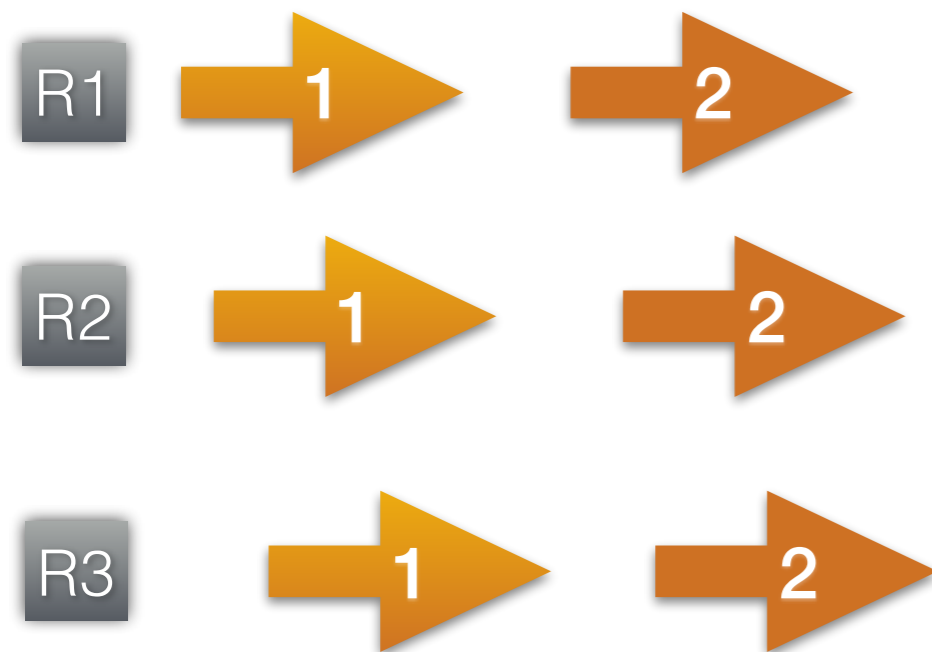
- Voluntary context switches after memory access
- No wait for completion - assume latency is hidden



Cimple Co-routines

Co-operative Scheduling

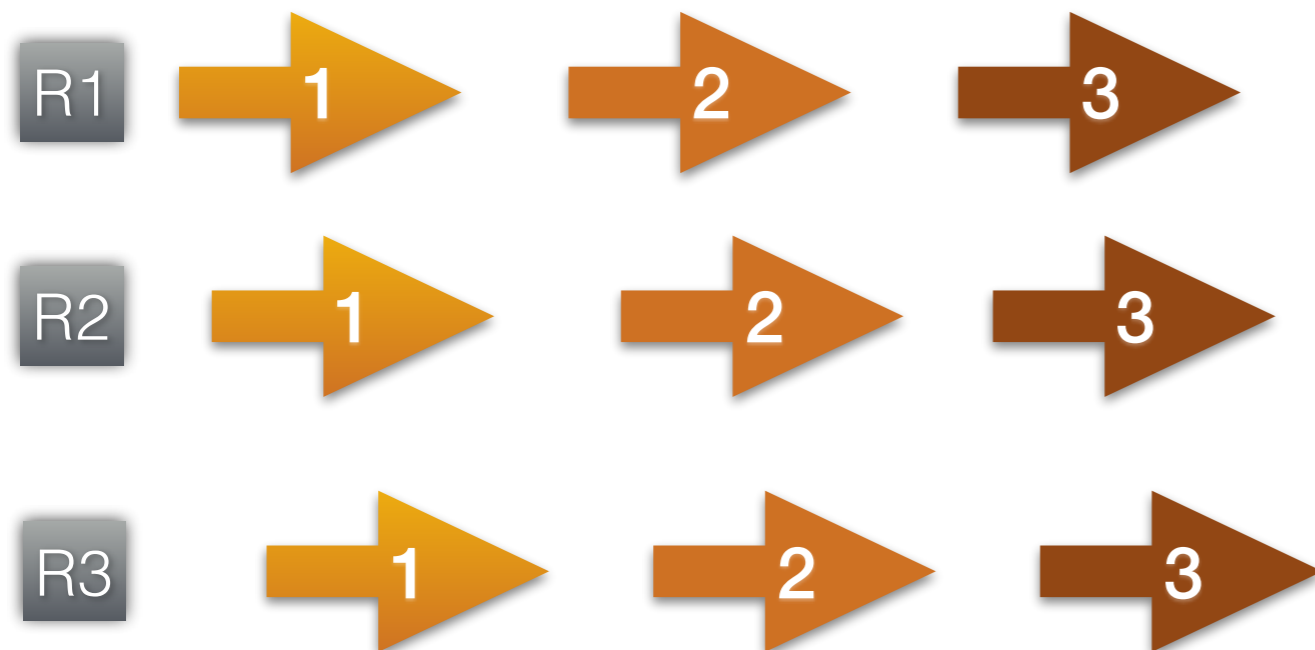
- Voluntary context switches after memory access
- No wait for completion - assume latency is hidden



Cimple Co-routines

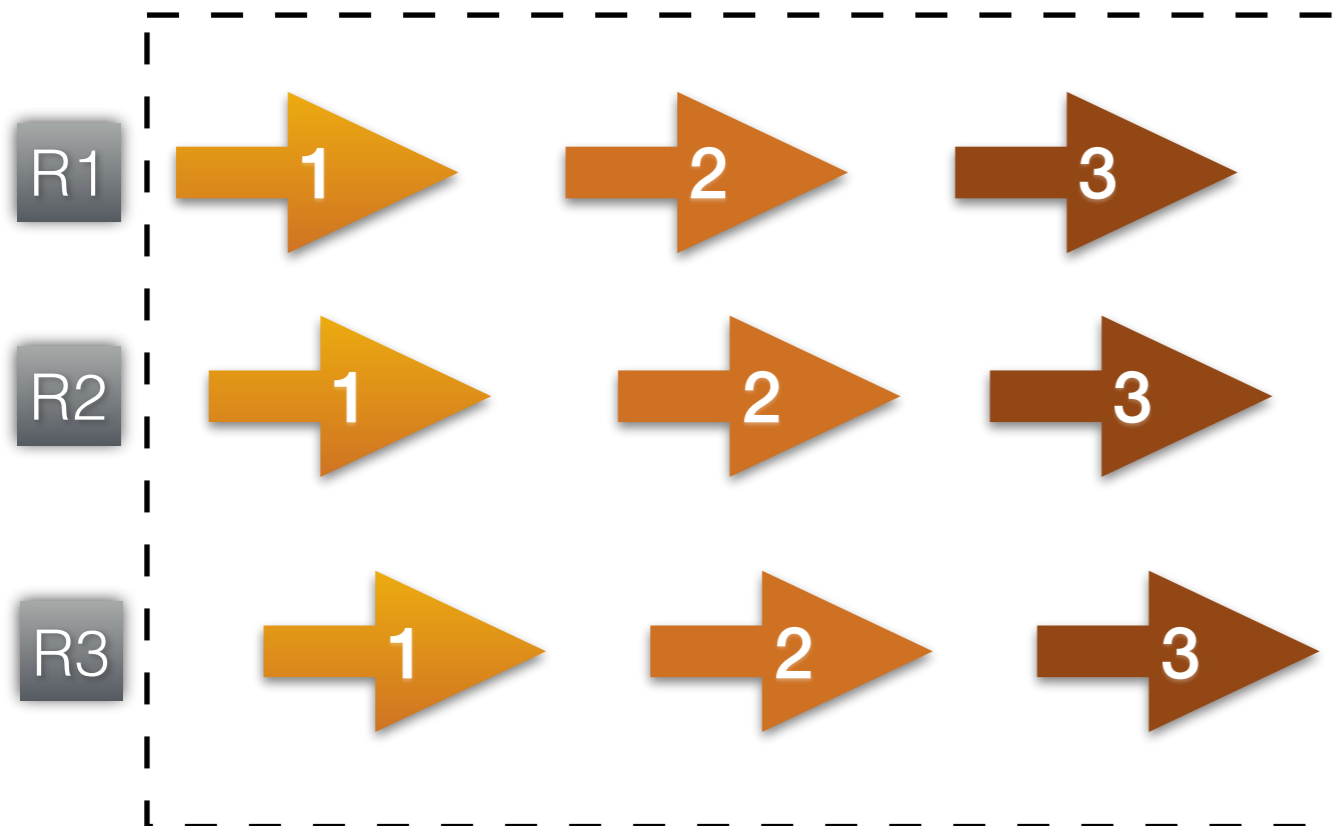
Co-operative Scheduling

- Voluntary context switches after memory access
- Mark explicitly with **yield**



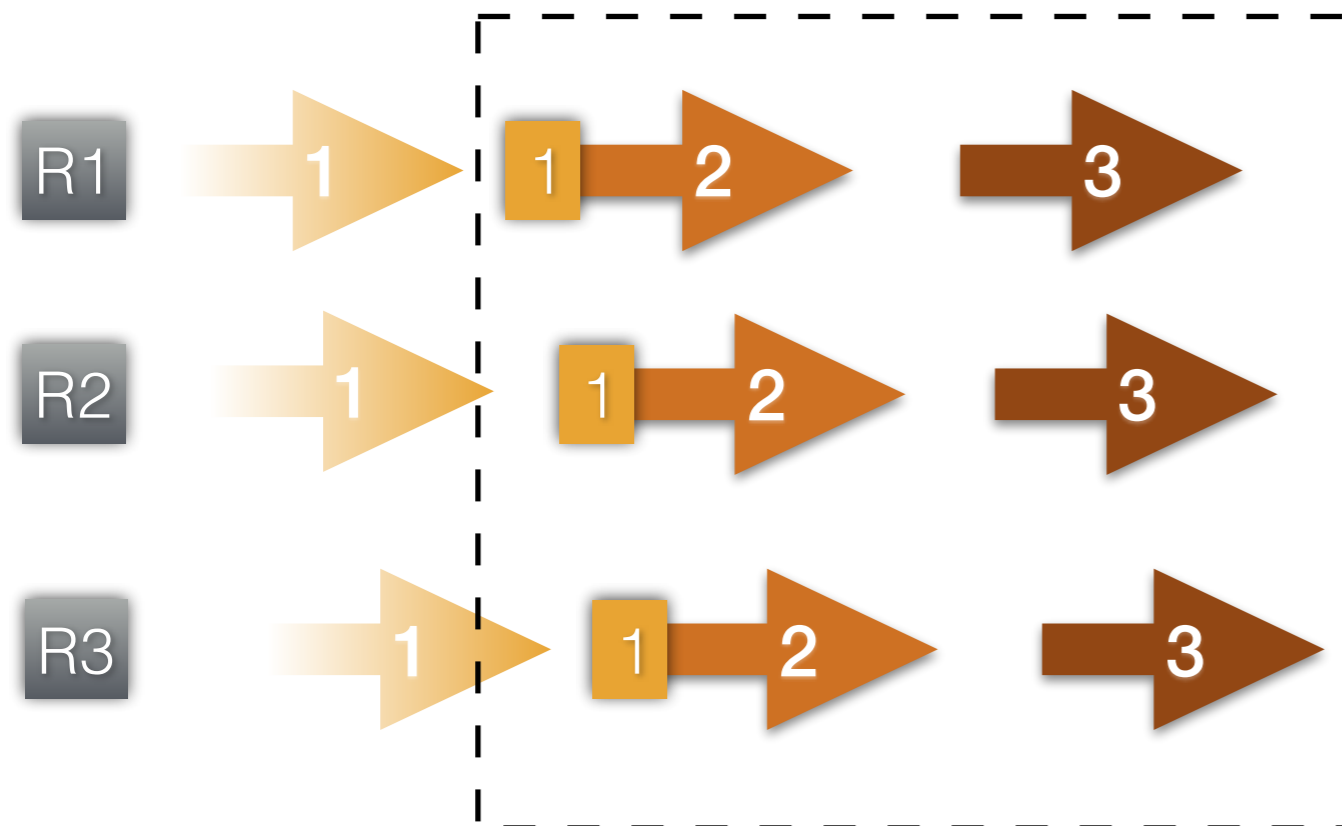
Co-routine Yield

- Mark voluntary context switches with **yield**
- Must fit all in instruction window



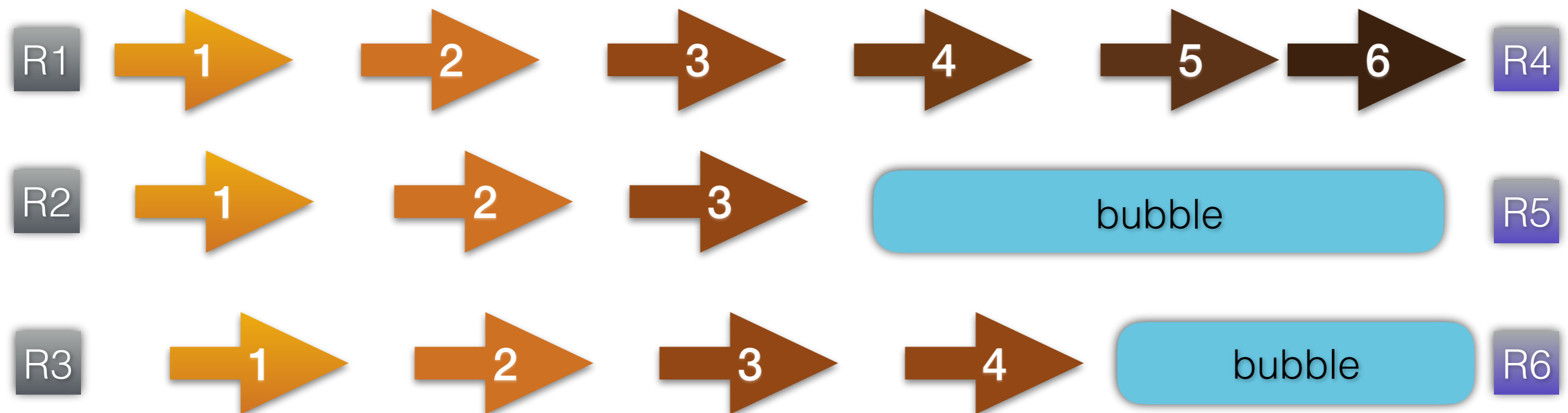
Co-routine + Prefetch

- **Prefetch** - Overlaps loads and computation
- More requests fit the instruction window



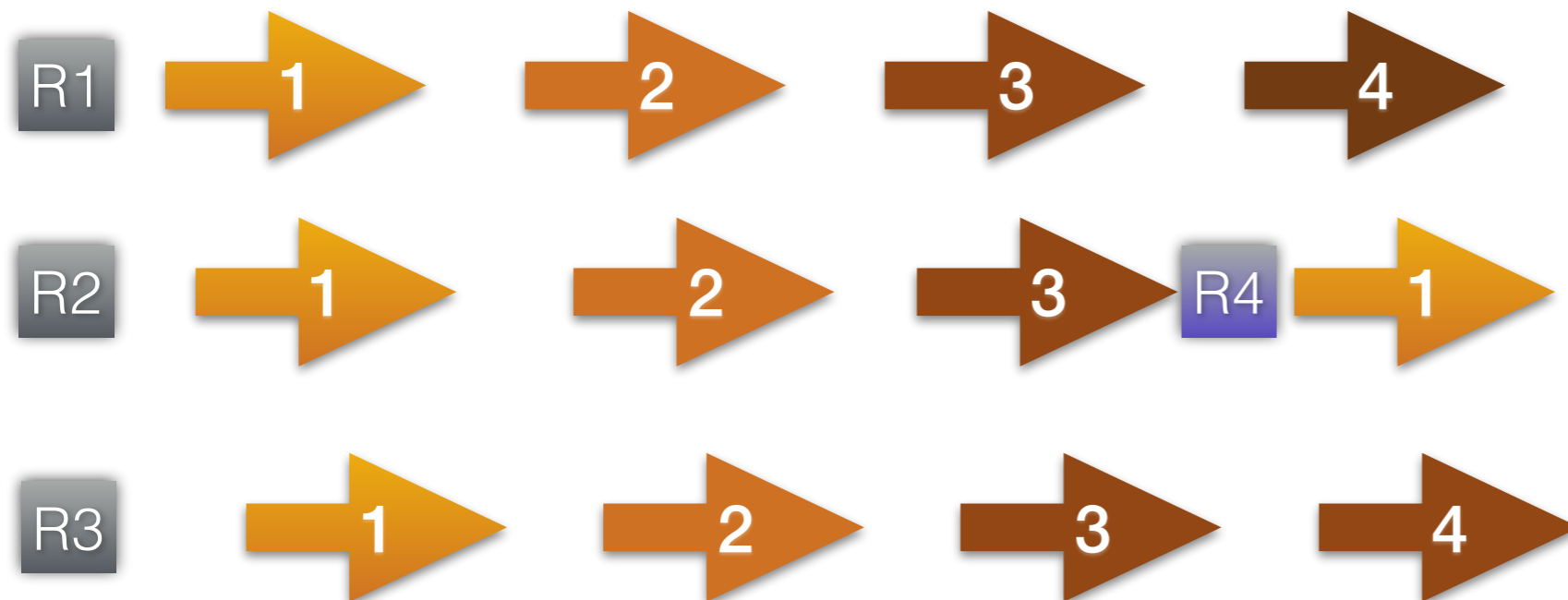
Co-routine Static Scheduling

- Execute a group at a time
- Wait until all tasks complete



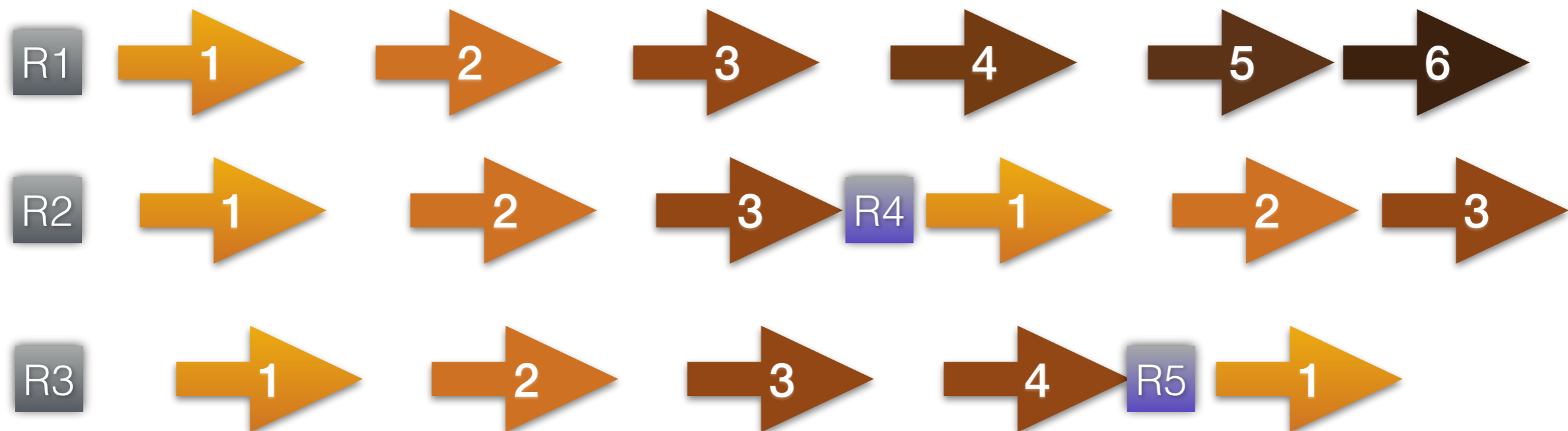
Co-routine Dynamic Scheduling

- Refill one task at a time
- Refill **R4** as soon as **R2** completes



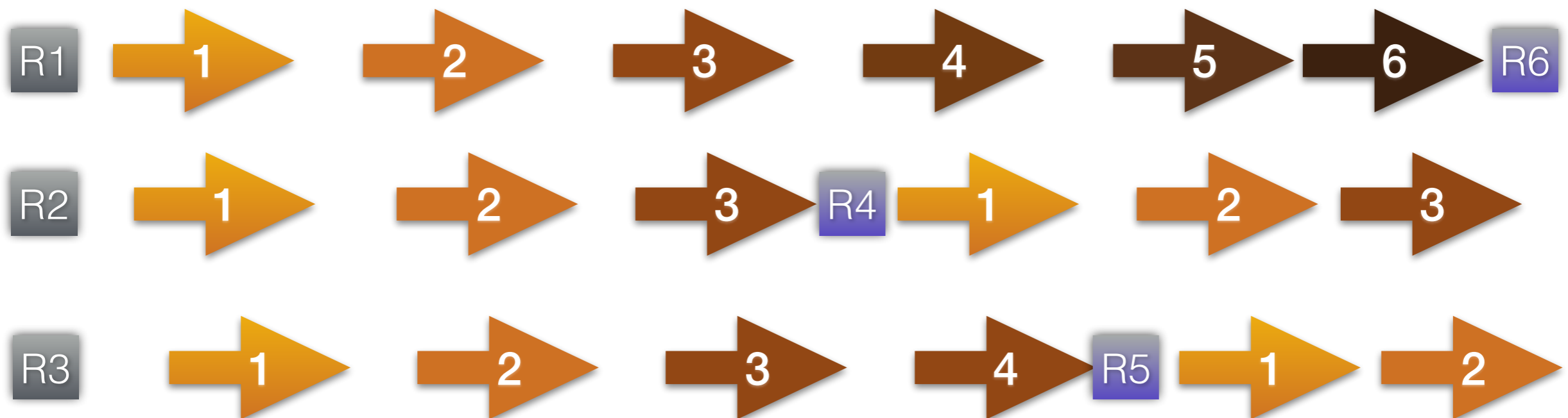
Co-routine Dynamic Scheduling

- Refill one task at a time
- Refill **R5** as soon as **R3** completes



Co-routine Dynamic Scheduling

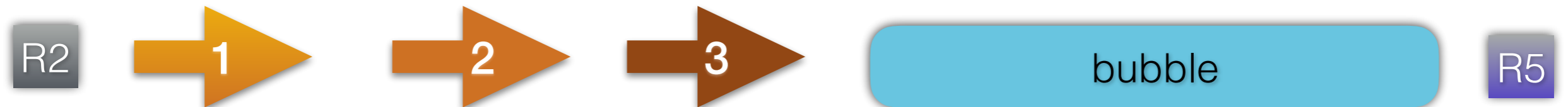
- Refill one task at a time
- Refill **R6** as soon as **R1** completes



Static vs Dynamic

- Is Dynamic always better?

Static

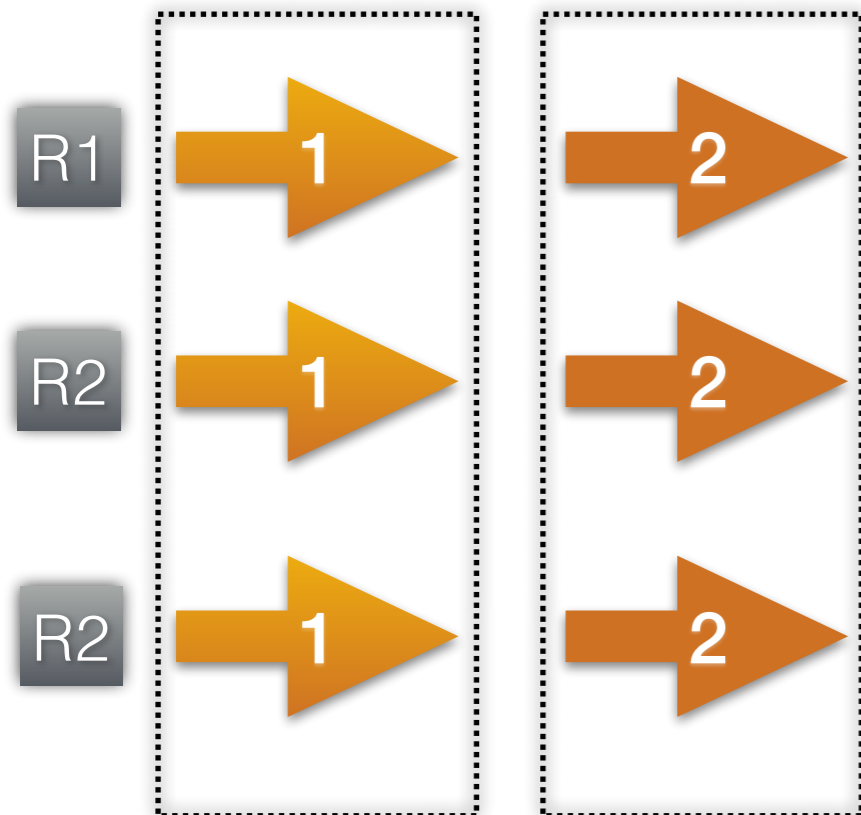


Dynamic



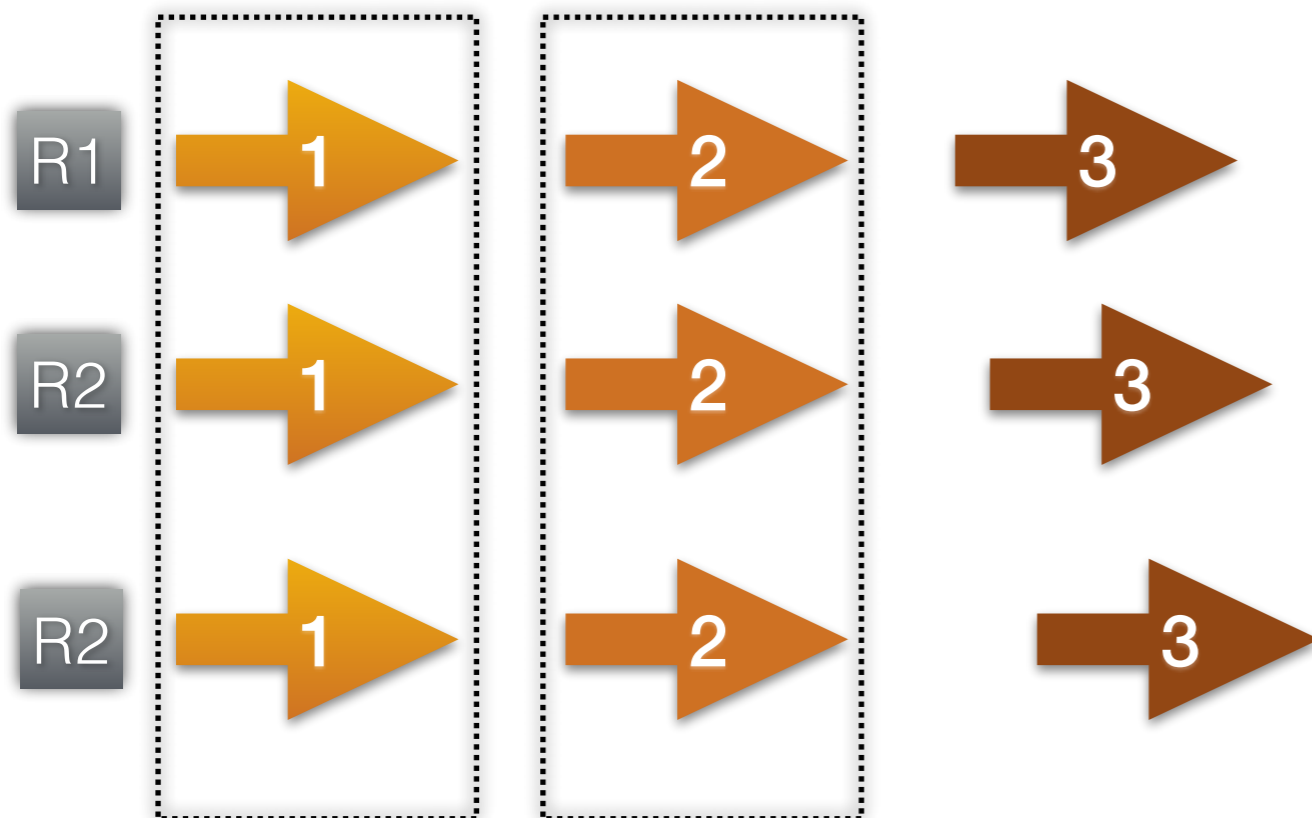
Co-routine Vectorization

- Static Scheduling






Co-routine Vectorization

- Hybrid Static+Dynamic Scheduling

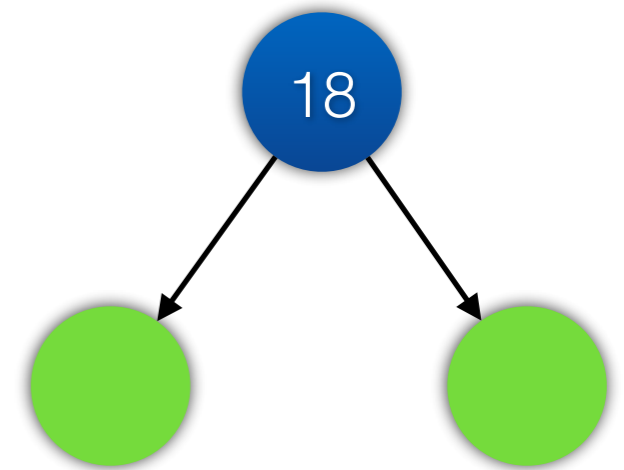


Three Keys to High MLP

-  Cooperative scheduling of co-routines
Yield at memory requests
-  Non-blocking loads overlap with computation
Prefetch avoids instruction window overflow
-  Branch misprediction penalty minimized
If/Switch grouping, and *branchless* code

Binary Tree Lookup

```
node* BinaryTree::find(node* n, Key key) {  
    while (n) {  
        if ( n->key == key )  
            return n;  
  
        if ( n->key < key )  
            n = n->right;  
        else  
            n = n->left;  
    }  
    return n;  
}
```



Binary Tree Hotspots

```
node* BinaryTree::find(node* n, Key key) {  
    while (n) {  
        if ( n->key == key ) // 1. cache miss  
            return n;  
  
        if ( n->key < key )  
            n = n->right;  
        else  
            n = n->left;  
    }  
    return n;  
}
```


Binary Tree Hotspots

```
node* BinaryTree::find(node* n, Key key) {  
    while (n) {  
        if ( n->key == key ) // 1. cache miss  
            return n;  
  
        if ( n->key < key ) // 2. branch  
            n = n->right;    // misprediction  
        else  
            n = n->left;  
    }  
    return n;  
}
```

Binary Tree Branchless

```
node* BinaryTree::find(node* n, Key key) {  
    while (n) {  
        if ( n->key == key )  
            return n;  
  
        n = n->child[ n->key < key ] ;  
  
    }  
    return n;  
}
```

Today: Cimple DSL for Experts

If it is fast and ugly, they will use it and curse you;
if it is slow, they will not use it.

- *David Cheriton*

[Jain, *The Art of Computer Systems Performance Analysis*]

Today: Cimple DSL for Experts

If it is fast and ugly, they will use it and curse you;
if it is slow, they will not use it.

- *David Cheriton*

- Performance critical database indices:
Replace LLVM IR builders in JIT query engines
- C++ Standard Template Library replacement

Past: Related Work

- GP: Group prefetching - [Chen et al'04]
manual *static scheduling* for hash-join
- AMAC: Asynchronous Memory Access Chaining
[Kocberber et al, VLDB'15]
manual *dynamic scheduling*

Concurrent: C++20 co_routines

- SAP Hana [Psaropoulos et al, VLDB'18]
- Microsoft SQLServer [Jonathan et al, VLDB'18]
automated *dynamic scheduling*
- Slower than manual GP!
Pretty front-end, high-overhead backend
- Dynamic schedule only, no vectorization

Binary Tree Lookup

```
node* BinaryTree::find(node* n, Key key) {  
    while (n) {  
        if ( n->key == key )  
            return n;  
  
        n = n->child[ n->key < key ] ;  
  
    }  
    return n;  
}
```

Cimple DSL: Binary Tree

```
6   While (n) .Do (
8       If ( n->key == key ) .
9       Then ( Return (n) ) .
10      Stmt ( n = n->child[n->key < key]; )
11  ) .
12  Return (n) ;
```


Cimple DSL: Binary Tree

```
6   While (n) .Do (
7       Prefetch (n) .Yield () .
8       If ( n->key == key ) .
9           Then ( Return (n) ) .
10          Stmt ( n = n->child[n->key < key]; )
11      ) .
12  Return (n) ;
```

Cimple DSL: Binary Tree

```
1 auto c = Coroutine(BST_find);
2 c.Result(node*).
3 Arg(node*, n).
4 Arg(KeyType, key).
5 Body().
6   While(n).Do(
7     Prefetch(n).Yield() .
8     If( n->key == key ) .
9     Then( Return(n) ) .
10    Stmt( n = n->child[n->key < key]; )
11  ) .
12 Return(n);
```

Co-routine State

```
Arg(node*, n) .  
Arg(KeyType, key) .
```

- **Arguments,
Variables**

```
1  struct Coroutine_BST_Find {  
2      node* n;  
3      KeyType key;
```

Co-routine State

```
c.Result (node*) .  
Arg (node*, n) .  
Arg (KeyType, key) .
```

- **Result**

```
1  struct Coroutine_BST_Find {  
2      node* n;  
3      KeyType key;  
4      node* _result;
```

Co-routine State

```
c.Result (node*) .  
Arg (node*, n) .  
Arg (KeyType, key) .
```

- Dynamic Schedule
Finite State Machine
_state

```
1  struct Coroutine_BST_Find {  
2      node* n;  
3      KeyType key;  
4      node* _result;  
5      int _state = 0;
```

Dynamic Schedule: Co-routine with **switch**



```
8      bool Step() {
9          switch(_state) {
10         case 0:
11
12             return false;
13         case 1:
14
15             return true;
16         case _Finished:
17             return true;
18     }
```

Dynamic Schedule: Co-routine with **switch**



```
8      bool Step() {
9          switch(_state) {
10             case 0:
11
12                 return false;
13             case 1:
14
15                 return true;
16             case _Finished:
17                 return true;
18         }
19     }
```

Duff's device
co-routine

Scheduler Width

- **Width** high to hide latency, low to fit state in L1

```
1  template<int  Width = 48>
4      using Next = CoroutineState_SkipList_next_limit;
5      SimplestScheduler<Width, Next>(len,
6          [&](Next* cs, size_t i) {
7          *cs = Next(&answers[i], IterateLimit,
8                  iter[i]);
9      });
```


Static Schedule: **for**

Vectorization friendly Struct-of-Arrays



```
1  bool SuperStep() {  
2      for(int _i = 0; _i < _Width ; _i++) {  
3          KeyType& k = _soa_k[_i];  
4          HashType& hash = _soa_hash[_i];  
  
7      }
```

Static Schedule: **for**

Vectorization friendly Struct-of-Arrays



```
1  bool SuperStep() {  
2      for(int _i = 0; _i < _Width ; _i++) {  
3          KeyType& k = _soa_k[_i];  
4          HashType& hash = _soa_hash[_i];
```



```
7      }  
8      for(int _i = 0; _i < _Width ; _i++) {  
9          KeyType& k = _soa_k[_i];  
10         HashType& hash = _soa_hash[_i];  
  
12     }
```

Applications

Binary Search

```
1 Arg(ResultIndex*, result).
2 Arg(KeyType, k).
3 Arg(Index, l).
4 Arg(Index, r).
5 Body().
6 While( l != r ).Do(
7     Stmts(R"""( {
8         int mid = (l+r)/2;
9         bool less = (a[mid] < k);
10        l = less ? (mid+1) : l;
11        r = less ? r : mid;
12        } )""").
13    Prefetch(&a[(l+r)/2]).Yield()
14 ).
15 Stmt( *result = l; );
```



Binary Search

```
1 Arg(ResultIndex*, result).
2 Arg(KeyType, k).
3 Arg(Index, l).
4 Arg(Index, r).
5 Body().
6 While( l != r ).Do(
7     Stmts(R"""( {
8         int mid = (l+r)/2;
9         bool less = (a[mid] < k);
10        l = less ? (mid+1) : l;
11        r = less ? r : mid;
12        } )""").
13    Prefetch(&a[(l+r)/2]).Yield()
14 ).
15 Stmt( *result = l; );
```

Binary Search

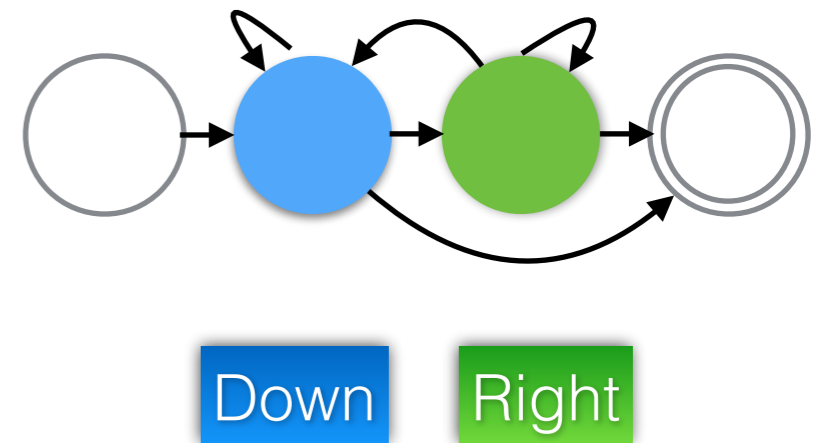
```
1 Arg(ResultIndex*, result).
2 Arg(KeyType, k).
3 Arg(Index, l).
4 Arg(Index, r).
5 Body().
6 While( l != r ).Do(
7     Stmts(R"""( {
8         int mid = (l+r)/2;
9         bool less = (a[mid] < k);
10        l = less ? (mid+1) : l;
11        r = less ? r : mid;
12        } )""").
13    Prefetch(&a[(l+r)/2]).Yield()
14 ).
15 Stmt( *result = l; );
```

Skip List Lookup

```
1 VariableInit(SkipListNode*, n, {}).
2 VariableInit(uint8, ht, {pred->height}).
3 While(true).Do(
4     While(ht > 0).Do( // down
5         Stmt( n = pred->skip[ht - 1]; ).
6         Prefetch(n).Yield().
7         If(!less(k, n->key)).Then(Break()).
8         Stmt( --ht; )
9     ).
10    If (ht == 0).Then( Return( nullptr )).
11    Stmt( --ht; ).
12    While (greater(k, n->key)).Do(
13        Stmt( pred = n; n = n->skip[ht]; ).
14        Prefetch(n).Yield().
15    ).
16    If(!less(k, n->key)).Then(
17        Return( n ));
```

Skip List Lookup

```
1 VariableInit(SkipListNode*, n, {}).
2 VariableInit(uint8, ht, {pred->height}).
3 While(true).Do(
4     While(ht > 0).Do( // down
5         Stmt( n = pred->skip[ht - 1]; ).
6         Prefetch(n).Yield().
7         If(!less(k, n->key)).Then(Break()).
8         Stmt( --ht; )
9     ).
10    If (ht == 0).Then( Return( nullptr ) ).
11    Stmt( --ht; ).
12    While (greater(k, n->key)).Do(
13        Stmt( pred = n; n = n->skip[ht]; ).
14        Prefetch(n).Yield().
15    ).
16    If(!less(k, n->key)).Then(
17        Return( n ));
```



Skip List Iteration

```
1   While( limit-- ).Do(  
2       Prefetch(n).Yield().  
3       Stmt( n = n->skip[0]; )  
4   ).  
5   Prefetch(n).Yield().  
6   Return( n->key );
```

- Pointer chasing

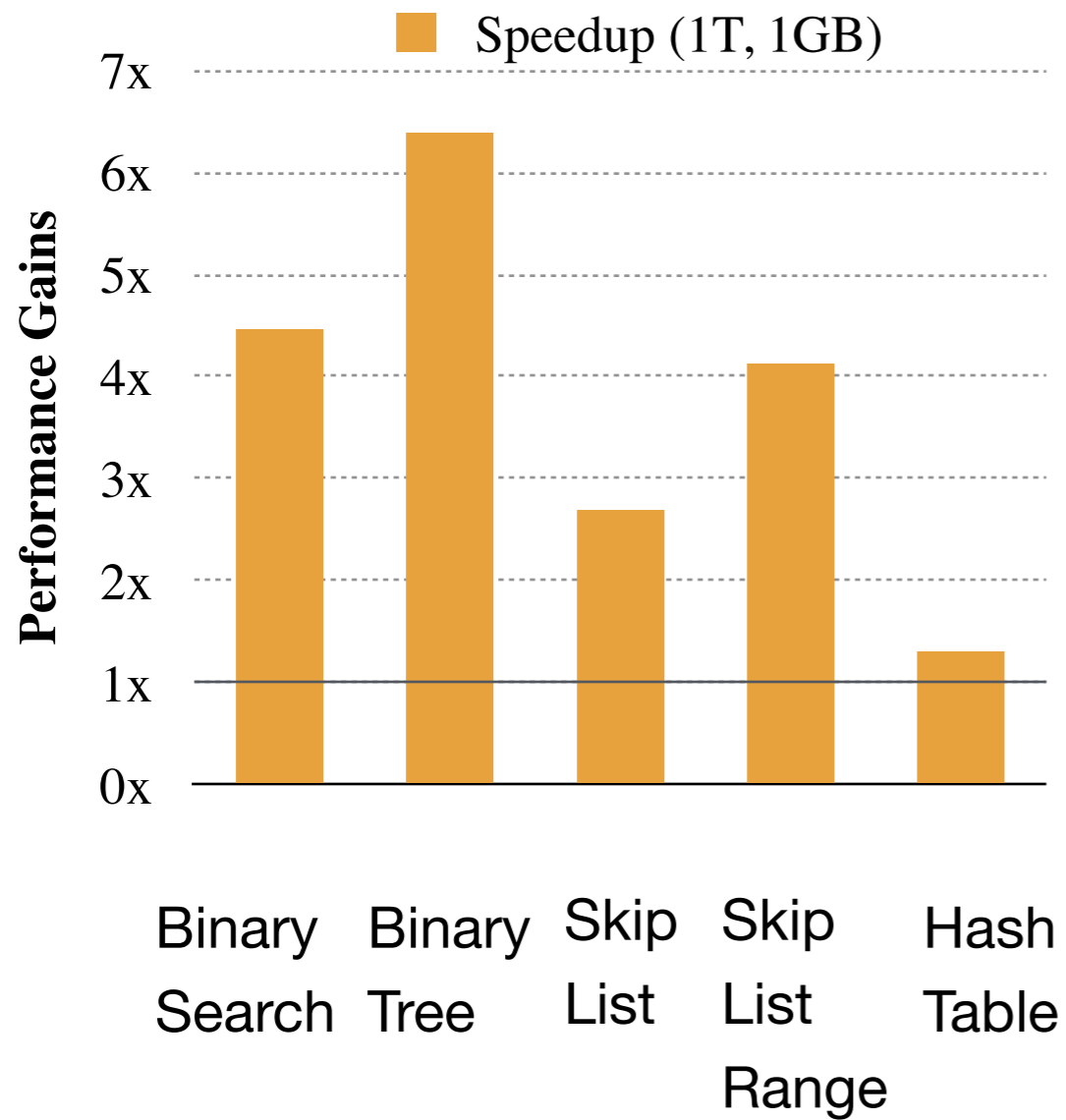
Hash Table Lookup (Linear Probing)

```
1 Result (KeyValue*) .
2 Arg (KeyType, k) .
3 Variable (HashType, hash) .
4 Body () .
5   Stmt ( hash = Murmur3::fmix(k); ) .
6   Stmt ( hash &= this->size_1; ).Yield() .
7   Prefetch( &ht[hash] ).Yield()
8   << R"" (
9     while (ht[hash].key != k &&
10           ht[hash].key != 0) {
11       hash++;
12       if (hash == size) hash = 0;
13     } )"" <<
14   Return( &ht[hash] );
```

- SIMD
- One cache line

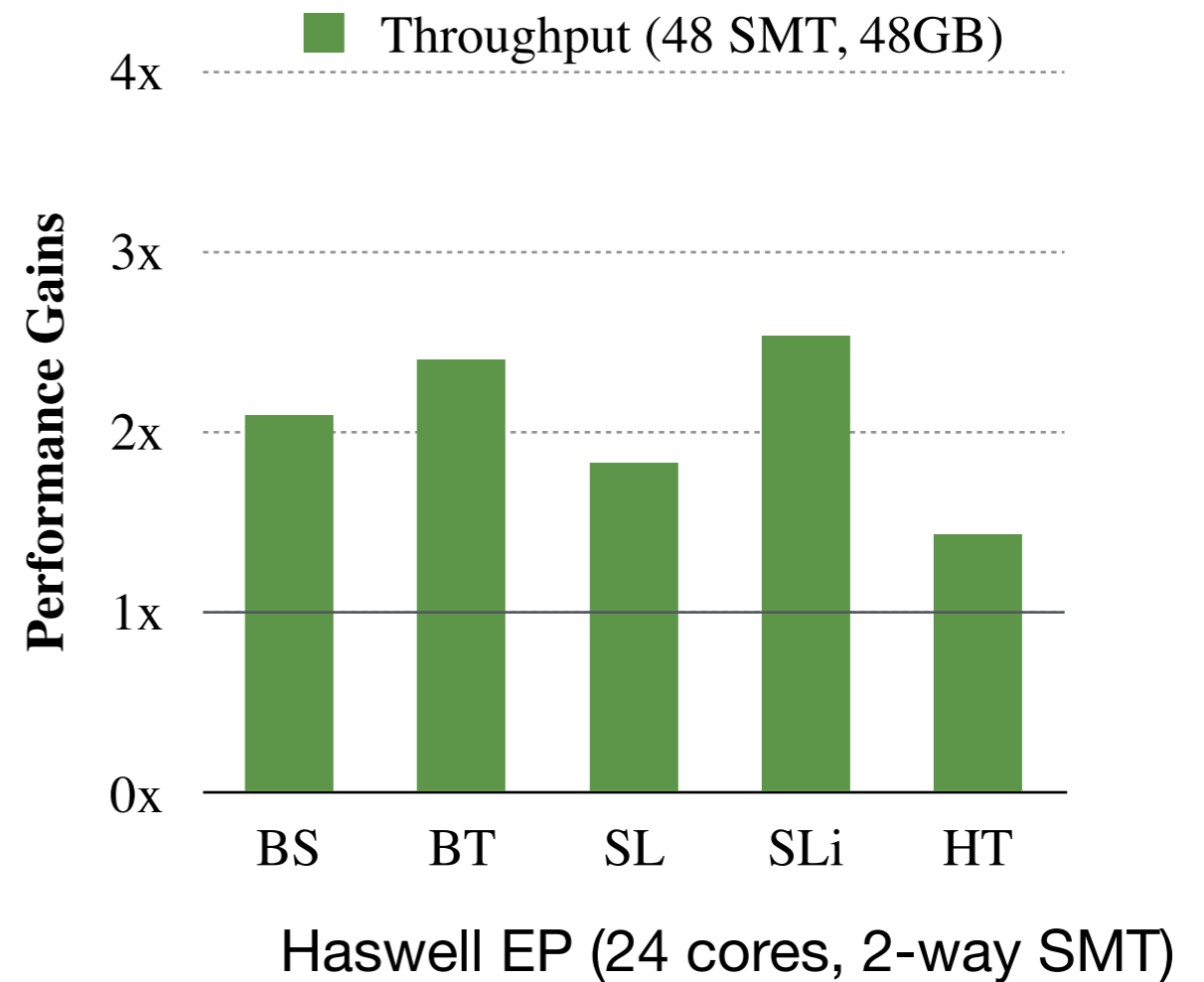
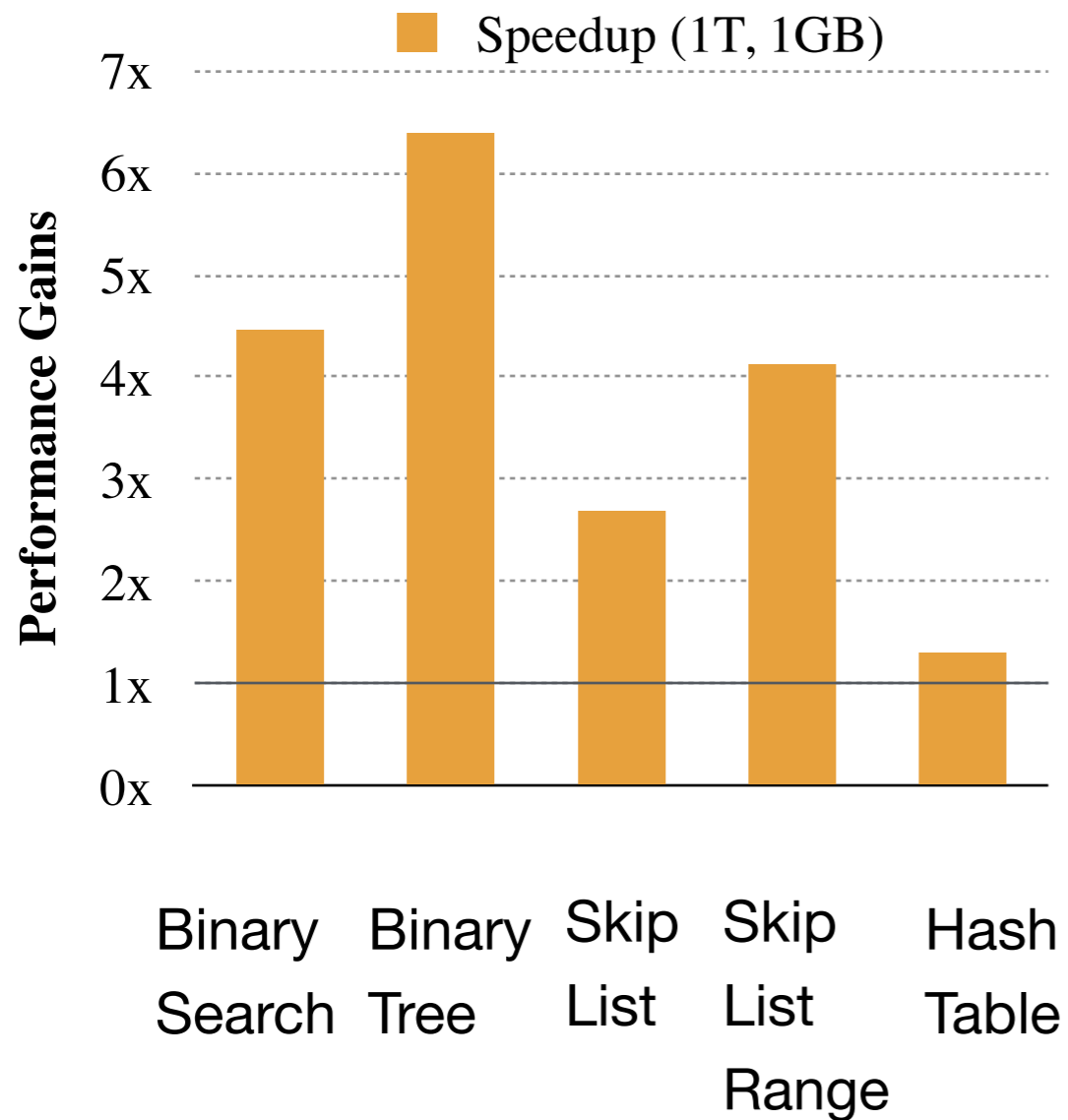
Performance Evaluation

Performance



[Default indices of
VoltDB/RocksDB]

Performance



[Default indices of VoltDB/RocksDB]

Thread Level Parallelism



Multi-core



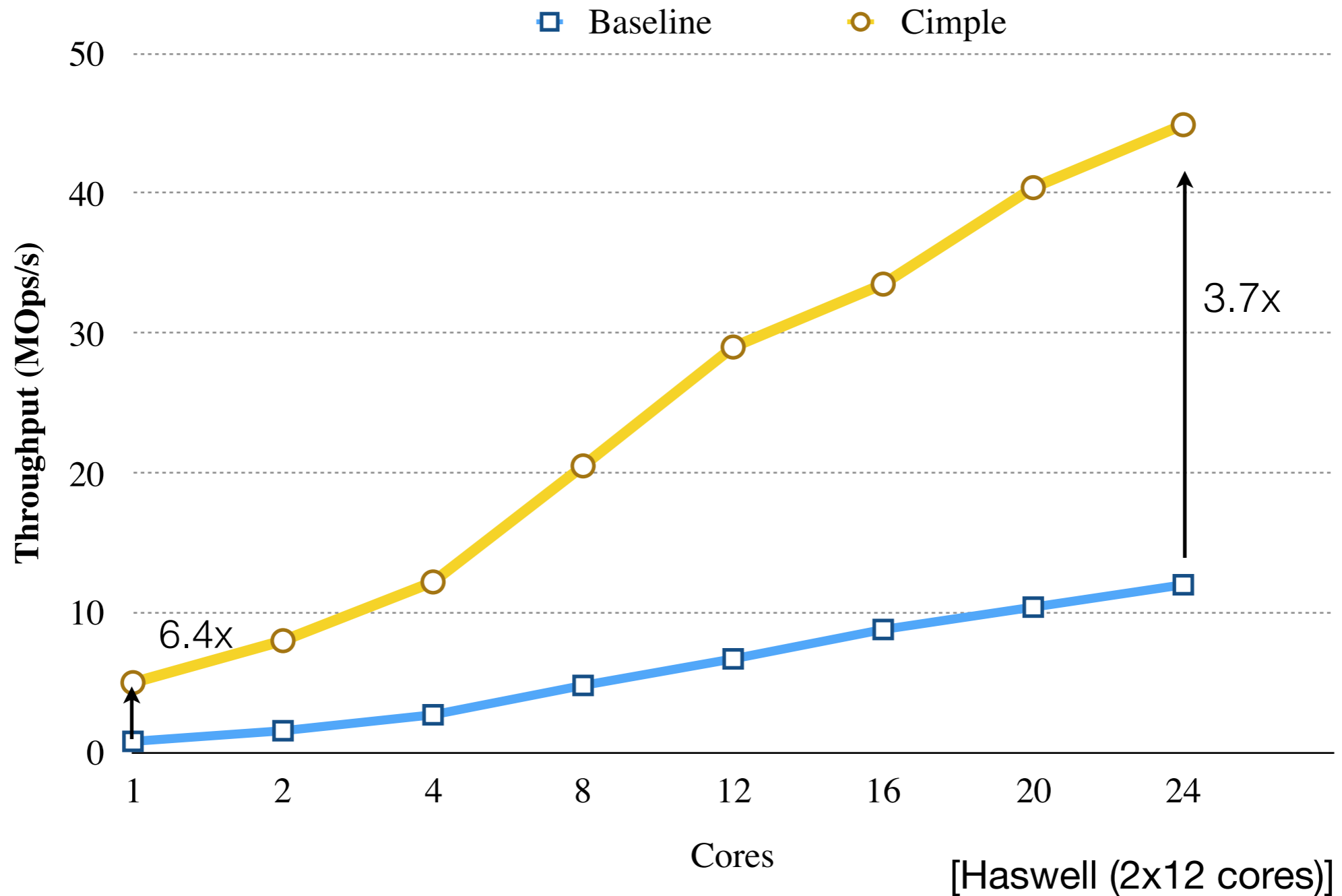
SMT hardware thread



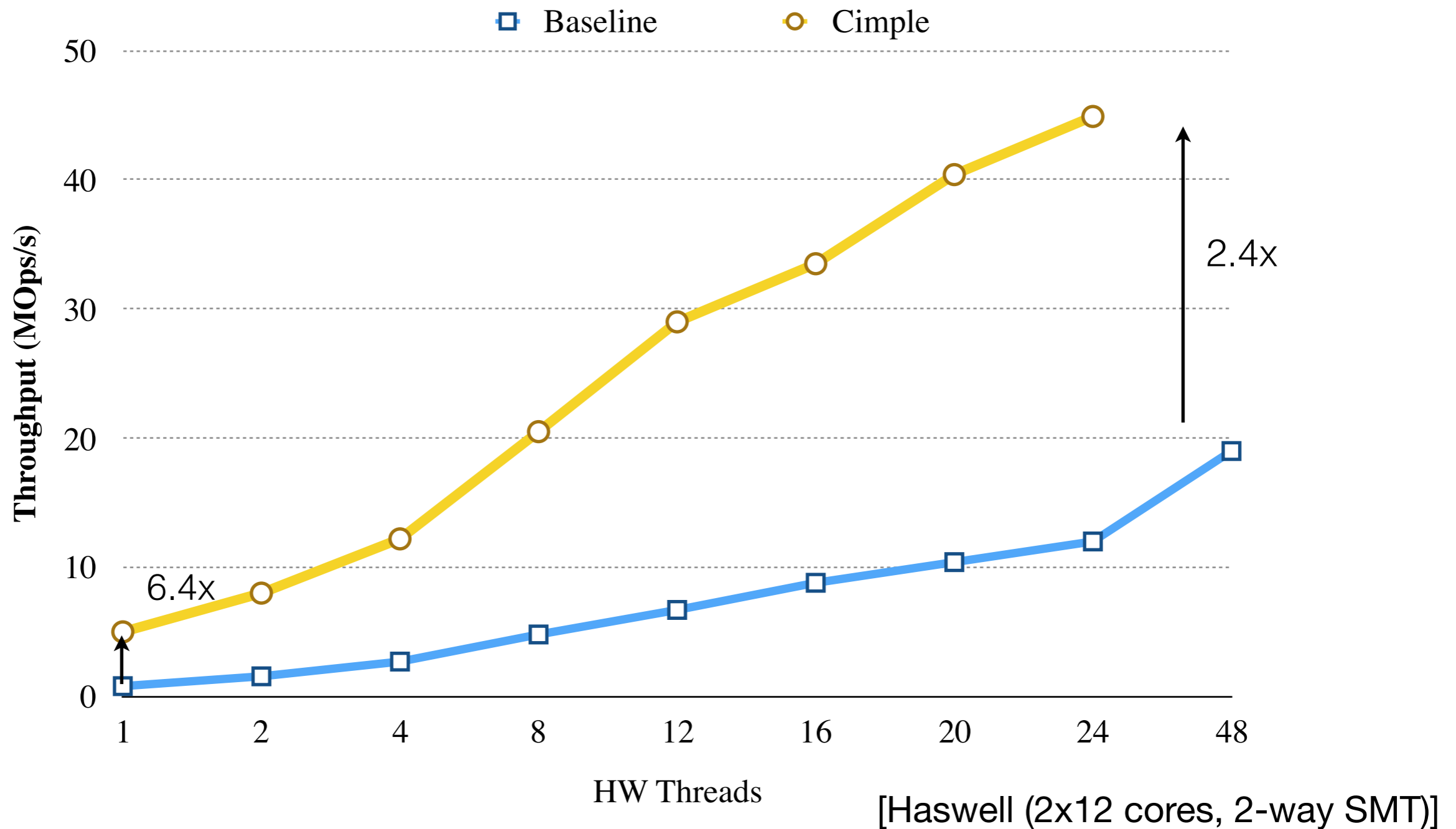
Single-thread OS context switching?

50x slower

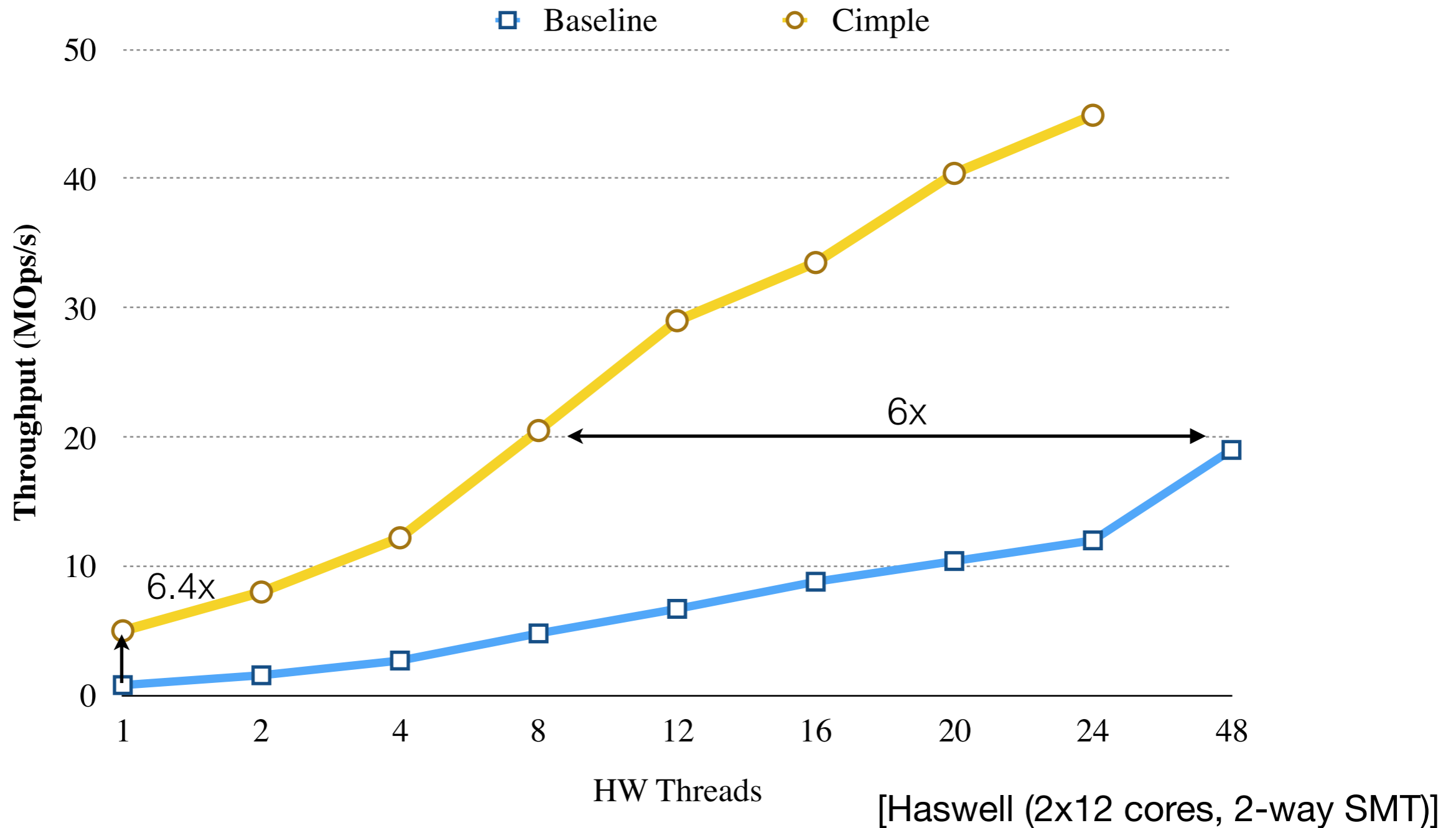
Cimple Throughput Gains on Multicore



Cimple Throughput Gains vs Hyper-threading



Cimple Throughput Gains vs Hyper-threading

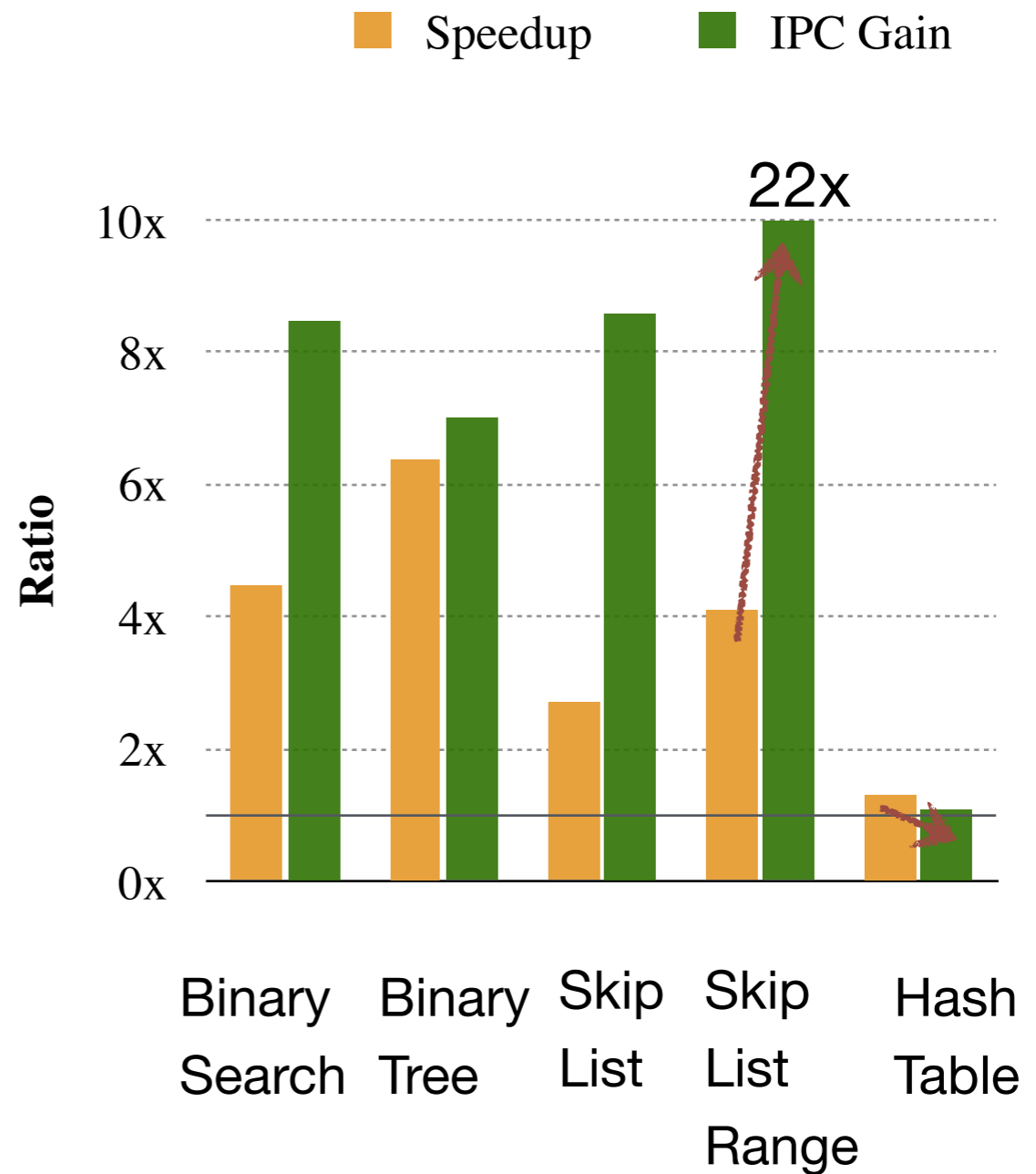


IPC Analyzed

- SkipList Range:
scheduler overhead

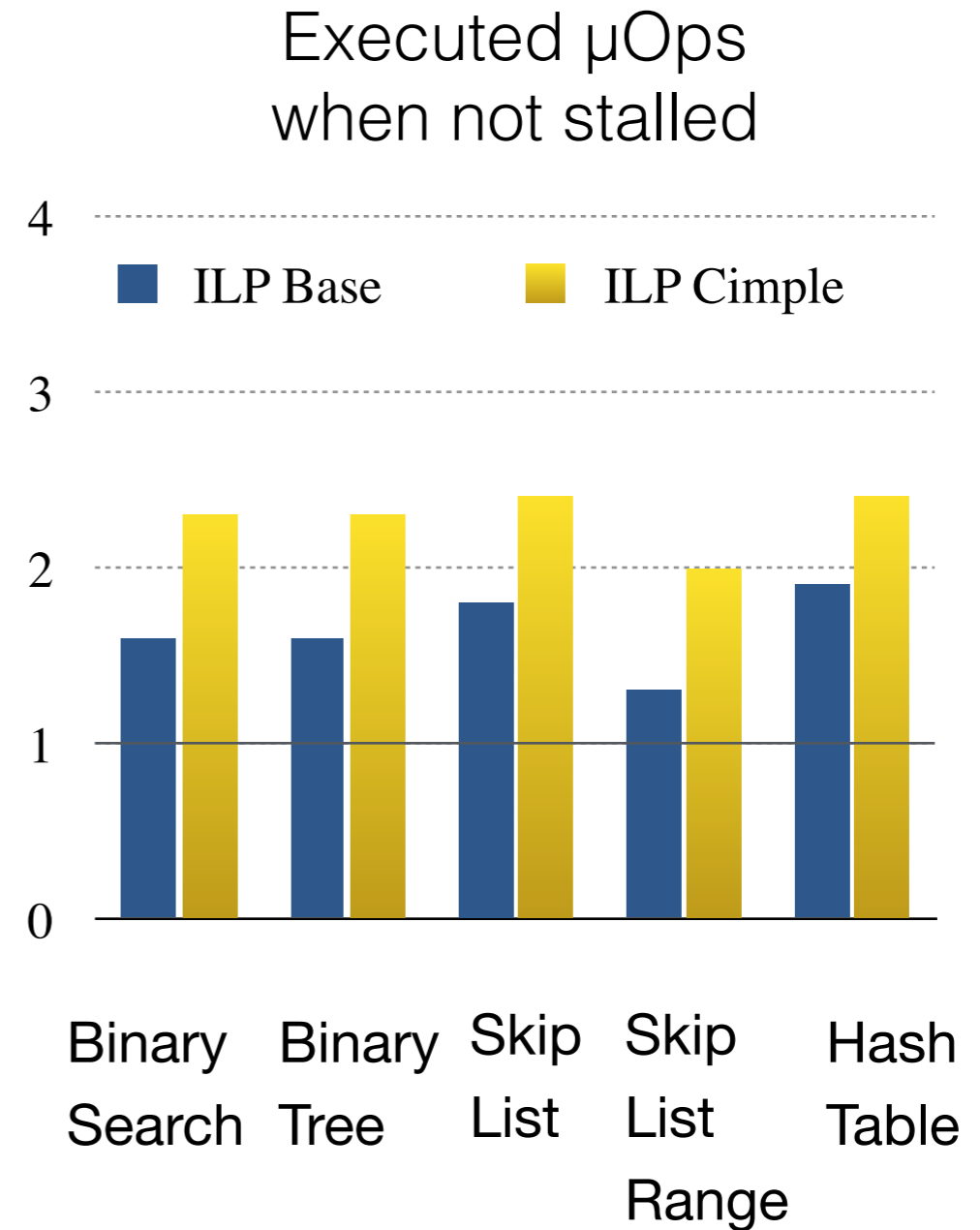
```
p = p->next;
```

- HashTable:
SIMD vectorization



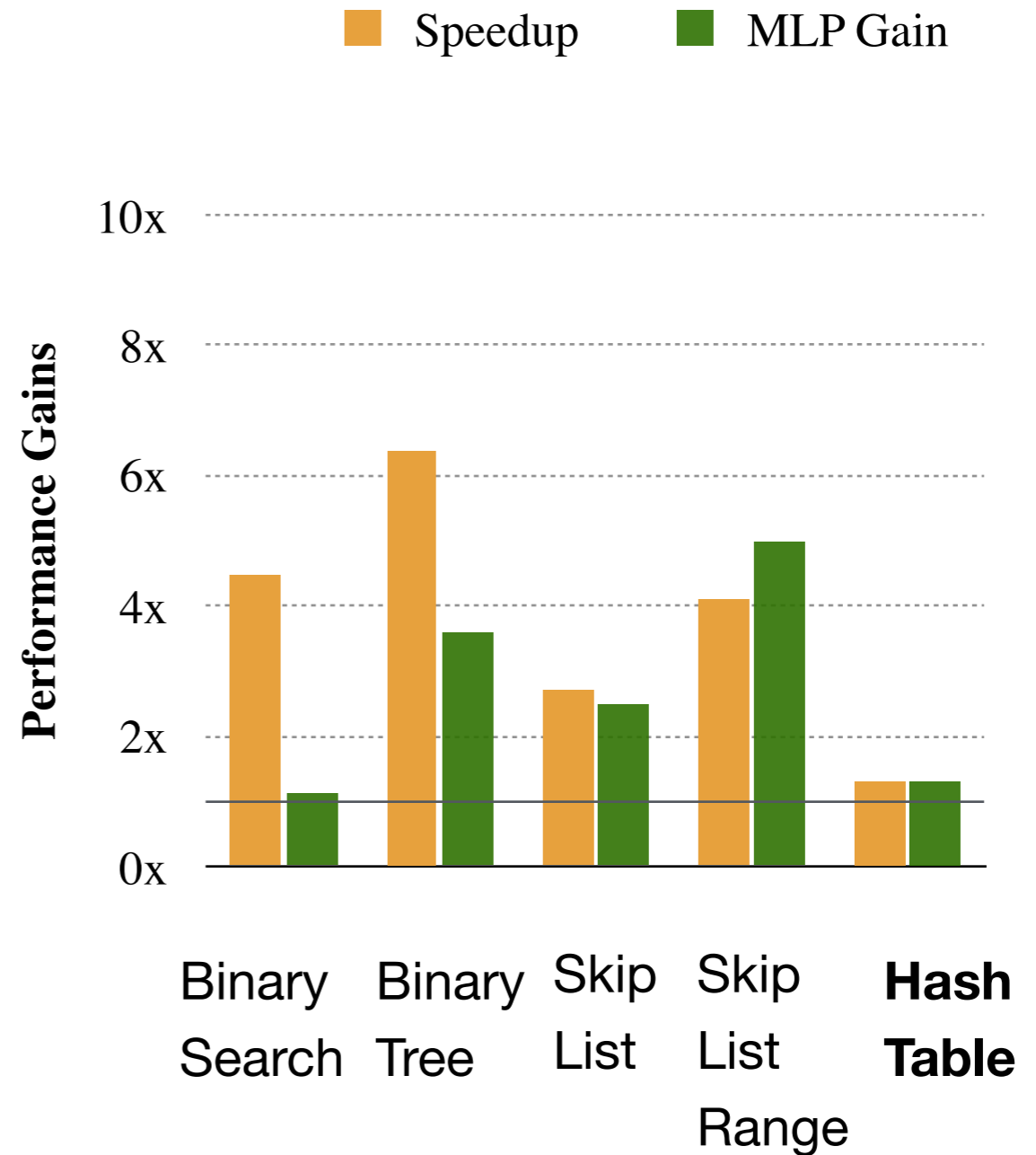
ILP Analyzed

- Uncovered more parallelism
- Absorbed scheduler overhead



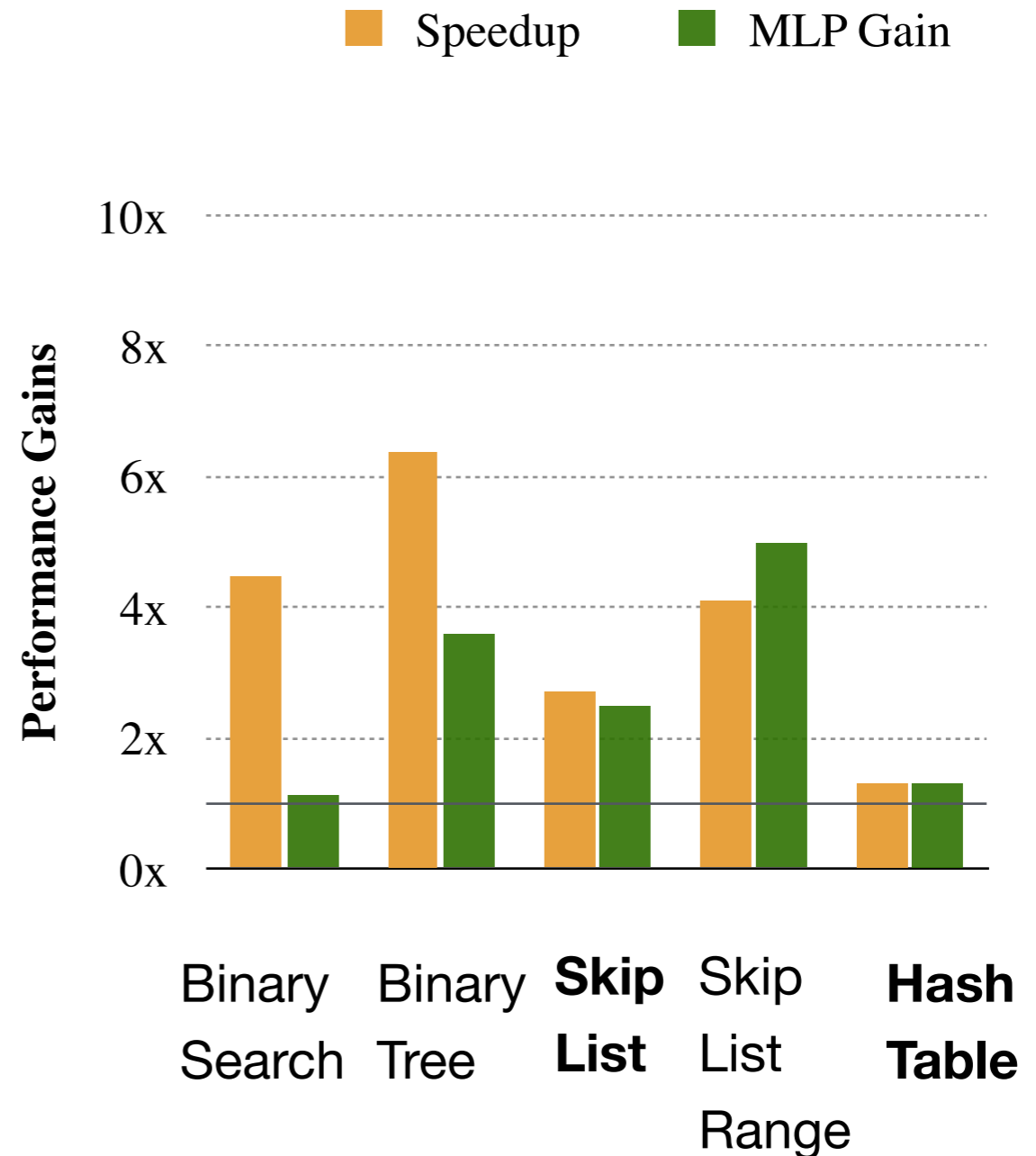
MLP Analyzed

- HashTable -
OoO HW extracts MLP

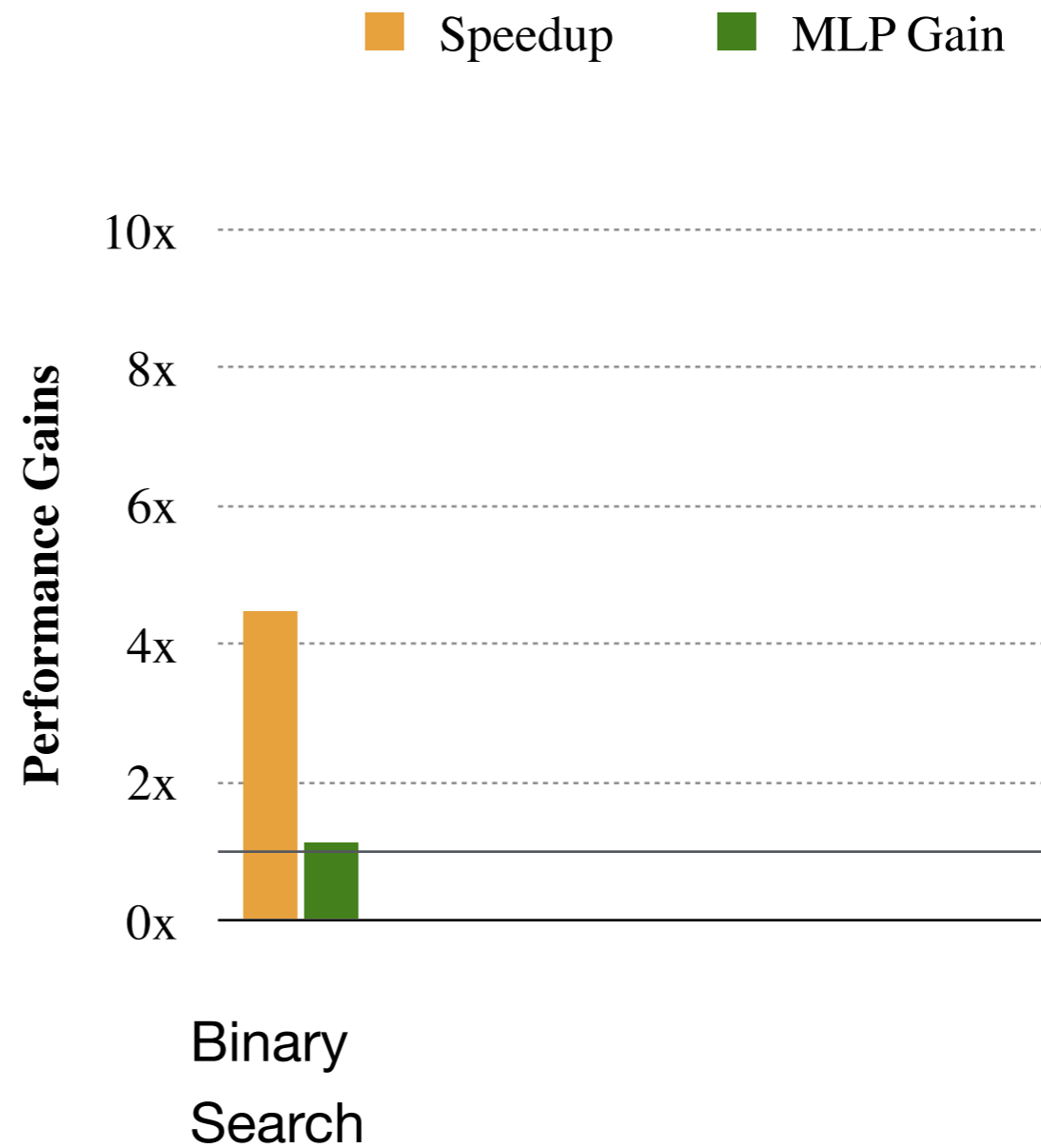
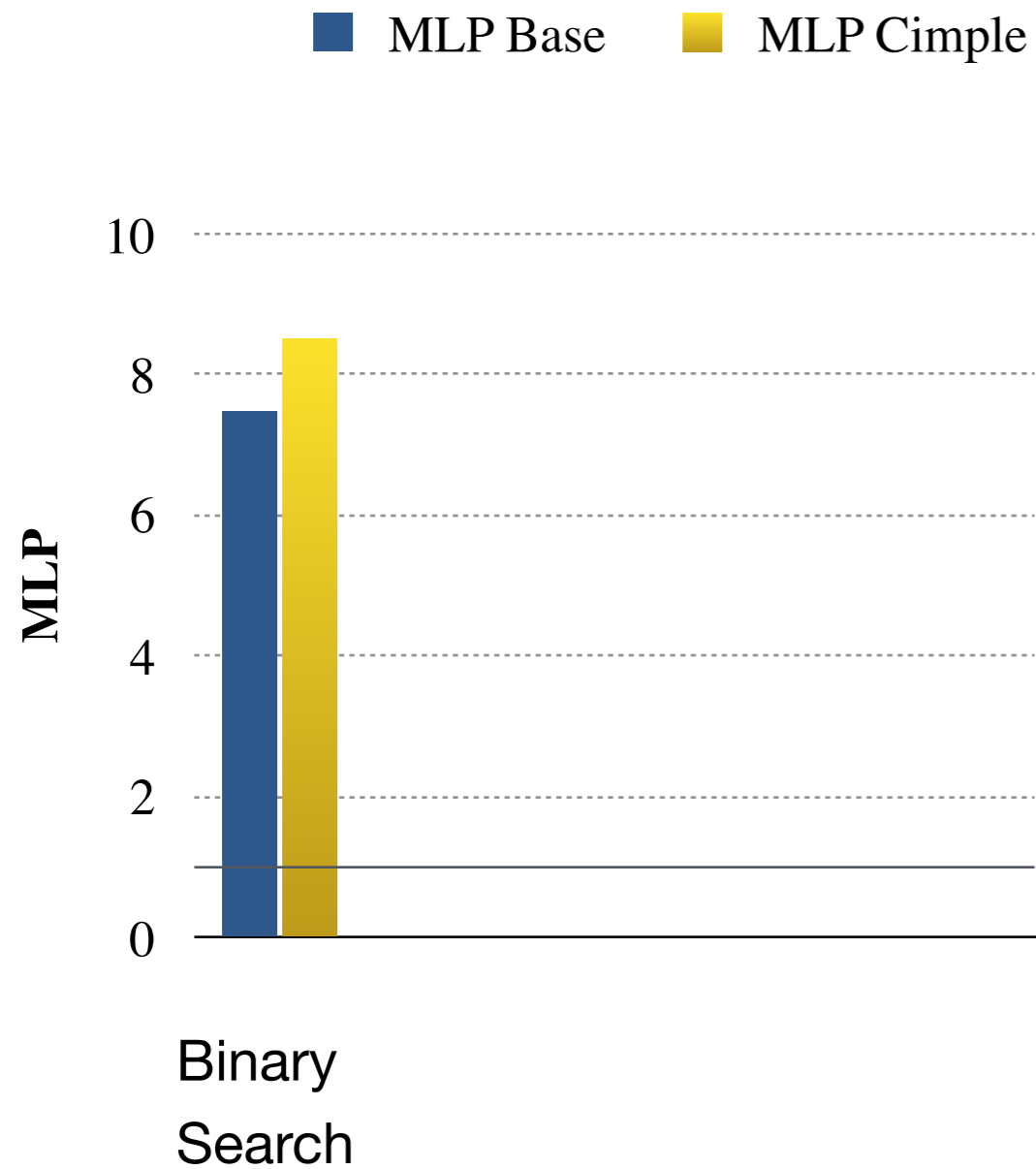


MLP Analyzed

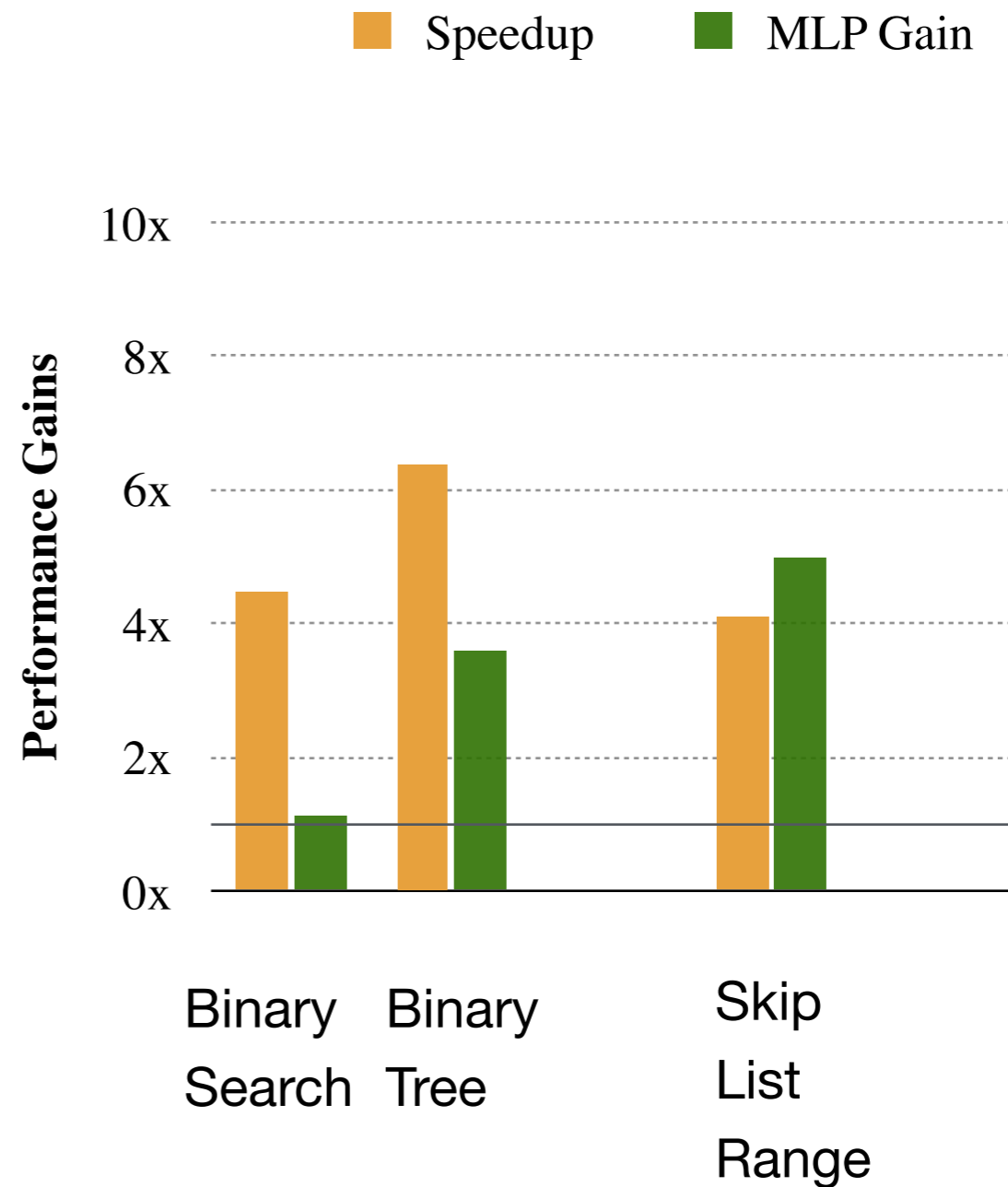
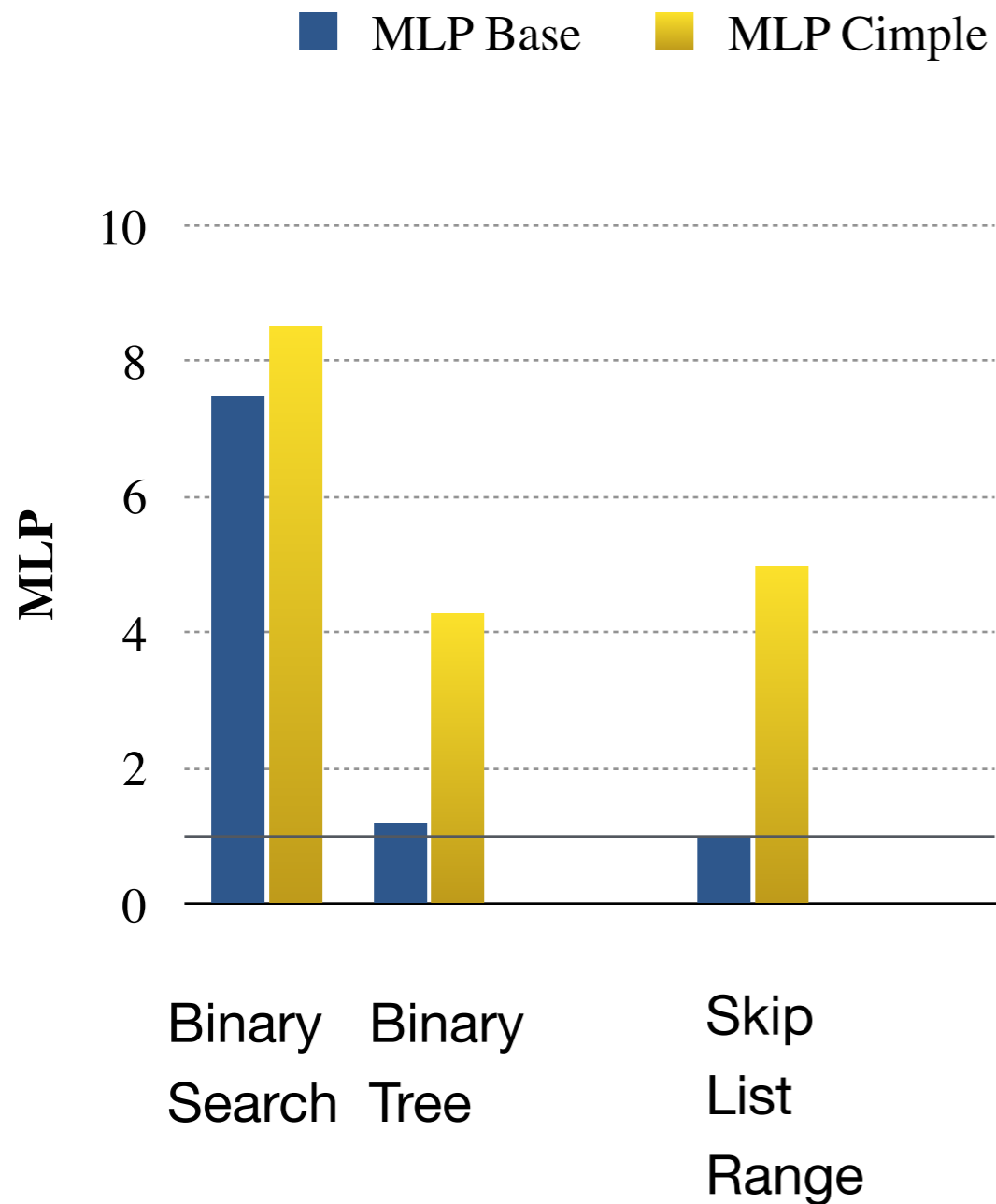
- HashTable -
OoO HW extracts MLP
- SkipList - speedup
matching MLP gains



MLP: Ineffective or Low



MLP: Ineffective or Low

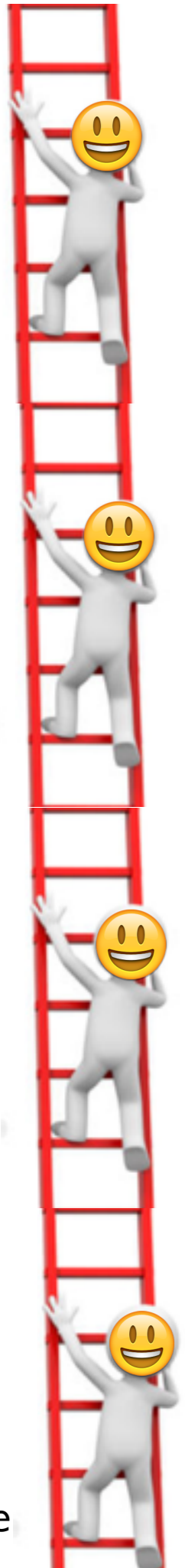


MLP Improvement Paths

- Increased Efficiency
 - Static scheduling
 - Vectorization
- Increased Effectiveness
 - Dynamic scheduling: no bubbles
 - Branchless code

Conclusion

- Fast
 - up to 6.4× speedup
- Portable DSL (*for Stephanies*)
 - template libraries
 - database query engines
- Next: C++ standards (*for Joes*)



Thanks

Vladimir Kiriansky
vlk@csail.mit.edu

