

Building Efficient Query Engines in a High-Level Language

Yannis Klonatos, Christoph Koch, Tiark Rompf, Hassan
Chafi

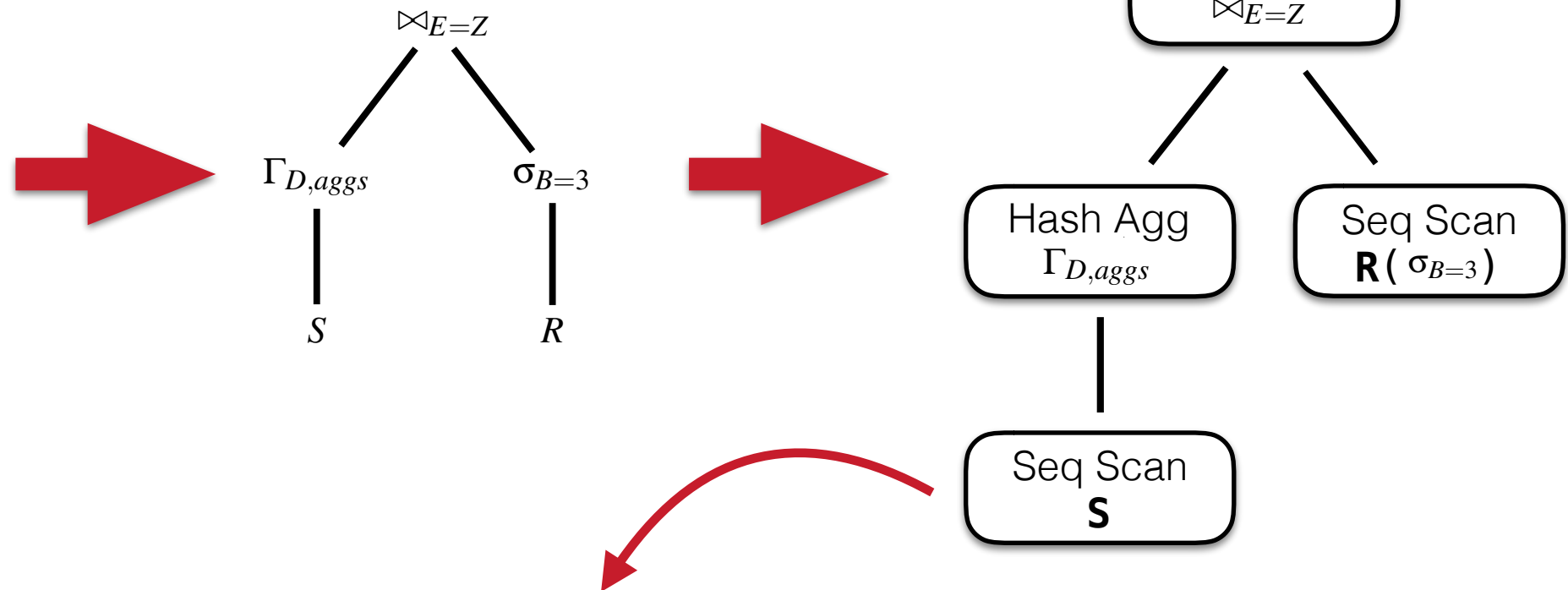
Prashanth Menon, Reading Group Fall 2015



Background

- What happens to your SQL query?

```
1 select *
2 from R, (select S.D,
3          sum(1-S.B) as E,
4          sum(S.A*(1-S.B)),
5          sum(S.A*(1-S.B)*(1+S.C))
6 from S group by S.D) T
7 where R.Z=T.E and R.B=3
```



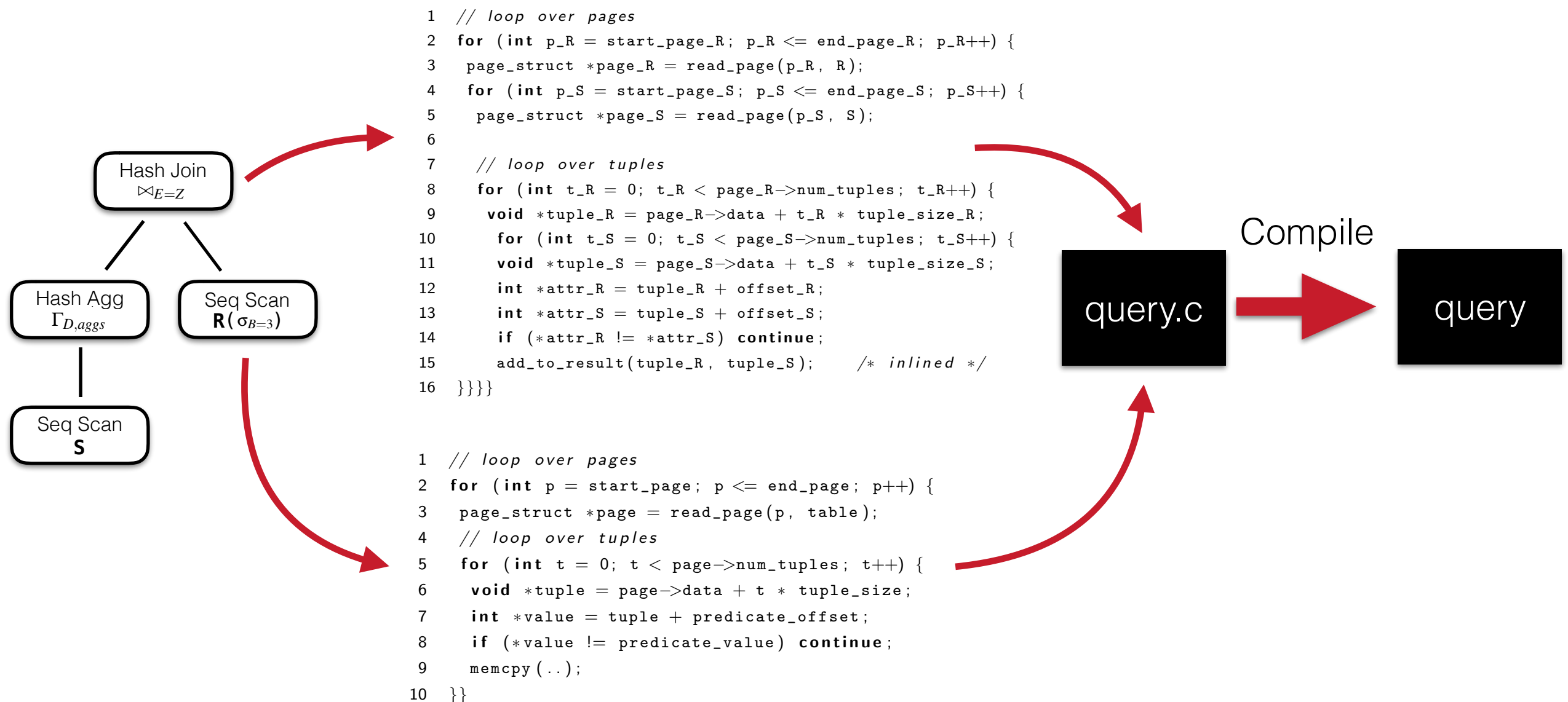
```
1 class Operator {
2     virtual void open() = 0;
3     virtual Tuple* next() = 0;
4     virtual void close() = 0;
5 };
```

Background (2)

- Volcano model is powerful, generic and composable
- Designed in an era where disk I/O dominated overhead
- If all data stored in main memory, it doesn't perform well
 - All **next()** calls are virtual (i.e., vtable lookup)
 - Single function call overhead for each tuple, for each operator!
 - Pretty poor cache utilization
- Can we do better?

Background (3)

- Generate a per-query execution engine!



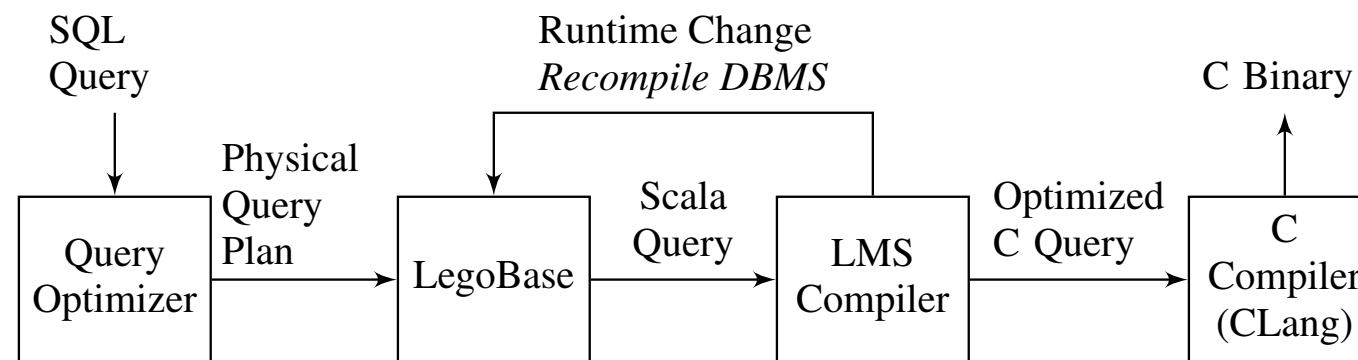
Background

- Compiling queries will yield better performance
- However, template expansion is:
 - Brittle
 - Very low level (i.e., hard to implement)
 - Limited scope of compilation
 - Limited adaptivity

Goal

- Performance of low-level hand-written query code
- Productivity of high-level language with rich type system guarantees

LegoBase



- Query engine written in Scala
- Cross-compiles Scala query plans into optimized C code
- Four steps:
 1. Convert pre-assembled physical query plan to naive Scala-based operator tree
 2. Use Lightweight Modular Staging (LMS) to convert operator tree into Scala IR
 3. Execute multiple optimization passes on IR
 4. Output optimized Scala or C query plan
- Optimizations are written in Scala, operate on Scala types
- **Programmatic removal of abstraction overhead**

Optimizations

- Optimizations are performed in LMS passes
 - Similar to LLVM where passes are independent
- Optimizations include:
 - Inter-operator optimizations
 - Eliminating redundant materializations
 - Data structure specialization
 - Data layout changes
 - Traditional compiler optimizations (DCE, loop unrolling)

Inter-Operator Optimizations

- Convert query plan from pull-based to push-based (à la HyPer)
 - Operators **push** data to consumer operators
 - Better cache locality (no function calls, tuples remain in registers)

```
1 case class HashJoin[B](leftChild: Operator,  
2   rightChild: Operator, hash: Record=>B,  
3   cond: (Record,Record)=>Boolean) extends Operator {  
4   val hm = HashMap[B,ArrayBuffer[Record]]()  
5   var it: Iterator[Record] = null  
6   def next() : Record = {  
7     var t: Record = null  
8     if (it == null || !it.hasNext) {  
9       t = rightChild.findFirst { e =>  
10         hm.get(hash(e)) match {  
11           case Some(hl) => it = hl.iterator; true  
12           case None => it = null; false  
13         }  
14       }  
15     }  
16     if (it == null || !it.hasNext) return null  
17     else return it.collectFirst {  
18       case e if cond(e,t) => conc(e, t)  
19     } get  
20   }  
21 }
```



```
1 case class HashJoin[B](leftChild: Operator,  
2   rightChild: Operator, hash: Record=>B,  
3   cond: (Record,Record)=>Boolean) extends Operator {  
4   val hm = HashMap[B,ArrayBuffer[Record]]()  
5   var it: Iterator[Record] = null  
6   def next(t: Record) {  
7     var res: Record = null  
8     while (res = {  
9       if (it == null || !it.hasNext) {  
10         hm.get(hash(t)) match {  
11           case Some(hl) => it = hl.iterator  
12           case None => it = null  
13         }  
14       }  
15       if (it == null || !it.hasNext) null  
16       else it.collectFirst {  
17         case e if cond(e,t) => conc(e, t)  
18       } get  
19     } != null) parent.next(res)  
20   }  
21 }
```

Inter-Operator Optimizations (2)

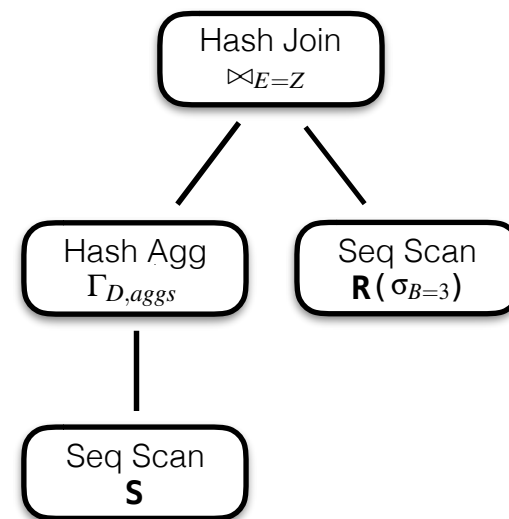
- Convert query plan from pull-based to push-based (à la HyPer)
 - Operators **push** data to consumer operators
 - Better cache locality (no function calls, tuples remain in registers)

```
1 case class HashJoin[B](leftChild: Operator,  
2   rightChild: Operator, hash: Record=>B,  
3   cond: (Record,Record)=>Boolean) extends Operator {  
4   val hm = HashMap[B,ArrayBuffer[Record]]()  
5   var it: Iterator[Record] = null  
6   def next(t: Record) {  
7     if (it == null || !it.hasNext) {  
8       hm.get(hash(t)) match {  
9         case Some(hl) => it = hl.iterator  
10        case None => it = null  
11      }  
12    }  
13    while (it!=null && it.hasNext) it.collectFirst {  
14      case e if cond(e,t) => parent.next(conc(e,t))  
15    }  
16  }  
17 }
```



```
1 case class HashJoin[B](leftChild: Operator,  
2   rightChild: Operator, hash: Record=>B,  
3   cond: (Record,Record)=>Boolean) extends Operator {  
4   val hm = HashMap[B,ArrayBuffer[Record]]()  
5   def next(t: Record) {  
6     hm.get(hash(t)) match {  
7       case Some(hl) => hl.foreach { e =>  
8         if (cond(e,t)) parent.next(conc(e,t))  
9       }  
10      case None => {}  
11    }  
12  }  
13 }
```

Redundant Materialization



- Not necessary to materialize aggregations
- Can bypass aggregation node and perform aggregation in build phase of join
- Difficult (probably not impossible) to express when using code templates

Redundant Materialization

- Implemented as IR pass
- If we see an HJ node whose left child is Agg grouping on same join attribute, merge them

```
1 def optimize(op: Operator): Operator = op match {
2   case hj@HashJoin(aggOp:AggOp,_,h,eq) =>
3     new HashJoin(aggOp.child,hj.rightChild,h,eq) {
4       override def open() {
5         // leftChild is now the child of aggOp
6         leftChild foreach { t =>
7           val key = hj.leftHash(aggOp.grp(t))
8           // Get aggregations from hash map of HJ
9           val aggs = hm.getOrElseUpdate(key,
10              new Array[Double](aggOp.aggFuncs.size))
11           aggOp.processAggs(aggs,t)
12         }
13       }
14     }
15   case x: Operator =>
16     x.leftChild = optimize(x.leftChild)
17     x.rightChild = optimize(x.rightChild)
18   case null => null
19 }
```

Data-Structure Specialization

- Use schema and query knowledge to specialize hash maps
 - Remove abstraction overhead of generic hash maps
- Three main problems:
 1. Redundant data storage (key is usually subset of value)
 2. Lookups require virtual calls to hashing functions
 3. Hash maps require resizing during runtime
- LegoBase solutions:
 1. Convert hash map to contiguous array (of buckets)
 2. Only store values in nodes
 3. Inline hash and equality functions
 4. Use runtime statistics to predict and allocate size of hash map at compile time

Changing Data Layout

- Possible to switch between row and column form at runtime
 - Does not require rewriting query engine
- Implemented as an IR optimization pass
 - Triggered when we see Array[Record] (array of type record) in IR
 - Possible to implement any new data storage layout as an IR optimization

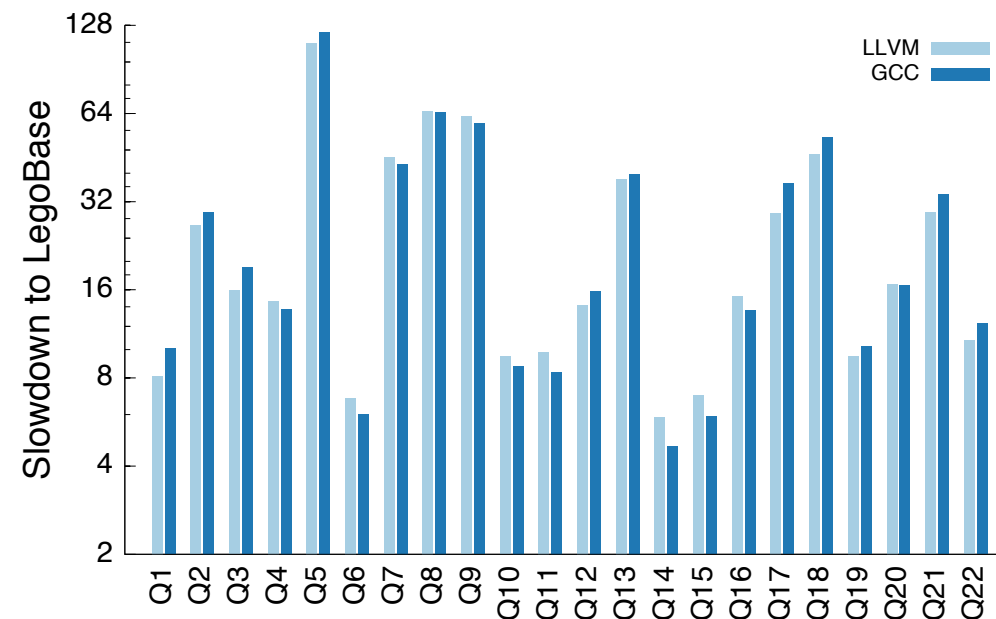
```
3  override def array_new[T:Manifest](n:Int) =  
4    manifest[T] match {  
5      case Record(attrs) =>  
6        // Create a new array for each attribute  
7        val arrays = for (tp<-attrs) yield array_new(n)(tp)  
8        // Pack everything in a new record  
9        record(attrs, arrays)  
10     case _ => super.array_new(n)  
11  }
```

```
24  override def array_apply[T:Manifest](ar:Array[T],  
25                                     n:Int) =  
26    manifest[T] match {  
27      case Record(attrs) =>  
28        val arrays = for (l <- attrs) yield field(ar, l)  
29        val elems = for (a <- arrays) yield a(n)  
30        // Perform record reconstruction  
31        record(attrs, elems)  
32      case _ => super.array_apply(ar, n)  
33    }
```

Evaluation Setup

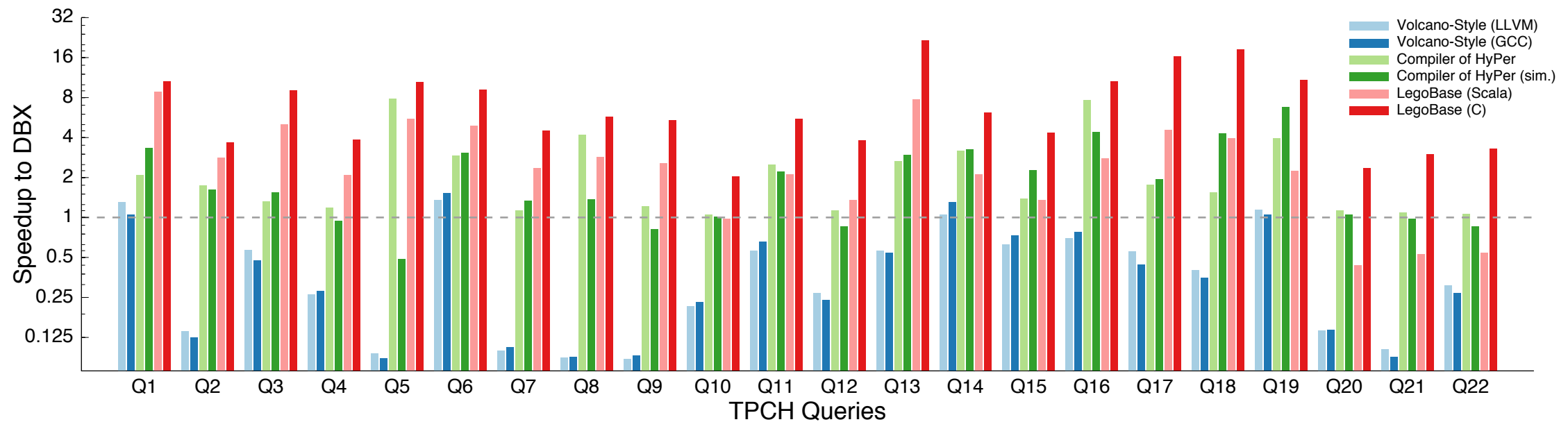
- 2x Intel Xeon 2GHz, 256GB RAM, 2TB HDD
- Scala 2.10.3, Clang 2.9
- Evaluate against DBX (in-memory row-store) and HyPer
- All systems get 192 GB RAM
- Run TPCH

Optimizing Query Plans



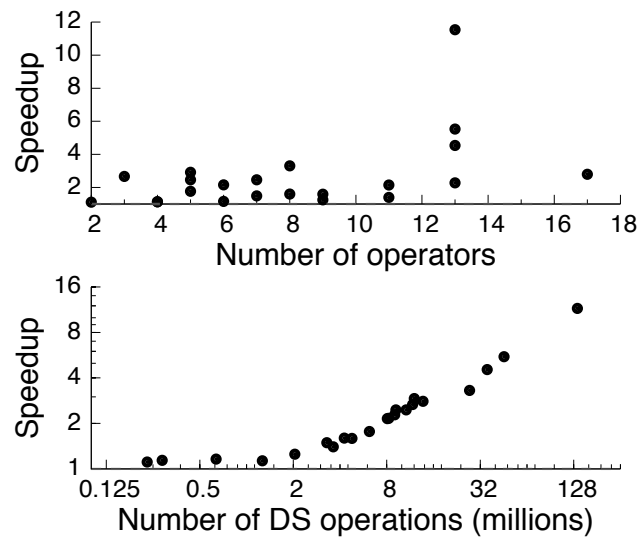
- LegoBase Volcano-style query engine compiled to C
- Compare code compiled with GCC and LLVM against fully optimized LegoBase
- Not all that interesting ... really a comparison of GCC and LLVM

TPCH Query Optimization

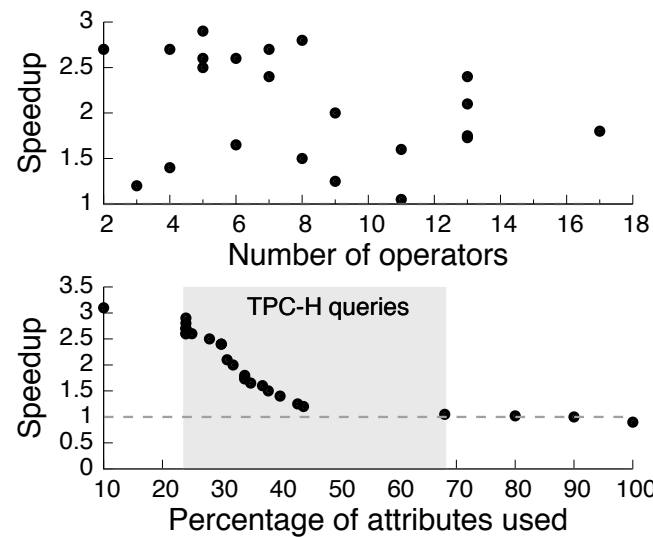


- Simulated HyPer is faster than HyPer
 - Due to data-structure specialization
- LegoBase is 5.3x-7.7x faster than HyPer
 - Due to data-structure specialization, data layout optimization
 - Better cache locality, branch prediction, fewer instructions executed
- LegoBase Scala 2.5x slower than LegoBase C
 - 1.3x-1.4x more branch mispredictions
 - 1.1x-1.8x more LLC misses
 - 5.5x more CPU instructions

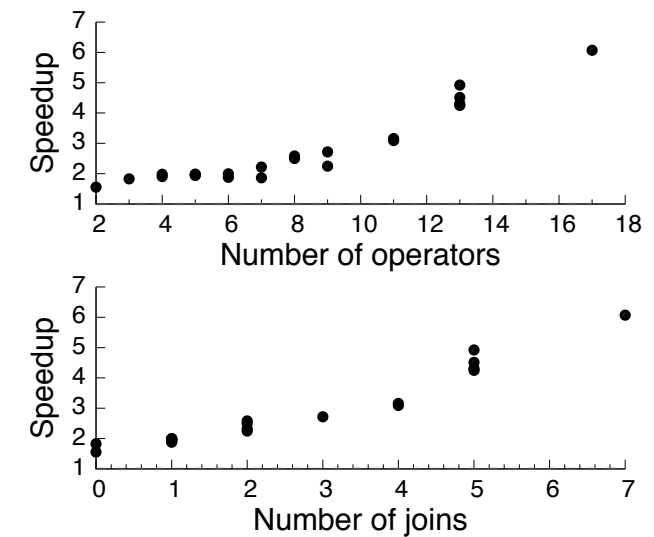
Impact of Compiler Optimizations



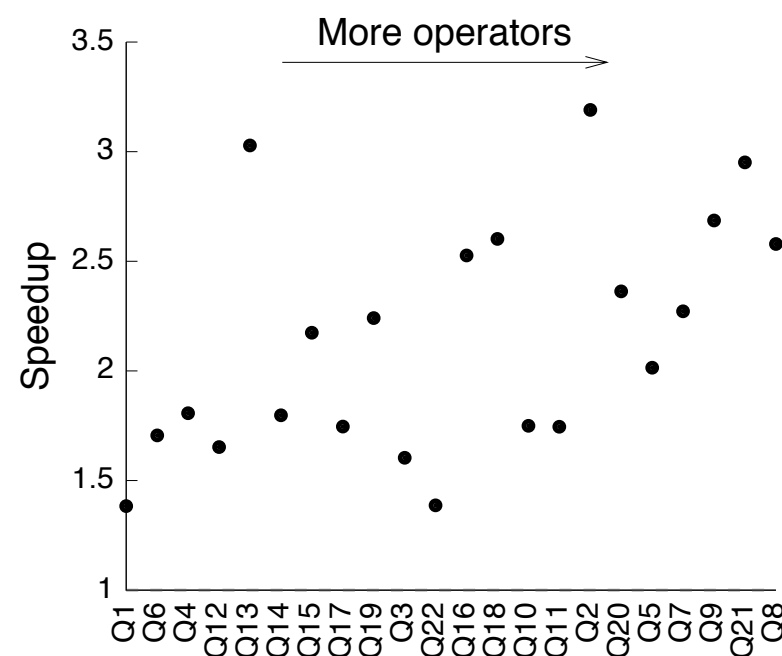
(a) Data Structure Opt.



(b) Change Data Layout



(c) Operator Inlining



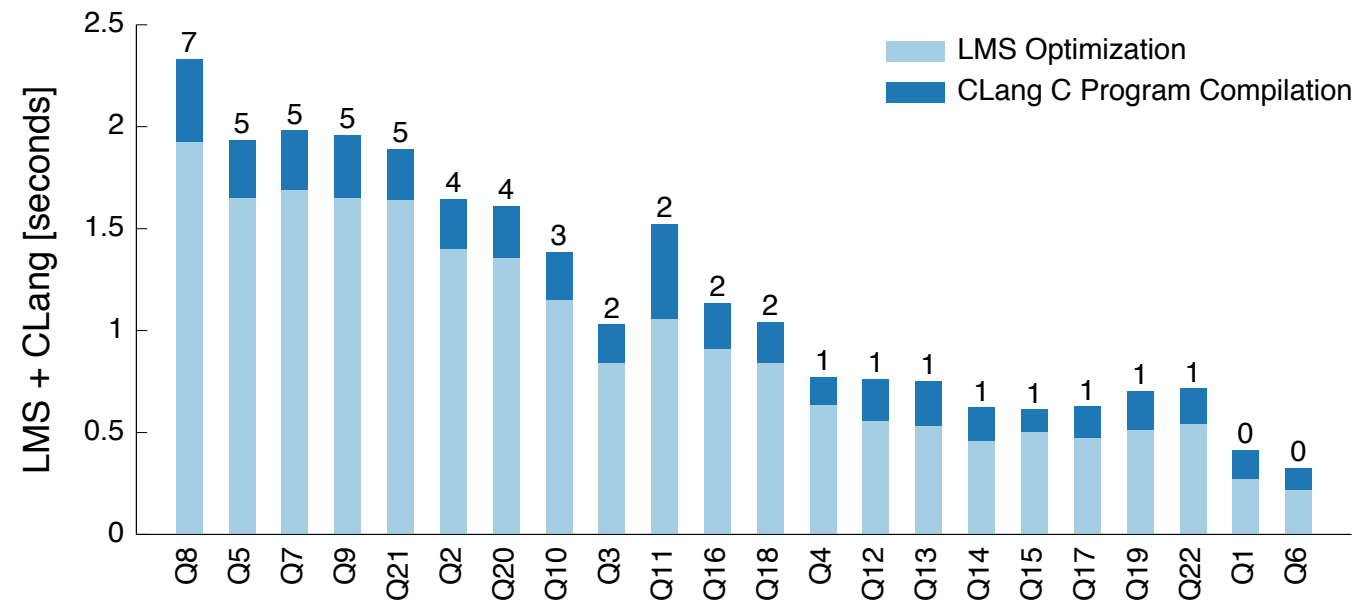
Productivity

	Coding Effort	Scala LOC	Average Speedup
Operator Inlining	–	0	2.07×
Push Engine Opt.	1 Week	~400 ^[6]	2.26×
Data Structure Opt.	4 Days	259	2.16×
Change Data Layout	3 Days	102	1.81×
Other Misc. Opt.	3 Days	124	– ¹⁰
LegoBase Operators	1 Month	428	–
LMS Modifications	2 Months	3953	–
Various Utilities	1 Week	538	–
Total	~4 Months	5831	7.7×

Table 1: Programming effort required for each LegoBase component along with the average speedup obtained from using it.

- Optimizations all done in a high-level language
- Easier to program, fewer lines of code
- High speedup-per-line-of-code

Compilation Overhead



- Compilation time ~ 2.5 seconds

Conclusions

- Possible to build query engine in high-level language with performance of hand-written low-level C
- Use LMS to transform naive query engine to IR
 - Optimize IR in independent stages
 - Specialize types, change data layouts at runtime
 - Emit optimize C code
- Performance beats existing main-memory DBMS and modern query compiler HyPer