

Paper Assigned: Runtime Code Generation in Cloudera Impala

Presentation by: Shivam Patel





What is Cloudera Impala?

- Impala is a MPP(Massively Parallel Processing) database for use in the Hadoop ecosystem.
- Uses HDFS and HBase(Hadoop version of Big Table) as storage managers.
- Uses Map Reduce for processing.
- Claims to execute queries in 10-100x speed compared to out of the box solutions.



How is the Speed up Achieved?

- (In part) By using run time information and JIT compilation to create optimized machine code that attempts to reduce instruction level latencies.
- Uses LLVM to compile and load "fully-optimized query specific" versions of functions at runtime for use in queries.



Instruction level latencies?

Query Specific?

- For a particular query that uses a function, the function will typically act based on specific information that can only be known at runtime, hence when the function is written, it needs to handle general cases.
- When we write code to handle general cases, it involve excessive instructions to check things(types, number of elements, increment things)
- Example: Branch misprediction, Stalling over previous instruction reading/writing not being complete.



How does Impala(want to) solve this?

- When a query is about to be executed, invoke the JIT compiler LLVM offers to create specialized machine code using the runtime information and then use that instead of interpreting.
- Helps to:
 - Remove Conditionals
 - Remove Loads
 - Inline Virtual Functions

Q: What is gained?



Optimization 1: removing conditionals, Can YOU write code without If statements or <, <=, >, >=, ?

- Using runtime information, create machine code that does things like loop unrolling.
- If loops are unrolled, then there doesn't exist the need to increment a counter or check if its less than a bound.
- Typically, looping over data where the data has indeterminate type requires checking its type before doing an operation.
- Modern CPUs implement a feature called "Branch Prediction" which assume a conditional is true, hence they "conditional jump" to the segment of machine code that it would imply.
- Only a few clock cycles later would it realize it made a mistake and reverse, hence clogging up the throughput.

Loop unrolling example

```
shivam@debian: ~/Desktop/CS595
```

File Edit View Search Terminal Help

```
#include <stdio.h>
#include <stdlib.h>
```

```
int main(int argc, char **argv){
    int n = atoi(argv[1]);
    int c = 0;
    for(int i = 0; i < n; i++){
        c++;
    }

    return c;
}
```

12, 1

ALL

shivam@debian: ~/...

File Edit View Search Terminal Help

```
(base) shivam@debian:~/Desktop/CS5
95$ scrot asm_code.jpeg
(base) shivam@debian:~/Desktop/CS5
95$ scrot normal_code.jpeg
```

[illegible]

11,7

Bot

Loop unrolling continued

shivam@debian: ~/Desktop/CS595
File Edit View Search Terminal Help

loop_unroll_n.o: file format elf64-x86-64

Disassembly of section .text:

```
0000000000000000 <main>:
0: 55          push    %rbp
1: 48 89 e5    mov     %rsp,%rbp
4: 48 83 ec 20 sub     $0x20,%rsp
8: 89 7d ec    mov     %edi,-0x14(%rbp)
b: 48 89 75 e0 mov     %rsi,-0x20(%rbp)
f: 48 8b 45 e0 mov     -0x20(%rbp),%rax
13: 48 83 c0 08 add     $0x8,%rax
17: 48 8b 00    mov     (%rax),%rax
1a: 48 89 c7    mov     %rax,%rdi
1d: e8 00 00 00 00 callq   22 <main+0x22>
22: 89 45 f4    mov     %eax,-0xc(%rbp)
25: c7 45 fc 00 00 00 00 movl    $0x0,-0x4(%rbp)
2c: c7 45 f8 00 00 00 00 movl    $0x0,-0x8(%rbp)
33: eb 08      jmp     3d <main+0x3d>
35: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
39: 83 45 f8 01 addl    $0x1,-0x8(%rbp)
3d: 8b 45 f8    mov     -0x8(%rbp),%eax
40: 3b 45 f4    cmp     -0xc(%rbp),%eax
43: 7c f0      jl      35 <main+0x35>
45: 8b 45 fc    mov     -0x4(%rbp),%eax
48: c9        leaveq %eax
49: c3        retq

21,49-53 All
```

shivam@debian: ~/...
File Edit View Search Terminal Help

```
(base) shivam@debian:~/Desktop/CS5
95$ scrot normal_code.jpeg
(base) shivam@debian:~/Desktop/CS5
95$ scrot asm_code.jpeg
```

shivam@debian: ~/Desktop/CS595
ch Terminal Tabs Help

shivam@debian: ... shivam@debian: ...

file format elf64-x86-64

Disassembly of section .text:

```
0000000000000000 <main>:
0: 55          push    %rbp
1: 48 89 e5    mov     %rsp,%rbp
4: 89 7d ec    mov     %edi,-0x14(%rbp)
7: 48 89 75 e0 mov     %rsi,-0x20(%rbp)
b: c7 45 fc 00 00 00 00 movl    $0x0,-0x4(%rbp)
12: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
16: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
1a: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
1e: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
22: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
26: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
2a: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
2e: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
32: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
36: 83 45 fc 01 addl    $0x1,-0x4(%rbp)
3a: 8b 45 fc    mov     -0x4(%rbp),%eax
3d: 5d        pop     %rbp
3e: c3        retq

12,51-55 All
```




Optimization 2: removing memory loads, when can you get around memory access penalty?

- Typically, when memory is written into a variable, and an instruction immediately follows that wants to read it, the processor has to stall until the writeback stage of the pipeline occurs. Similar for write to memory then read from memory.
- Based on runtime information, if we know the structure data, then we can add compiler optimizations to remove excessive memory operations, as well remove conditionals that do checks.
- Ex: suppose we know our data has 9 types, and our tables have at most 4, we can use the same function but remove the checks as well as excessive reads to repeated intermediate variables.

```
void MaterializeTuple(char* tuple) {  
    for (int i = 0; i < num_slots_; ++i)  
    {  
        char* slot = tuple + offsets_[i];  
        switch(types_[i]) {  
            case BOOLEAN:  
                *slot = ParseBoolean();  
                break;  
            case INT:  
                *slot = ParseInt();  
                break;  
            case FLOAT: ...  
            case STRING: ...  
            // etc.  
        }  
    }  
}
```

interpreted

```
void MaterializeTuple(char* tuple) {  
    *(tuple + 0) = ParseInt();    // i = 0  
    *(tuple + 4) = ParseBoolean();// i = 1  
    *(tuple + 5) = ParseInt();    // i = 2  
}
```

codegen'd

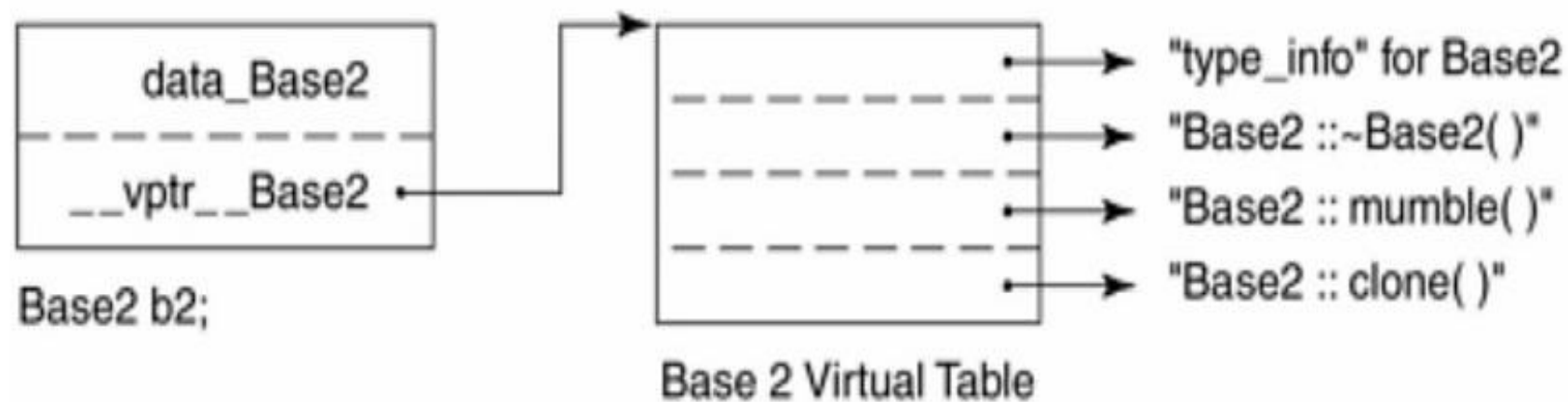
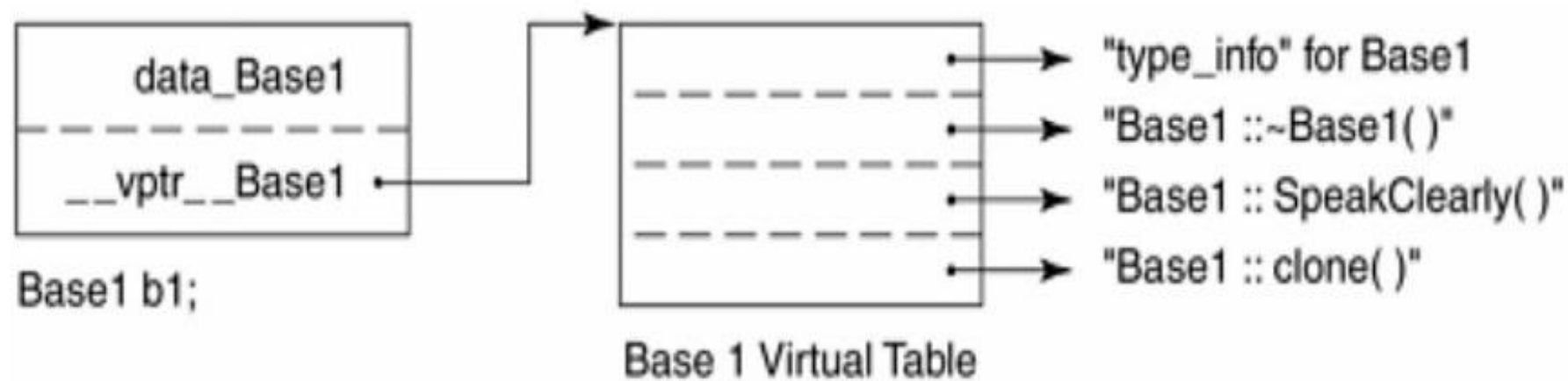
Basic five-stage pipeline

Instr. No.	Clock cycle	1	2	3	4	5	6	7
1		IF	ID	EX	MEM	WB		
2			IF	ID	EX	MEM	WB	
3				IF	ID	EX	MEM	WB
4					IF	ID	EX	MEM
5						IF	ID	EX



Optimization 3: Inline virtual functions

- In C++ `virtual` (abstract in java) is a keyword for class methods that allow for a child class to implement it and use at a later point in time.
- The issue is, suppose the method depends on a template (generic) that isn't determined until runtime, then without code generation the compiler will have to maintain a vtable (virtual function table) in the object or globally to index at runtime to jump to the "real" function address.
- So to prevent excessive instructions (indirection and conditional jumps), a JIT compiled version that has runtime information to determine "which version" can be used to inline a particular instantiation of a virtual function.





Putting it all together (kind of)

- LLVM has what's called IR (intermediate representation) which is a pseudo assembly code with its own virtual instruction set.
- LLVM IR can be generated via either IRBUILDER api to do so programmatically (more difficult, error prone, potentially less optimized) or via compiler front ends that compile to IR for later use in JIT compilation.
- For C++ clang allows compilation of C++ code to LLVM IR (does gcc?)
- Impala does the second, i.e. compiling to LLVM via clang then injecting JIT compiled into process at runtime with runtime optimizations made that classically would only be possible at compile time.



Idea sounds nice in theory but...

- The paper claims their implementation has so far managed to replace interpreted (virtual) function calls with inlined specific ones and that they haven't implemented the other 2 optimizations.
- Arguably optimization #3 is the weakest of the 3, since in principle its preventing only a few extra instructions, compared to what would be saved by removing inner loops, conditional jumps, and memory read/write stalls.
- They claim they're working on optimization 1 and 2 and the ability to do this to pre compiled functions(they didnt explain how, perhaps functions from a shared library, so code can be borrowed and modified with runtime info and then dynamically linked? .so/.dll/.dylib depending on linux/windows/mac)



Support for user defined functions

- Impala provides support for users to provide UDFs
- rather than writing a function, compiling it to a shared object file, and then dynamically linking at runtime,
- They suggest it be compiled to IR via clang, and then JIT compiled with runtime info to allow UDFs to have same performance as native functions. Perhaps a detail they left out about the "precompiled" functions from the previous slide, is they be compiled to IR on the system when Impala is installed, and then are susceptible to the inlining and future optimizations.

Performance

	# Instructions	# Branches
Code generation disabled	72,898,837,871	14,452,783,201
Code generation enabled	19,372,467,372	3,318,983,319
Speedup	4.29x	3.76x

Table 4: Instruction and branch counts for TPC-H-Q1

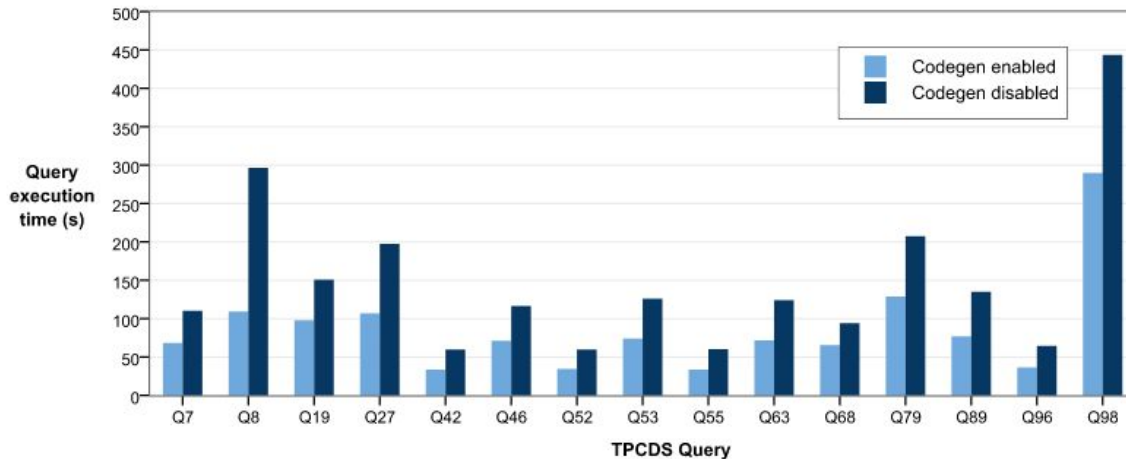


Figure 4: TPC-DS query execution with and without code generation



Interesting final thought.

- Clearly this paper is old and not fully implemented.
- In a fully implemented version of Impala with loops, conditional jumps, and pipeline stalling memory reads/writes optimized away, thing how its achieved.
- In principle its by adding runtime constants to individual instructions, and writing more instructions to get rid of loops and conditional jumps
- This results in a large number of instructions being added. Think in terms of opcodes, register ids and runtime constants. At minimum $\log_2(\#Opcodes) + 2$ for the registers plus possibly upto 8 bytes each for a instruction constant at best? So code segment and data segment of process probably increase massively.



The End.

Questions?