Transpiling PL/pgSQL UDF in DuckDB

Adrian Abedon, Yuchen Liu, Arvin Wu

Goal: We want to generate C++ code for PL/pgSQL UDFs



Not supported by DuckDB Supported by DuckDB



Support DuckDB vectorization

 Inline DuckDB operators during compilation



Project Roadmap



Phase I

Transpile PL/pgSQL into a scalar UDF

- Essentially the control flow in the UDF
- Pass SQL queries not parsed to be handled by DuckDB

Phase II

Transpile PL/pgSQL into a vectorized UDF

• Still calling DuckDB to handle queries

Phase III

Compile SQL queries through operator inlining

- Gradually support the compilation of certain queries
 - Comparators >, <
 - Numeric Operators +, -, *
 - DuckDB built-in functions

Phase I: Transpile PL/pgSQL into a scalar UDF

• Focus was on transpiling control structures

- Function body and arguments
- If statements

Ο

• For/While loops

return a+10:

• Let DuckDB handle queries

return con.Query("select \$1+10", a);

• Issues with DuckDB scalar UDF API

- Only supports 3 arguments
- Loose precision with DECIMAL type
- Difficult to interpret certain types such as DATE since they are passed as int32

Phase II: Transpile PL/pgSQL into a vectorized UDF Motivation



- All the drawbacks of scalar UDF
- The bottleneck of performance: **HUGE** overhead of calling DuckDB's query
 - $\sim 1 \text{ ms} \approx 10^6 \text{ CPU}$ cycles for each simple query: select l_quantity+1

By switching to the vectorized UDF framework, we can share that overhead in the vector

- Most DuckDB operators are vectorized
 - void Equals(Vector &left, Vector &right, Vector &result, idx_t count);
 - Important when operators are inlined in Phase III

Phase II: Transpile PL/pgSQL into a vectorized UDF

• Rewrite control structures to support vectorization

- Function body and arguments
- If statements
- For/While loops
- Use vectorized prepared statements

• Explore how to write vectorized UDF

- Interact directly with duckdb vectors
 - Flat vector
 - Constant vector
 - Dictionary vector

Phase II: Transpile PL/pgSQL into a vectorized UDF

- Have to rewrite control structures to support vectorization
 - Function body and arguments
 - If statements

```
    For/While loops
    2. template<typename TR, typename... Args>
    Us
    Void CreateVectorizedFunction(string name, vector<SQLType> args, SQLType ret_type, scalar_function_t udf_func, SQLType varargs = SQLType::INVALID)
```

Phase II: Transpile PL/pgSQL into a vectorized UDF

- Have to rewrite control structures to support vectorization
 - Function body and arguments
 - If statements



Vectorized control structures

- Vectorizing arbitrary control structures is difficult
- Our approach was to keep track of "active lanes" (std::vector<bool>)
 - Set upon if condition, return, break, or continue
- If a lane is inactive, no computation happens on that lane

Example:

```
a = [5,10,4,13]
b = [1,11,3,14]
if(a < b): -> active lanes = [0,1,0,1]
a = a+b -> a = [5,21,4,27]
```

```
create function line count(oprio char, mode varchar)
returns int as $$
declare val int = 0;
begin
   if(mode = 'high') then
        if(oprio = '1-URGENT' OR oprio = '2-HIGH') then
            val = 1:
        end if;
   elsif(mode = 'low') then
        if(oprio <> '1-URGENT' AND oprio <> '2-HIGH') then
            val = 1;
        end if;
    end if;
    return val:
end $$
LANGUAGE PLPGSQL;
```

```
std::vector<duckdb::Value> tempvar116 = query26(shipdate, &active2, &returns1, NULL, NULL);
std::vector<bool> active3 = active2;
for (size_t tempvar115 = 0; tempvar115 < tempvar116.size(); tempvar115++)
{
    // if we are in a while loop, we don't need to check loop_active/continues since
    // the tempvar116 will be appropriately set to null.
    if (!tempvar116[tempvar115].IsNull())
    {
        active3[tempvar115] = active3[tempvar115] && (tempvar116[tempvar115].template GetValue<bool>());
     }
     }
        /** RETURN **/
```

```
Exa
```

```
for (size_t tempvar113 = 0; tempvar113 < args.size(); tempvar113++)
{
    if (active3[tempvar113] && !returns1[tempvar113])
    {
        // note that this lane has returned
        returns1[tempvar113] = true;
        // set its return value
        return_values[tempvar113] = tempvar114[tempvar113];
    }
    for (size_t tempvar115 = 0; tempvar115 < active3.size(); tempvar115++)
    {
        if (!tempvar116[tempvar115].IsNull())
        /
    }
}</pre>
```

std::vector<duckdb::Value> tempvar114 = const vector gen(0);

active3[tempvar115] = !active3[tempvar115];

Phase III: Compile queries

- Not yet implemented
- Support only queries that do not access a table
 - i.e. select (x*2.0) < 20;
- Inline DuckDB operators and functions
 - /**, *, <, +, =,...**
 - $\circ \quad date_add, date_part, cast$
- Current vectorized transpiler provides a great foundation to add compilation

Correctness Test

- Run TPC-H queries with and without UDF's and compared results
 - Queries 1, 3, 4, 5, 6, 7, 9, 10, 12, 14, 19 from <u>Froid Preprint</u>

Query 3 with no UDF's

select l orderkey. sum(l extendedprice * (1 - l discount)) as revenue, o orderdate, o shippriority from customer. orders, lineitem where c mktsegment = 'BUILDING' and c custkey = o custkey and 1 orderkey = o orderkey and o orderdate < date '1995-03-15' and 1 shipdate > date '1995-03-15' group by 1 orderkey, o orderdate, o_shippriority order by revenue desc. o orderdate limit 10:

Query 3 with UDF's



Performance without Query Compilation

Current version is significantly bottlenecked by making calls • TLDR: Its slow right now but it can be pretty fast

Performance without Query Compilation

TPC-H Performance Benchmarks





Query Type









Query Type

Conclusions from Benchmarks

- Calling DuckDB to execute queries in UDF's is way to slow
- Need to compile DuckDB operators
- Compiling UDFs makes sense if it is performing general computation, such as **inside the SELECT**
- Compiling UDFs does NOT make sense if it is used to **filter rows**
 - DuckDB (and other databases) will not be able to apply
 - block skipping
 - short circuiting predicate evaluations with AND
 - Better to use Froid/APFEL approach
- We can potentially improve performance with UDFs

Code Quality

- The framework we have developed for transpilation is extensible
 - Can choose to compile certain operators
 - For unsupported operators we can call DuckDB
- All C++ building blocks come from yaml templates
 - Can easily modify C++ output
- We feel that code is past the prototype stage but far from production ready
 - Clean up code
 - Generated C++ code is not formatted

Future Work

- Compile DuckDB operators
- Support more PL/pgSQL
 - nested blocks
 - block labels
- Support UDF calls within UDF's
 - Limitation of current approach of calling DuckDB for queries
 - \circ This will come with compilation

