

# Диалект MLIR для Общего Тайлинга Циклов

А.В. Левченко

2@lto.ru

Суперкомпьютерный центр СПбПУ

# Диалект MLIR Common Loop Tiling (CLT)

## Progressive lowering

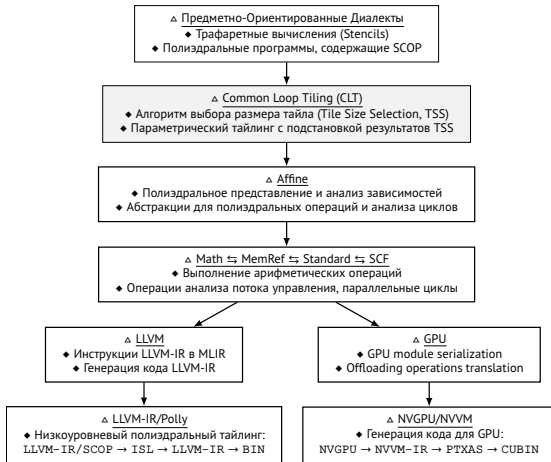


Figure 1 – Схема взаимодействия CLT с некоторыми диалектами инфраструктуры MLIR

# Диалект MLIR Common Loop Tiling (CLT)

Полиэдральный тайлинг с использованием конверсии диалектов Affine → GPU

```
1 // Циклы CLT/Affine
2 func.func @clt_loop(%arg0: memref<?x?x?x?xf32>,
3 % arg1: memref<?x?x?x?xf32> ) {
4   clt.launch [dev = "gpu"] {
5     clt.for %arg2 = 0 to 512 {
6       clt.for %arg3 = 0 to 256 {
7         clt.for %arg4 = 0 to 128 {
8           clt.for %arg5 = 0 to 64 {
9             % 0 = memref.load %arg0[%arg2, %arg3, %arg4,
10 % arg5] : memref<?x?x?x?xf32>
11             memref.store %0, %arg1[%arg2, %arg3, %arg4, %arg5]
12             : memref<?x?x?x?xf32>
13           [...]
14           clt.terminator
15           // Тайлинг ( размер тайла = 256, размер кеша = 1024 )
16           func.func @loop_tiling(
17 % arg0: memref<196608x1xf32>,
18 % arg1: memref<196608x1xf32>)
19 -> memref<196608x1xf32> {
20   clt.launch [dev = "gpu", tss = 256, cache = 1024] {
21     clt.for %arg2 = 0 to 196608 step 1024 {
22       clt.for %arg3 = 0 to 1 step 256 {
23         clt.for %arg4 = #map(%arg2) to #map1(%arg2) step 256 {
24           clt.for %arg5 = #map(%arg3) to #map2(%arg3) step 256 {
25             clt.if #set(%arg4) {
26               clt.for %arg6 = #map(%arg4) to #map2(%arg4) {
27                 clt.for %arg7 = #map(%arg5) to #map3(%arg5) {
28                   % 0 = clt.load %arg0[%arg6, %arg7]
29                   : memref<196608x1xf32>
30                   clt.store %0, %arg1[%arg6, %arg7]
31                   : memref<196608x1xf32>
32                 [...]
33               }
34             }
35           }
36         }
37       }
38     }
39   }
40 }
41 }
```

```
1 // Циклы GPU/Affine
2 module {
3   func.func @clt_loop(%arg0: memref<?x?x?x?xf32>,
4 % arg1: memref<?x?x?x?xf32> ) {
5     % c0 = arith.constant 0 : index
6     % c512 = arith.constant 512 : index
7     % 0 = arith.subi %c512, %c0 : index
8     % c1 = arith.constant 1 : index
9     % c0_0 = arith.constant 0 : index
10    % c256 = arith.constant 256 : index
11    % 1 = arith.subi %c256, %c0_0 : index
12    % c1_1 = arith.constant 1 : index
13    % c1_2 = arith.constant 1 : index
14    // Операция запуска ядра GPU
15    gpu.launch blocks(%arg2, %arg3, %arg4)
16    in (%arg8 = %0, %arg9 = %c1_2, %arg10 = %c1_2)
17    threads(%arg5, %arg6, %arg7)
18    in (%arg11 = %1, %arg12 = %c1_2, %arg13 = %c1_2) {
19      % 2 = arith.addi %c0, %arg2 : index
20      % 3 = arith.addi %c0_0, %arg5 : index
21      affine.for %arg14 = 0 to 128 {
22        affine.for %arg15 = 0 to 64 {
23          % 4 = memref.load %arg0[%2, %3, %arg14, %arg15]
24          : memref<?x?x?x?xf32>
25          memref.store %4, %arg1[%2, %3, %arg14, %arg15]
26          : memref<?x?x?x?xf32>
27        }
28      }
29      gpu.terminator
30    }
31    return
32  }
33 }
```

Figure 2 — Циклы с параметрическим тайлингом в MLIR CLT/Affine (листинг 1) и пример конверсии в GPU (листинг 2)

# Низкоуровневый полиэдральный тайлинг

## LLVM/Polly + Integer Set Library

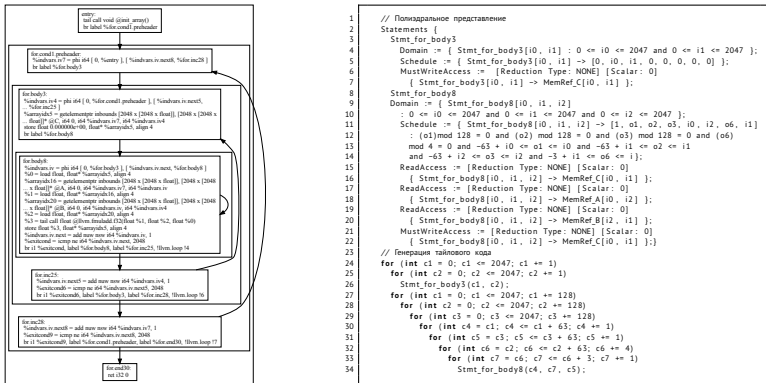


Figure 3 – Представление циклов в SCOP в LLVM-IR (листинг 1) и в полиэдральном описании с генерацией тайлового кода (листинг 2)

# Экспериментальные результаты трансляции SCOP

## Оценка эффективности операций

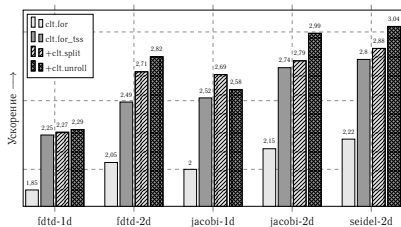


Figure 4 – CLT/Affine → GPU → NVGPU → NVVM-IR

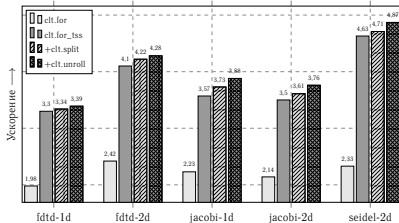


Figure 5 – CLT/Affine → LLVM → LLVM-IR → Polly + ISL

Figure 6 – Ускорение тайлового кода Stencils относительно версии NoOpt

- ◆ [🌐 ftp://ftp.lto.ru/clt.pdf](ftp://ftp.lto.ru/clt.pdf)
- ◆ Gysi, Müller, Zinenko, Herhut, Davis, Wicky, Fuhrer, Hoefler, Grosser. Domain-Specific Multi-Level IR Rewriting for GPU: The Open Earth Compiler for GPU-accelerated Climate Simulation (2021)
- ◆ Moses, Chelini, Zhao, Zinenko. Polygeist: Raising C to Polyhedral MLIR (2021)