# Sparse Matrix Compression Primitives With OpenCL Framework to Support Halide

Chao-Lin Lee[1], Chen-Ting Chao[1], Jenq-Kuen Lee[1],

Chung-Wen Huang[2], and Ming-Yu Hung[2]

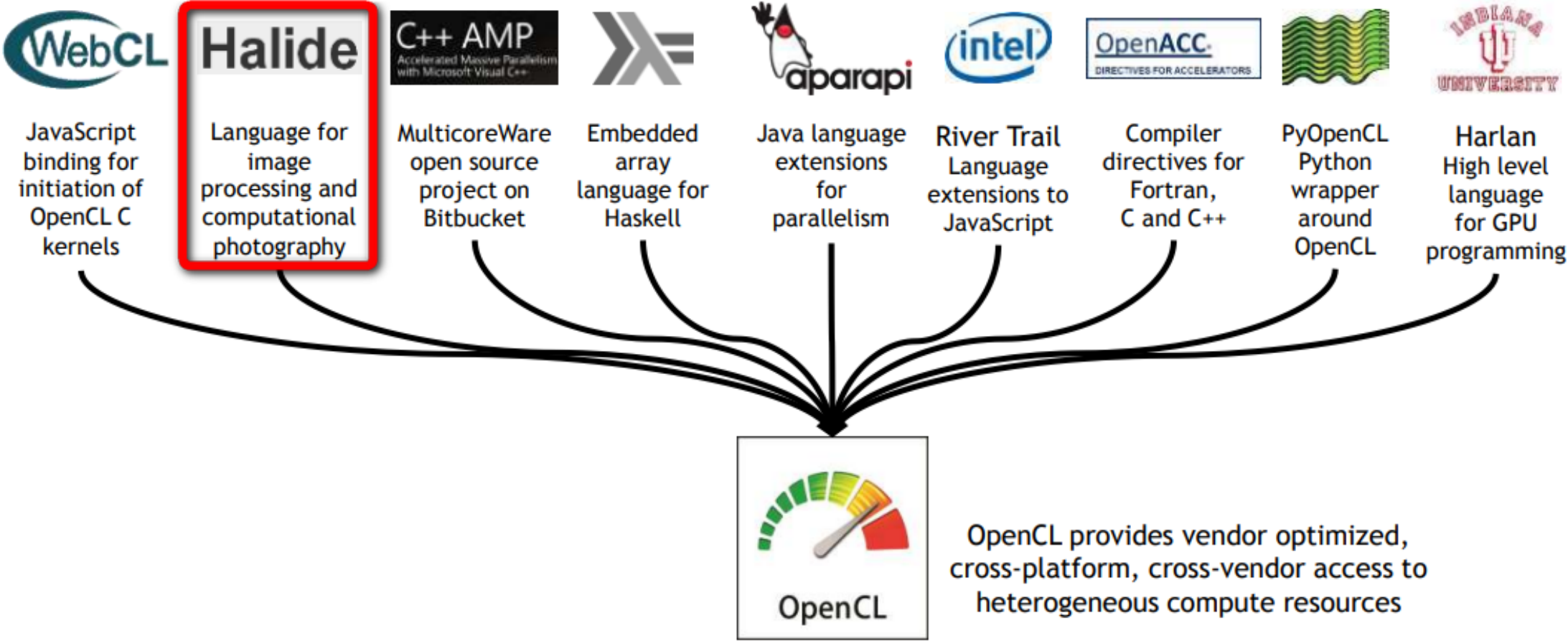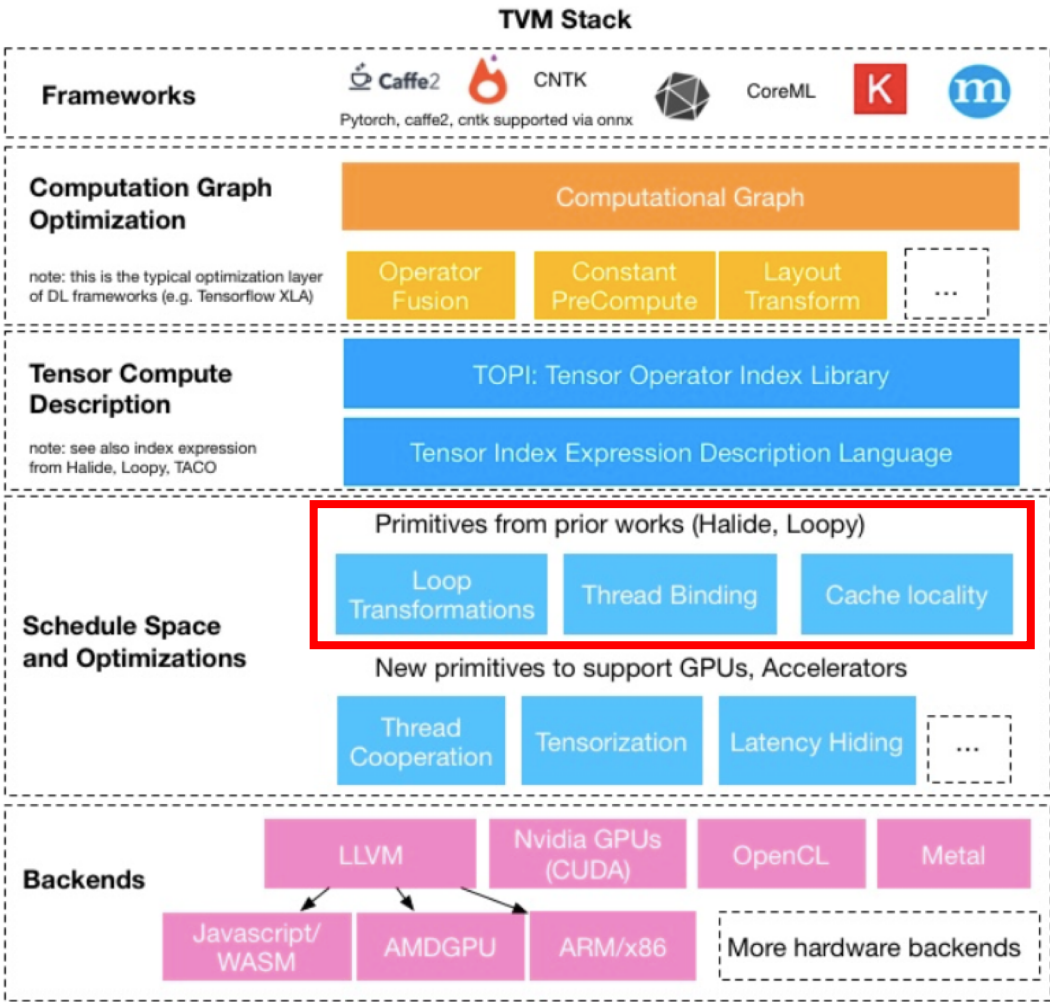National Tsing-Hua University,
Hsinchu, Taiwan[1]

MediaTek Inc.[2]

# Outline

- Background

- Halide overview

- Our Design in OpenCL with Halide

- Conclusion

# Background



WebCL — JavaScript binding for initiation of OpenCL C kernels

Halide — Language for image processing and computational photography

C++ AMP — Accelerated Massive Parallelism with Microsoft Visual C++ — MulticoreWare open source project on Bitbucket

Embedded array language for Haskell

aparapi — Java language extensions for parallelism

intel — River Trail Language extensions to JavaScript

OpenACC DIRECTIVES FOR ACCELERATORS — Compiler directives for Fortran, C and C++

PyOpenCL Python wrapper around OpenCL

INDIANA UNIVERSITY — Harlan High level language for GPU programming

OpenCL — OpenCL provides vendor optimized, cross-platform, cross-vendor access to heterogeneous compute resources

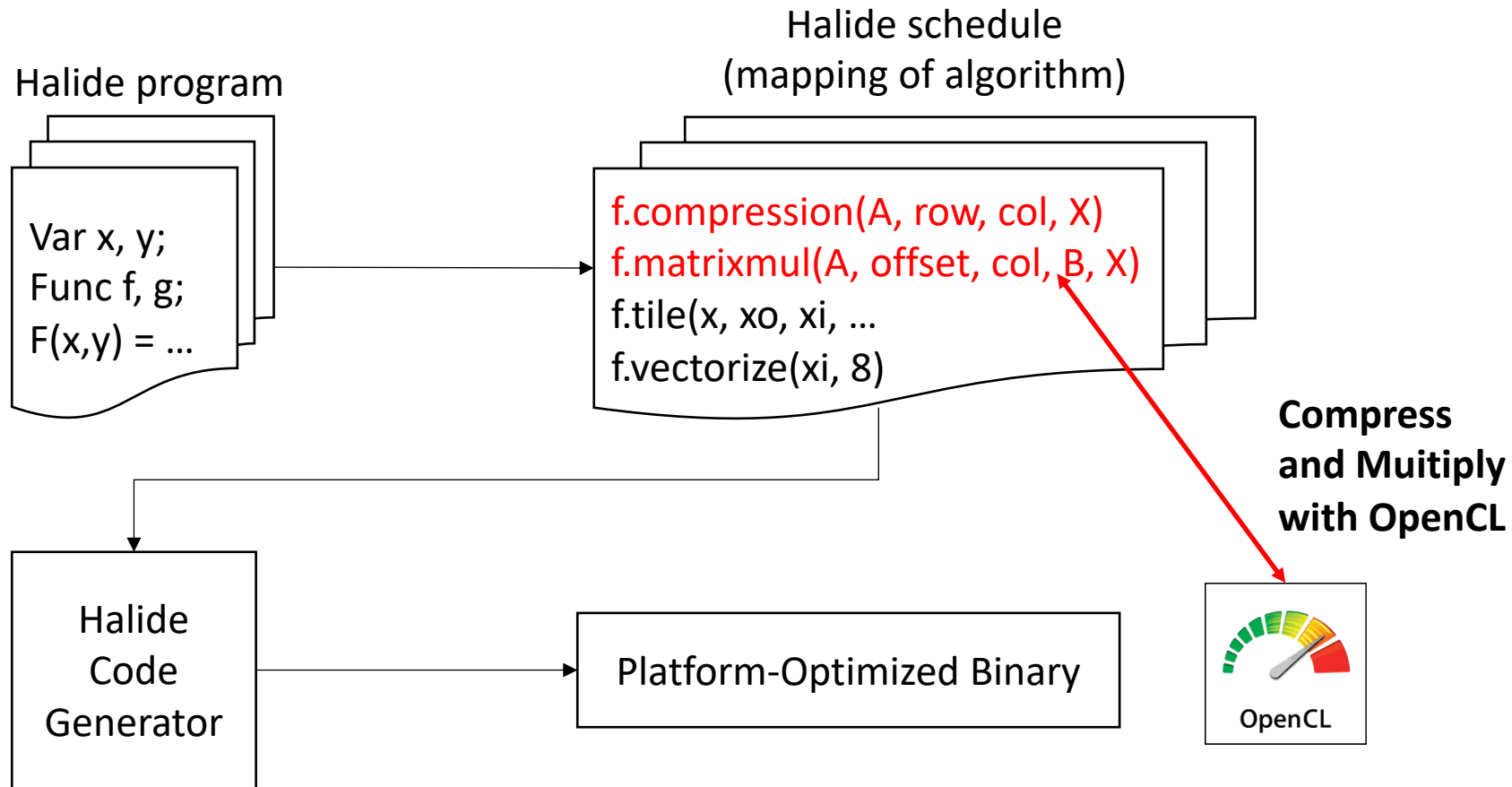Courtesy of Khronos, SIGGRAPH, Vancouver 2014

# Halide in TVM

# Halide and Halide IR

```
1   Func blur_3x3(Func input) {
2     Func blur_x, blur_y;
3     Var x, y, xi, yi;
4
5     // The algorithm - no storage or order
6     blur_x(x, y) = (input(x-1, y) + input(x, y) + input(x+1, y))/3;
7     blur_y(x, y) = (blur_x(x, y-1) + blur_x(x, y) + blur_x(x, y+1))/3;
8
9     // The schedule - defines order, locality; implies storage
10    blur_y.tile(x, y, xi, yi, 256, 32)
11          .vectorize(xi, 8).parallel(y);
12    blur_x.compute_at(blur_y, x).vectorize(x, 8);
13
14    return blur_y;
15  }
```

Reference: Halide – a language for fast portable portable computation on images and tensor

IWOCL    **The 7th International Workshop on OpenCL**

# The Flow for SPMM Enabled in Halide with OpenCL



Halide program

Var x, y;
Func f, g;
F(x,y) = …

Halide schedule
(mapping of algorithm)

f.compression(A, row, col, X)
f.matrixmul(A, offset, col, B, X)
f.tile(x, xo, xi, …
f.vectorize(xi, 8)

**Compress and Muitiply with OpenCL**

Halide Code Generator

Platform-Optimized Binary

OpenCL

# CSR

$$A = \begin{pmatrix} 7.5 & 2.9 & 2.8 & 2.7 & 0 & 0 \\ 6.8 & 5.7 & 3.8 & 0 & 0 & 0 \\ 2.4 & 6.2 & 3.2 & 0 & 0 & 0 \\ 9.7 & 0 & 0 & 2.3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 5.8 & 5.0 \\ 0 & 0 & 0 & 0 & 6.6 & 8.1 \end{pmatrix}$$

rowptr:  ( 0  4  7  10  12  14  16  )

colind:  ( 0  1  2  3  0  1  2  0  1  2  0  3  4  5  4  5  )

val:  ( 7.5  2.9  2.8  2.7  6.8  5.7  3.8  2.4  6.2  3.2  9.7  2.3  5.8  5.0  6.6  8.1  )

# Our Sparse GEMM in Halide

```
1   Halide::Buffer<int> row_ptr;
2   Halide::Buffer<int> col_idx;
3   Halide::Buffer<double> X;
4
5   Func compress;
6   compress(x, y) = 0;
7   compress.compression_csr(A, row_ptr, col_idx, X);
8   compress.realize(A.width(), A.height());
9
10  Halide::Buffer<double>C(B.width(), A.height());
11  Func mult;
12  mult(x,y) = cast<double>(0);
13  mult.computation_csr(A, row_ptr, col_idx, X, B, C);
14  mult.realize(B.width(), A.height());
```

# SPMM Parallelization in Halide with OpenCL Kernel

```
1    gpu_block<OpenCL> (f5.s1.r27$z.r30.__block_id_y, 0, 500) {
2      gpu_block<OpenCL> (f5.s1.r27$x.idxxo.__block_id_x, 0, 500) {
3        gpu_thread<OpenCL> (.__thread_id_y, 0, 8) {
4          gpu_thread<OpenCL> (.__thread_id_x, 0, 8) {
5            let f5.s1.r27$y.prologue.s = b20[((f5.s1.r27$x.idxxo.__block_id_x*8)
                  + .__thread_id_x)]
6            let f5.s1.r27$y.epilogue.s =
                  max(b20[((f5.s1.r27$x.idxxo.__block_id_x*8) + .__thread_id_x)],
                  b20[(((f5.s1.r27$x.idxxo.__block_id_x*8) + .__thread_id_x) +
                  1)])
7            let f5.s1.r27$y.new_min.s =
                  min(b20[((f5.s1.r27$x.idxxo.__block_id_x*8) + .__thread_id_x)],
                  max(min(f5.s1.r27$y.prologue.s, 28505), 0))
8            let f5.s1.r27$y.new_max.s =
                  max(min(b20[((((f5.s1.r27$x.idxxo.__block_id_x*8) +
                  .__thread_id_x) + 1)], max(min(f5.s1.r27$y.prologue.s, 28505),
                  0)), f5.s1.r27$y.new_min.s)
9            let t148 = (f5.s1.r27$z.r30.__block_id_y*8)
```
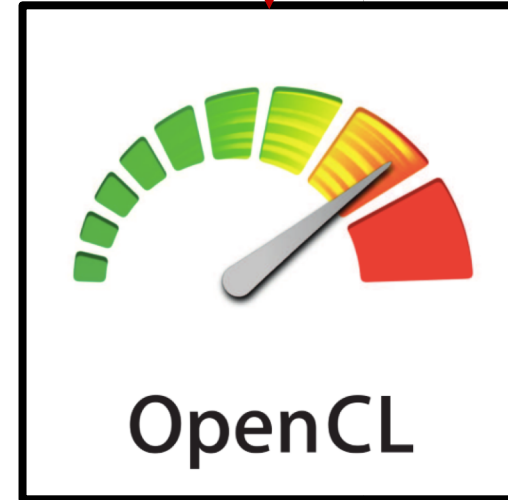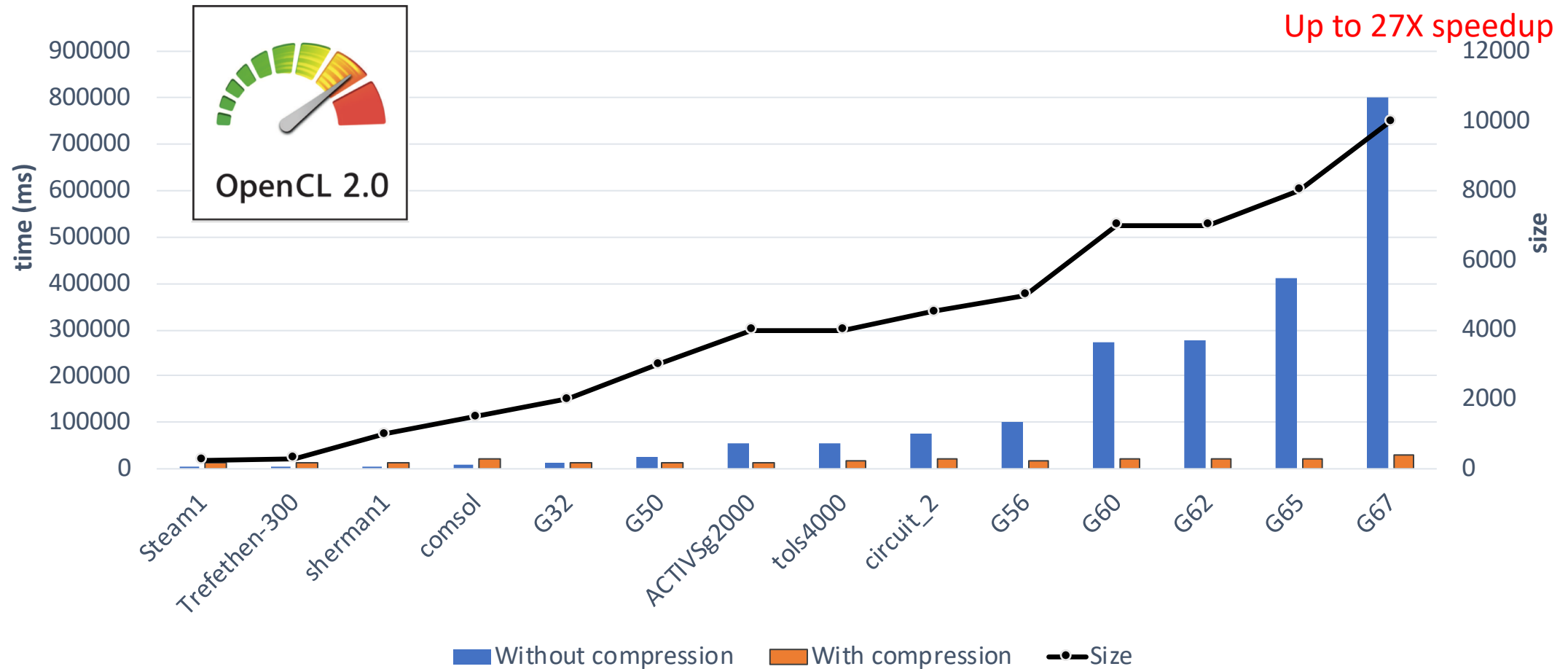
Halide

# SPMM Parallelization in Halide with OpenCL Kernel

```
1   __kernel void kernel_f5_s1_r27_z_r30___block_id_y(
2    const int _f5_stride_1,
3    const int _t143,
4    __address_space__b0 const double *restrict _b0,
5    __address_space__b22 const double *restrict _b22,
6    __address_space__f5 double *restrict _f5,
7    __address_space__b20 const int *restrict _b20,
8    __address_space__b21 const int *restrict _b21,
9    __address_space___shared int16* __shared)
10   {
11    int _f5_s1_r27__z_r30___block_id_y = get_group_id(1);
12    int _f5_s1_r27__x_idxxo___block_id_x = get_group_id(0);
13    int ___thread_id_y = get_local_id(1);
14    int ___thread_id_x = get_local_id(0);
15    int _55 = _f5_s1_r27__x_idxxo___block_id_x * 8;
16    int _56 = _55 + ___thread_id_x;
17    int _57 = _b20[_56];
18    int _58 = _56 + 1;
19    int _59 = _b20[_58];
20    int _60 = max(_57, _59);
21    int _61 = min(_57, 28505);
22    int _62 = max(_61, 0);
23    int _63 = min(_57, _62);
24    int _64 = min(_59, _62);
25    int _65 = max(_64, _63);
26    int _66 = _f5_s1_r27__z_r30___block_id_y * 8;
27    ...
```

# Sparse Matrix Compress and Multiplication with Halide / OpenCL

# Conclusion

- Matrix compression CSR in Halide with OpenCL

- Exploit the opportunities of optimizing GEMM convolution layers with sparse matrix compression

# Thank you