# CLBlast: A Tuned BLAS Library

Cedric Nugteren
May 16, 2018

http://github.com/cnugteren/clblast
http://cnugteren.github.io/clblast

IWOCL

# CLBlast?

BLAS: "Basic Linear Algebra Subpgrams"

- NVIDIA's cuBLAS is great, or is it?

- NVIDIA's cuBLAS is great, or is it?

- NVIDIA's cuBLAS is great, or is it?

  – Not portable, not customisable, not open-source, …

# But why a new BLAS Library?

- NVIDIA's cuBLAS is great, or is it?

  - Not portable, not customisable, not open-source, …

- Is AMD's clBLAS great?

  - Not performance portable, not well engineered, …

  - Discontinued, superseeded by incomplete ROCblas

# But why a new BLAS Library?

- NVIDIA's cuBLAS is great, or is it?
  - Not portable, not customisable, not open-source, ...

- Is AMD's clBLAS great?
  - Not performance portable, not well engineered, ...
  - Discontinued, superseeded by incomplete ROCblas

- Don't vendors ship their hardware with their own libraries?
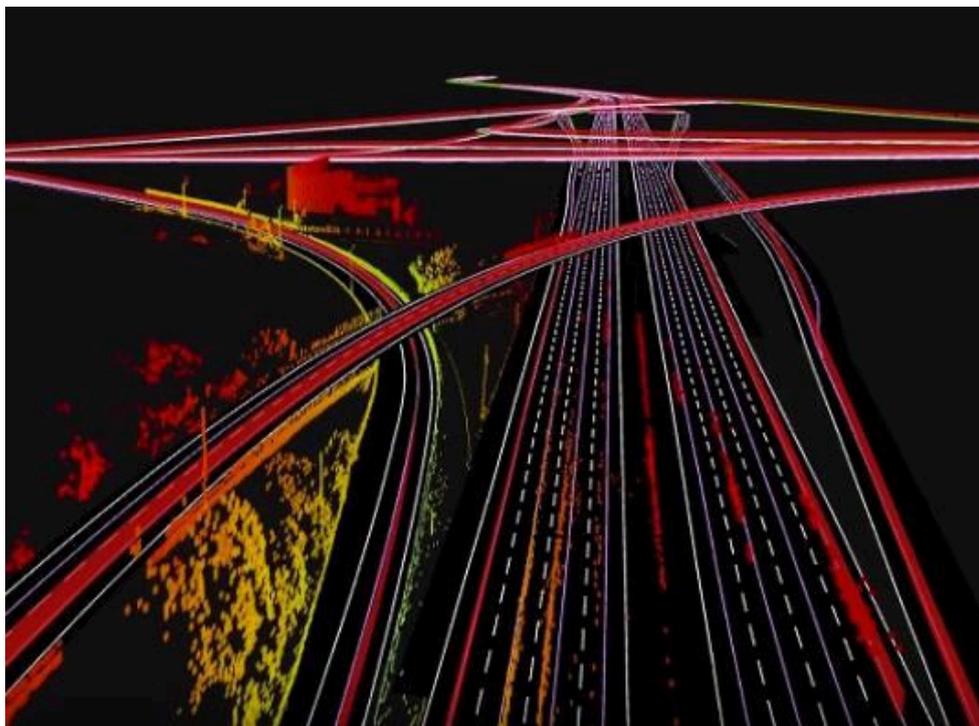  - Not a portable solution, and actually not always true...

clMathLibraries / **clBLAS**

<> Code    ⊙ Issues 63    Pull requests 0    Projects 0    Wiki

Filters ▾    🔍 is:issue is:open    Labels    Miles

⊙ 63 Open    ✓ 93 Closed    Author ▾

Not able to Build clBLAS from Sources
#315 opened 28 days ago by saviogeorge

⊙ make rebuilds every file even if no source files were changed
#310 opened on 22 Mar by cirosantilli

⊙ Leaked events in GEMM calls (and probably other functions)
#304 opened on 22 Feb by Oblomov

⊙ tests got segmentation fault
#299 opened on 25 Jan by arom4github

⊙ Error compiling during run-time
#293 opened on 5 Dec 2016 by aml5600

⊙ ixamax errors
#292 opened on 2 Nov 2016 by mgates3

⊙ test-short fails on Ubuntu with AMD Card

- HDMap making → Deep-learning

- Deep-learning → Fast BLAS libraries

- CLBlast: Modern C++11 OpenCL BLAS library

- Implements all BLAS routines for all precisions (S, D, C, Z)

- Accelerates all kinds of applications:

  - Fluid dynamics, quantum chemistry, linear algebra, finance, etc.

  - Some extra focus on deep learning

# Introducing CLBlast

- CLBlast: Modern C++11 OpenCL BLAS library

- Implements all BLAS routines for all precisions (S, D, C, Z)

- Accelerates all kinds of applications:

  - Fluid dynamics, quantum chemistry, linear algebra, finance, etc.

  - Some extra focus on deep learning

- Already integrated into various projects:

  - JOCLBlast (Java bindings)

  - PyCLBlast (Python bindings)

  - ArrayFire (GPU accelerated library and applications) {🔥}

  - OpenCL fork of Caffe (github.com/dividiti/ck-caffe)

  - OpenCL fork of TF (github.com/hughperkins/tensorflow-cl)

# Introducing CLBlast

- All kernels are generic and tunable thanks to integration of the CLTune auto-tuner (presented at MCSoC '15 and GTC '16)

# But… is it fast? 🔥

- All kernels are <span style="color:red">generic and tunable</span> thanks to integration of the CLTune auto-tuner (presented at MCSoC '15 and GTC '16)

```c
#define WGS 64  // The local work-group size
#define WPT 1   // The amount of work-per-thread
#define VW 1    // Vector width of vectors X and Y

typedef float dtype; // Example data-type
#if VW == 1
  typedef float dtypeV;
#elif VW == 2
  typedef float2 dtypeV;
#endif
```

```c
__kernel __attribute__(reqd_work_group_size(WGS))
void Xaxpy(const int n, const dtype alpha,
           const __global dtypeV* restrict xgm,
           __global dtypeV* ygm) {
  #pragma unroll
  for (int w=0; w<WPT; ++w) {
    int i = w*get_global_size(0)+get_global_id(0);
    ygm[i] = ygm[i] + alpha * xgm[i];
  }
}
```

- All kernels are generic and tunable thanks to integration of the CLTune auto-tuner (presented at MCSoC '15 and GTC '16)

```c
#define WGS 64   // The local work-group size
#define WPT 1    // The amount of work-per-thread
#define VW 1     // Vector width of vectors X and Y

typedef float dtype;  // Example data-type
#if VW == 1
  typedef float dtypeV;
#elif VW == 2
  typedef float2 dtypeV;
#endif
```

```c
__kernel __attribute__(reqd_work_group_size(WGS))
void Xaxpy(const int n, const dtype alpha,
           const __global dtypeV* restrict xgm,
           __global dtypeV* ygm) {
  #pragma unroll
  for (int w=0; w<WPT; ++w) {
    int i = w*get_global_size(0)+get_global_id(0);
    ygm[i] = ygm[i] + alpha * xgm[i];
  }
}
```

- Tuned out-of-the-box for 50+ common devices
  - For new devices: run the auto-tuner when installing CLBlast

# CLBlast Benchmark Results



AXPY
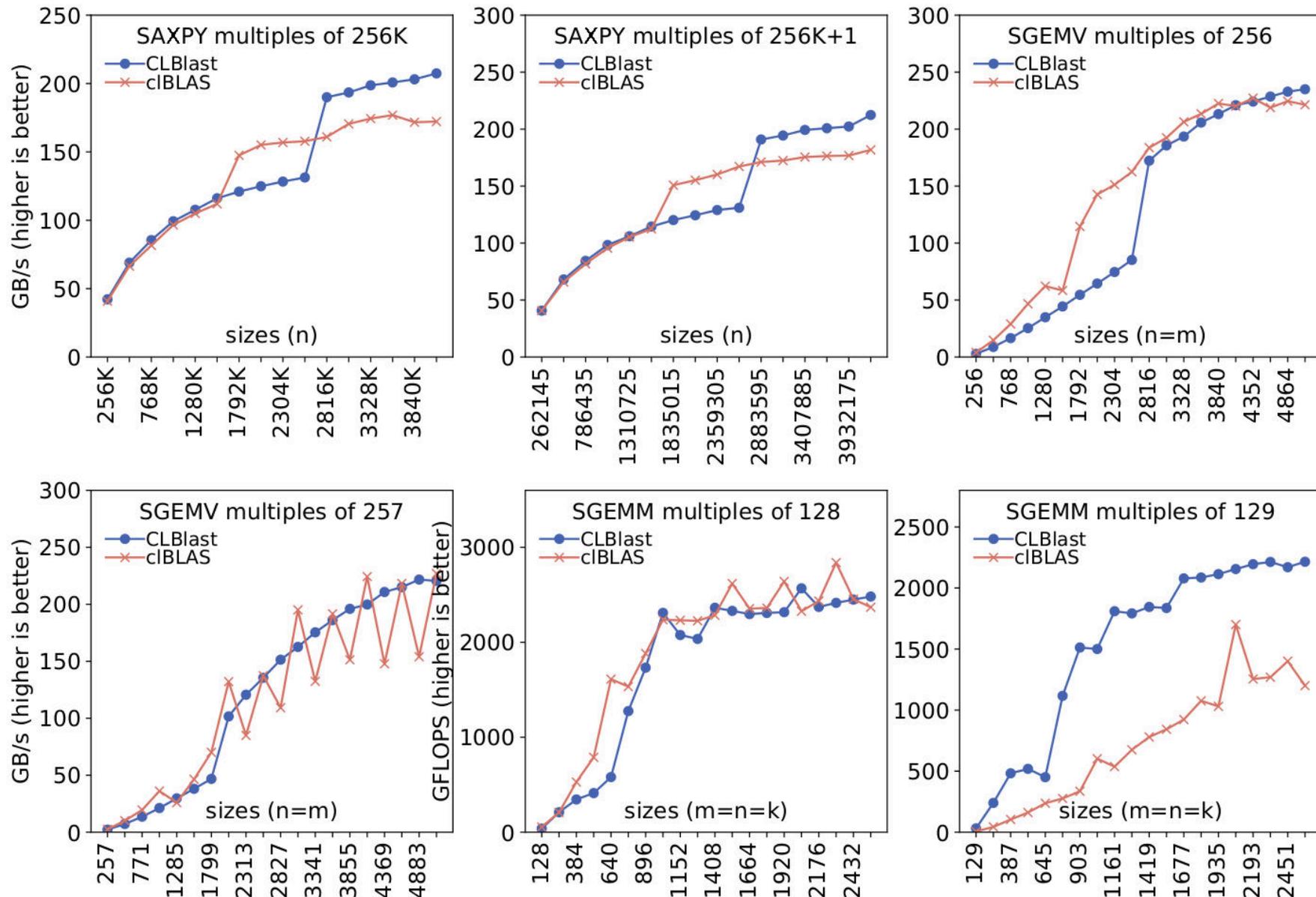regular
(in GB/s)

AXPY
odd
(in GB/s)

GEMV
regular
(in GB/s)

GEMV
odd
(in GB/s)

GEMM
regular
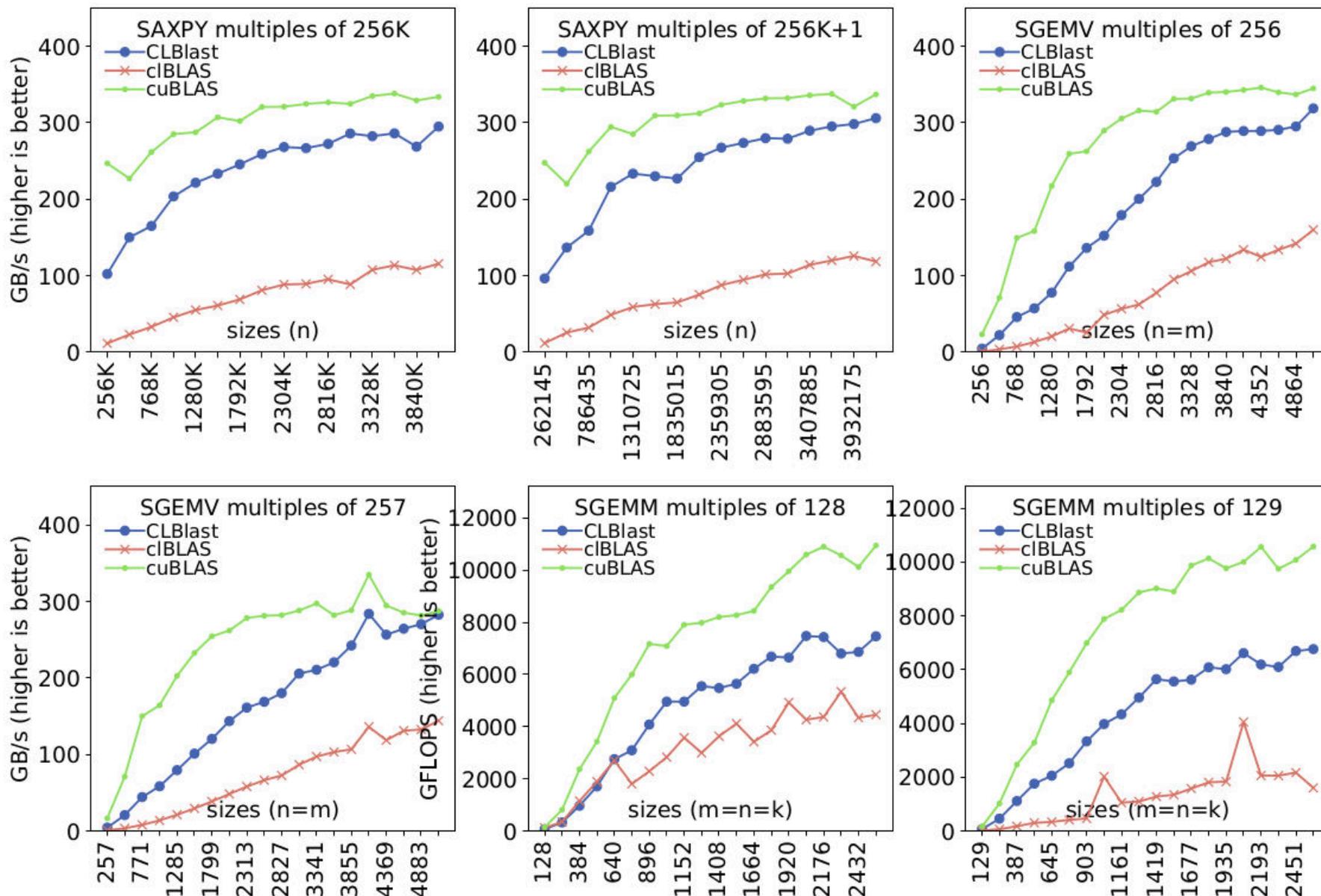(in GFLOPS)

GEMM
odd
(in GFLOPS)

- Higher is better

- More results at http://cnugteren.github.io/clblast
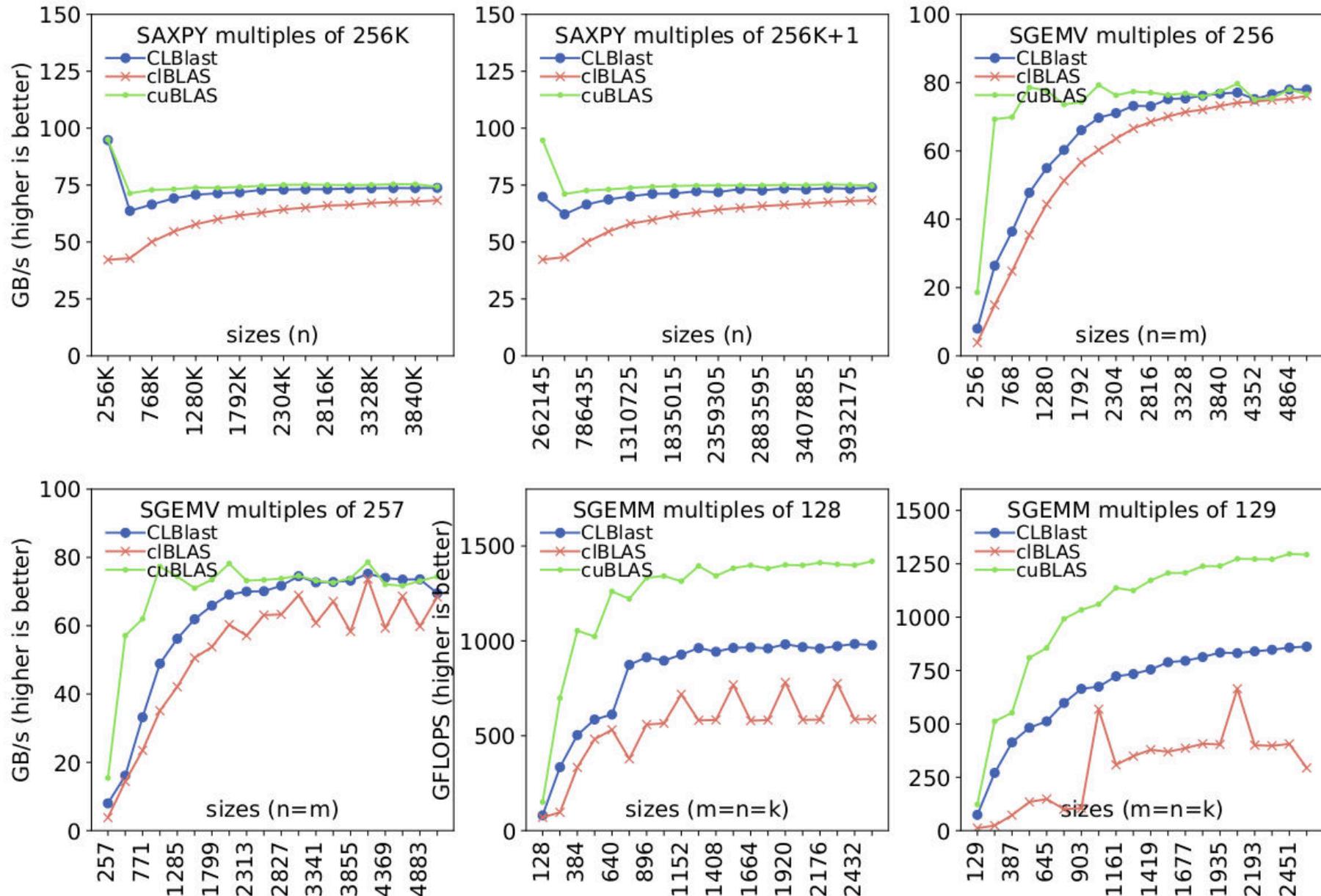
- On-par or better than clBLAS (especially for odd-sized GEMM)

- Much better than clBLAS, reasonably close to cuBLAS

- Much better than clBLAS, reasonably close to cuBLAS

- GEMM much better for CLBlast

- On-par or better than clBLAS

- On-par or better than clBLAS (especially for GEMM)

## Intel GEMM update April '18

| Library | GFLOPS |
|---|---|
| clBLAS | ~25 |
| CLBlast v1.3 | ~125 |
| Intel's GEMM snippet with subgroup shuffling | ~236 |

- On-par or bett...

## Intel GEMM update April '18

| Library | GFLOPS |
|---|---|
| clBLAS | ~25 |
| CLBlast v1.3 | ~125 |
| CLBlast v1.4 with 2D register tiling | ~180 |
| CLBlast v1.4 with Intel subgroup shuffling | ~230 |
| Intel's GEMM snippet with subgroup shuffling | ~236 |

- On-par or bett

- What can we do for the deep-learning community?

    - Problem-specific tuning

    - Half-precision floating-point (FP16)

    - Batched routines

- Default GEMM tuning:
  - 1024x1024 matrices

- Deep-learning:
  - <span style="color:red">Various but fixed matrix sizes</span> (dependent on network layout)
  - Typically smaller and/or rectangular matrices

- Default GEMM tuning:

  - 1024x1024 matrices

- Deep-learning:

  - <span style="color:red">Various but fixed matrix sizes</span> (dependent on network layout)

  - Typically smaller and/or rectangular matrices

- Potential for optimal performance in CLBlast:

  - <span style="color:red">Tuning for a custom size</span> possible

  - C++ API to change parameters at run-time

- SGEMM tuning for Radeon M370X GPU



Relative SGEMM performance on Radeon M370X

- SGEMM tuning for Radeon M370X GPU

- Best on the diagonal

- >100% due to random tuning



Relative SGEMM performance on Radeon M370X

- SGEMM tuning for Radeon M370X GPU

- Best on the diagonal

- >100% due to random tuning

- Gain of ~2x for some cases



Relative SGEMM performance on Radeon M370X

- Double-precision (FP64) not needed for deep-learning

- Even FP32 is too much → introducing half-precision FP16

- Implemented in low-power devices (ARM Mali, Intel GPUs) and deep-learning specific GPUs (Tesla P100, V100)

- Double-precision (FP64) not needed for deep-learning

- Even FP32 is too much → introducing half-precision FP16

- Implemented in low-power devices (ARM Mali, Intel GPUs) and deep-learning specific GPUs (Tesla P100, V100)

- Potential for 2x savings in:
  bandwidth, storage, compute, energy

# Half-precision floating-point (FP16)
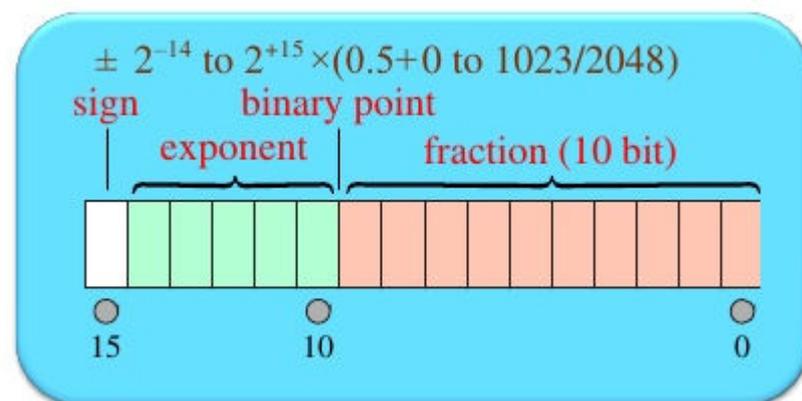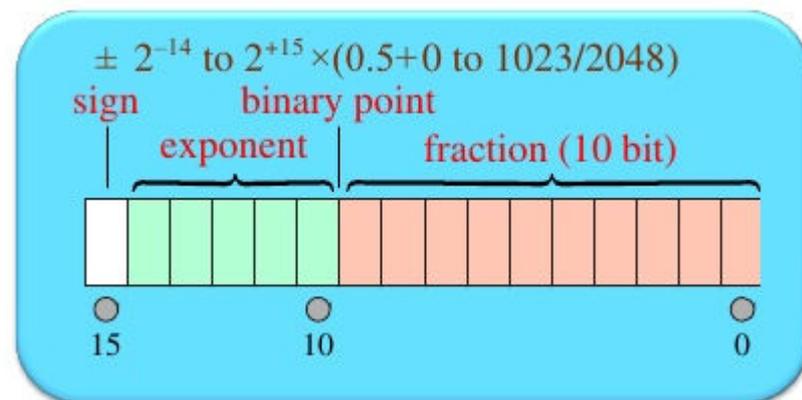
- Double-precision (FP64) not needed for deep-learning

- Even FP32 is too much → introducing half-precision FP16

- Implemented in low-power devices (ARM Mali, Intel GPUs) and deep-learning specific GPUs (Tesla P100, V100)

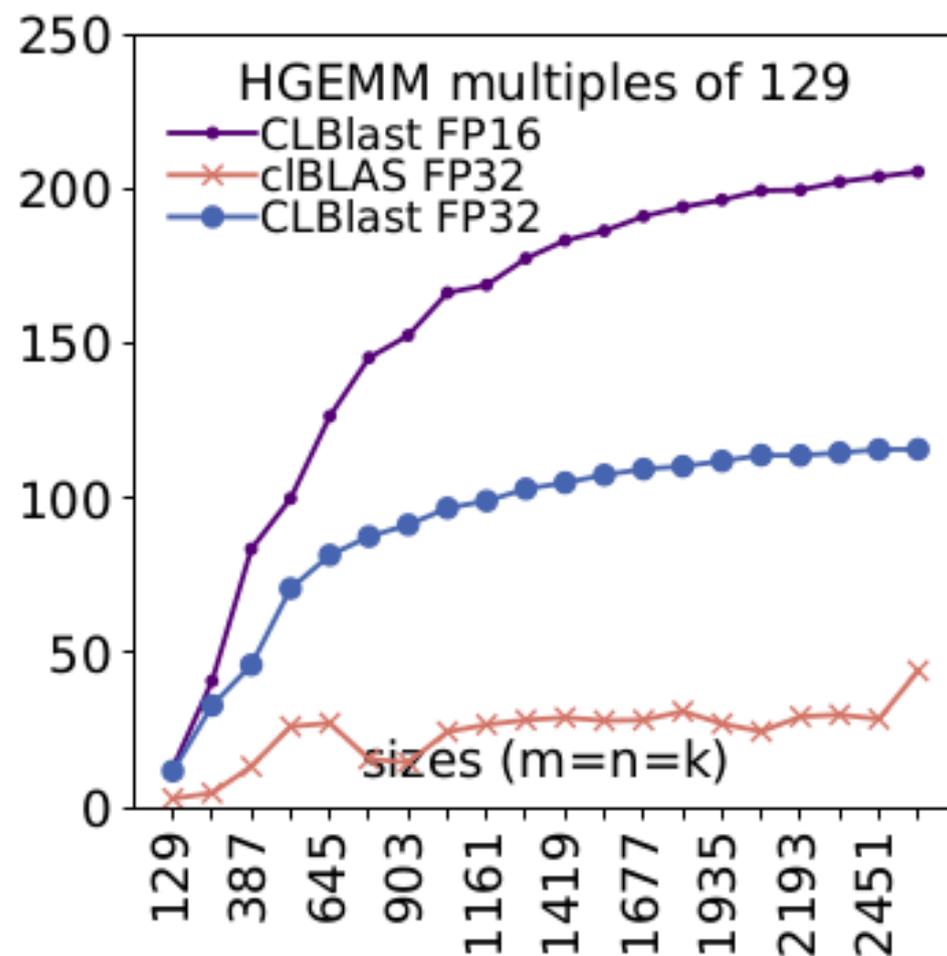- Potential for 2x savings in:
  bandwidth, storage, compute, energy

- Current FP16 support for GPUs:
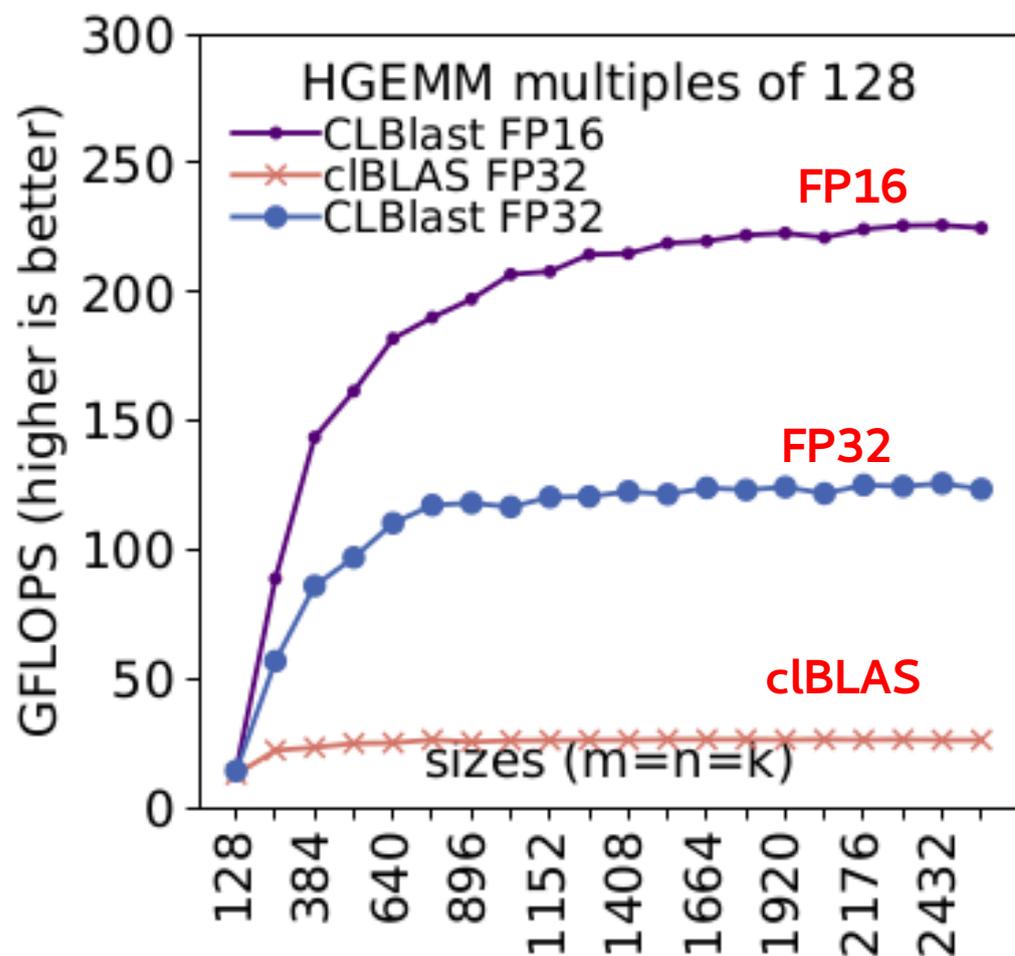  - cuBLAS: HGEMM only
  - clBLAS: no FP16 at all
  - CLBlast: all routines!

- FP16 ~1.9x faster across the board!

- FP16 2+ times faster across the board!

- ## Small-sized GEMM is super slow

  - Not enough work-groups

  - Not enough threads

- Small-sized GEMM is super slow

  - Not enough work-groups

  - Not enough threads



- Let's make it fast again:

  - Combine multiple small GEMM operations into a single kernel

  - Use offsets to indicate where the next matrices start

- SGEMM 128x128x128:
  - Regular: ~60 GFLOPS
  - Batched: ~20 GFLOPS (1 GEMM) up to ~500 GFLOPS (4K)!

- Significant benefits for larger sizes as well
  - mostly beneficial in the range n=64 till 512

- Strided (left) versus regular (right):

  – More assumptions, smaller overhead

  – Better performance for smaller batches

# What's next?

- More features for deep learning (convolution as GEMM, …)

- Auto-tuning using learned models:

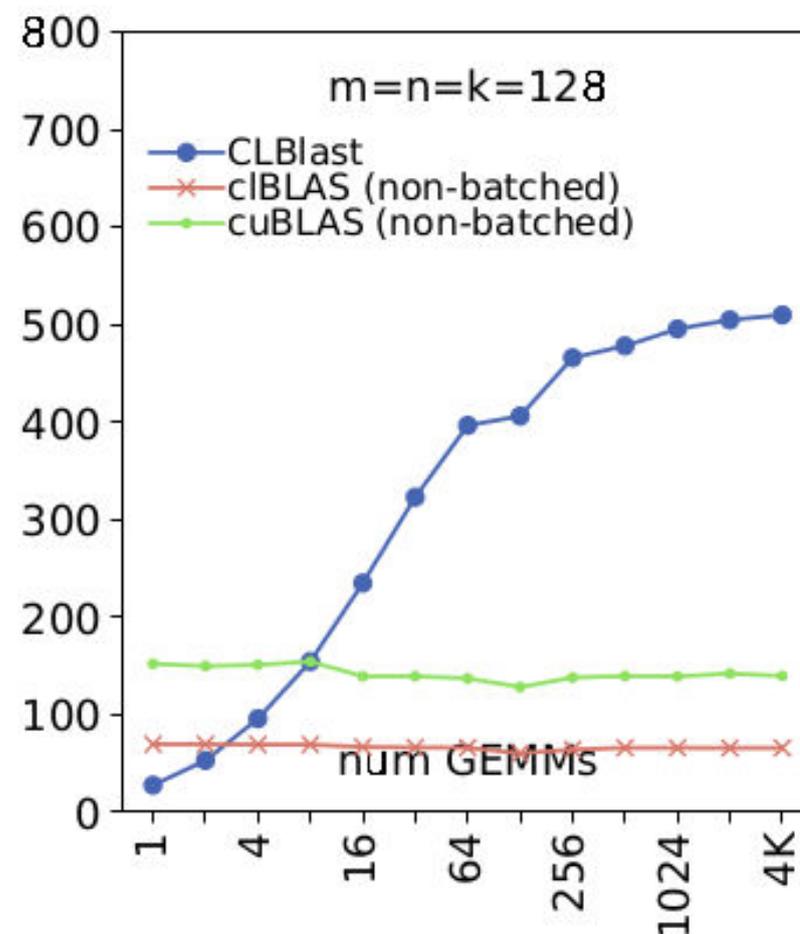  - Similar to the 'ISAAC' library

  - Cooperation with:

    - Rafael Ballester-Ripoll (University of Zürich)
      - Based on tensor trains
    - Flavio Vella (dividiti)
      - Matrix-size aware tuning

- Suggestions?

- **CLBlast**: a modern C++11 OpenCL BLAS library

- Performance portable thanks to generic kernels and **auto-tuning**

- Has important features to accelerate deep-learning:

  - **Problem-size specific tuning**:

    - Up to 2x in an example experiment

  - **Half-precision FP16** support:

    - Up to 2x benefit in speed and memory savings

  - **Batched GEMM** routines:

    - Order of magnitude benefit depending on the use-case

- Time to checkout the album…. Or clone the repository!

# CLBlast: A Tuned BLAS Library

Cedric Nugteren
May 16, 2018

http://github.com/cnugteren/clblast
http://cnugteren.github.io/clblast