

The Ginkgo Sparse Linear Algebra Package



Overview

- Open-source C++ framework for sparse linear algebra.
- Sparse linear solvers, preconditioners, SpMV etc.
- **Generic algorithm implementation:**
 - + reference kernels for checking correctness;
 - + architecture-specific highly optimized kernels.
- Focused on Multicore and Manycore accelerators;
- Software quality and sustainability efforts guided by xSDK community policies:



<https://xsdk.info/>

- Static polymorphism for templating precisions
 - **ValueType** (By default: Z,C,D,S), Integer
 - Smart pointers to avoid memory leaks
 - Runtime polymorphism for **operators** and **kernels**
 - Kernels have the same signature for different architectures
 - **Executor** determines which kernel is used
- Determines where Data lives & operations are executed
- **LinOp** class for any linear operator:
 - Matrices
 - Solvers
 - Preconditioners
 - ...
- generate apply ...



Core Concept: Separate Algorithm from Kernels



Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

- SpMV
- Solver kernels
- Precond kernels
- ...

Core

- Library Infrastructure
- Algorithm Implementations
 - Iterative Solvers
 - Preconditioners
 - ...

Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
- SpMV
- Solver kernels
- Precond kernels
- ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;

Core

- Library Infrastructure
- Algorithm Implementations
 - Iterative Solvers
 - Preconditioners
 - ...

Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Runtime polymorphism selects the right kernel depending on the target architecture;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
- SpMV
- Solver kernels
- Precond kernels
- ...

CUDA

- NVIDIA-GPU kernels
- SpMV
- Solver kernels
- Precond kernels
- ...

OpenMP

- OpenMP-kernels
- SpMV
- Solver kernels
- Precond kernels
- ...

HIP

- AMD-GPU kernels
- SpMV
- Solver kernels
- Precond kernels
- ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;

Optimized architecture-specific kernels;

Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Runtime polymorphism selects the right kernel depending on the target architecture;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

CUDA

- NVIDIA-GPU kernels
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

OpenMP

- OpenMP-kernels
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

HIP

- AMD-GPU kernels
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;

Core

Library Infrastructure
Algorithm Implementations

- Iterative Solvers
- Preconditioners
- ...

Optimized architecture-specific kernels;



Core Concept: Linear Operators

We express everything as Linear Operator and leverage C++ class inheritance.



Core Concept: Linear Operators



We express everything as Linear Operator and leverage C++ class inheritance.

Matrix-Vector Product

$$x := A \cdot b$$

Core Concept: Linear Operators



We express everything as *Linear Operator* and leverage C++ class inheritance.

Matrix-Vector Product

$$x := A \cdot b$$

Preconditioner (for matrix A)

$$x := M^{-1} \cdot b$$

$$M^{-1} \approx A^{-1}$$

Core Concept: Linear Operators

We express everything as *Linear Operator* and leverage C++ class inheritance.



Matrix-Vector Product

$$x := A \cdot b$$

Preconditioner (for matrix A)

$$x := M^{-1} \cdot b$$

$$M^{-1} \approx A^{-1}$$

Solver (for system $Ax = b$)

$$x := S \cdot b$$

$$S \approx A^{-1}$$

Core Concept: Linear Operators



We express everything as *Linear Operator* and leverage C++ class inheritance.

Matrix-Vector Product

$$x := A \cdot b$$

Preconditioner (for matrix A)

$$x := M^{-1} \cdot b$$

$$M^{-1} \approx A^{-1}$$

$$M^{-1} = \Pi(A)$$

Solver (for system $Ax = b$)

$$x := S \cdot b$$

$$S \approx A^{-1}$$

$$S = \Sigma(A)$$

Core Concept: Linear Operators



We express everything as *Linear Operator* and leverage C++ class inheritance.

Matrix-Vector Product

$$x := A \cdot b$$

Preconditioner (for matrix A)

$$x := M^{-1} \cdot b$$

Solver (for system $Ax = b$)

$$x := S \cdot b$$

$$M^{-1} \approx A^{-1}$$

$$S \approx A^{-1}$$

$$M^{-1} = \Pi(A)$$

$$S = \Sigma(A)$$

All of them can be expressed as

- Application of a linear operator* (LinOp) $L : \mathbb{F}^m \rightarrow \mathbb{F}^m$

Core Concept: Linear Operators



We express everything as *Linear Operator* and leverage C++ class inheritance.

Matrix-Vector Product

$$x := A \cdot b$$

Preconditioner (for matrix A)

$$x := M^{-1} \cdot b$$

Solver (for system $Ax = b$)

$$x := S \cdot b$$

$$M^{-1} \approx A^{-1}$$

$$S \approx A^{-1}$$

$$M^{-1} = \Pi(A)$$

$$S = \Sigma(A)$$

All of them can be expressed as

- Application of a linear operator* (LinOp) $L : \mathbb{F}^m \rightarrow \mathbb{F}^m$
- (non-linear) transformation applied to linear operator (Factory) $\Phi : \mathbb{L}^{mn}(\mathbb{F}) \rightarrow \mathbb{L}^{kl}(\mathbb{F})$

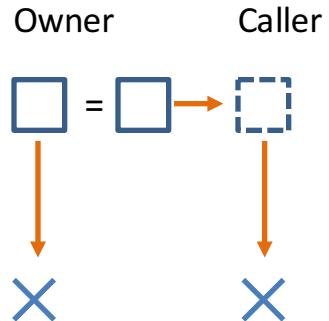
Memory Management: Smart Pointers



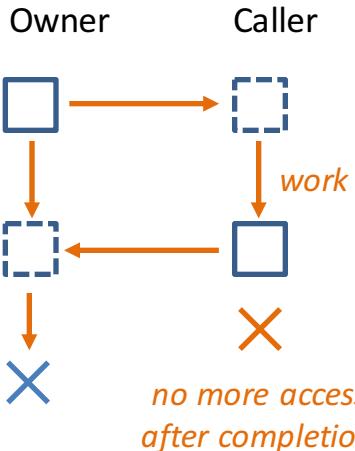
actual object

"empty container"

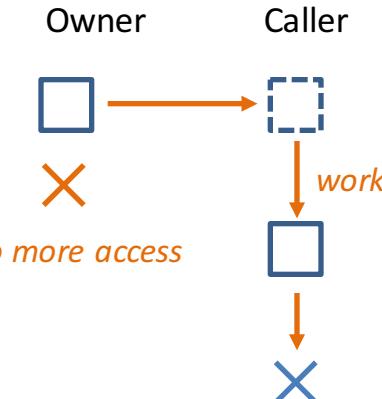
Clone



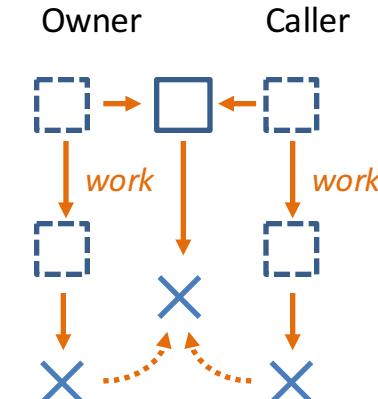
Lend



Give



Share



Ginkgo Functionality (version 1.1.0)



Matrix Formats

- Dense
- CSR
- COO
- ELL
- SELLP
- hybrid

Hardware Backends

- Reference (sequential, CPU)
- OpenMP
- CUDA
- HIP

Krylov Solvers

- BiCGSTAB
- CG
- CGS
- Flexible CG
- GMRES

Supported OS

- Linux
- MacOS
- Windows

Preconditioners

- (Block-) Jacobi
- Adaptive Precision Block-Jacobi
- ILU(0)

Software Integration

- MFEM
- deal.II



Ginkgo Functionality (version 1.1.0)



Include Ginkgo functionality

```
1 #include <deal.II/lac/ginkgo_solver.h>
2 #include <deal.II/lac/sparse_matrix.h>
3 #include <deal.II/lac/vector.h>
4 #include <deal.II/lac/vector_memory.h>
5
6 #include "../testmatrix.h"
7 #include "tests.h"
8
9 #include <iostream>
10 #include <typeinfo>
11
12 int main(int argc, char **argv)
13 {
14     // Set solver parameters
15     SolverControl control(200, 1e-6);
16
17     const unsigned int size = 32;
18     unsigned int dim = (size - 1) * (size - 1);
19
20     // Setup a simple matrix
21     FDMatrix testproblem(size, size);
22     SparsityPattern structure(dim, dim, 5);
23     testproblem.five_point_structure(structure);
24     structure.compress();
25     SparseMatrix<double> A(structure);
26     testproblem.five_point(A);
27
28     Vector<double> f(dim);
29     f = 1.;
30     Vector<double> u(dim);
31     u = 0.;
32
33     // Instantiate a CUDA executor.
34     auto gpu = gko::CudaExecutor::create(0, gko::OmpExecutor::create());
35     // Use ginkgo to solve the system on the gpu using the CG solver.
36     GinkgoWrappers::SolverCG solver(control, gpu );
37     // Solves the system and copies the data back to deal.II's solution variable.
38     solver.solve(A, u, f);
39 }
```

Control Ginkgo solver

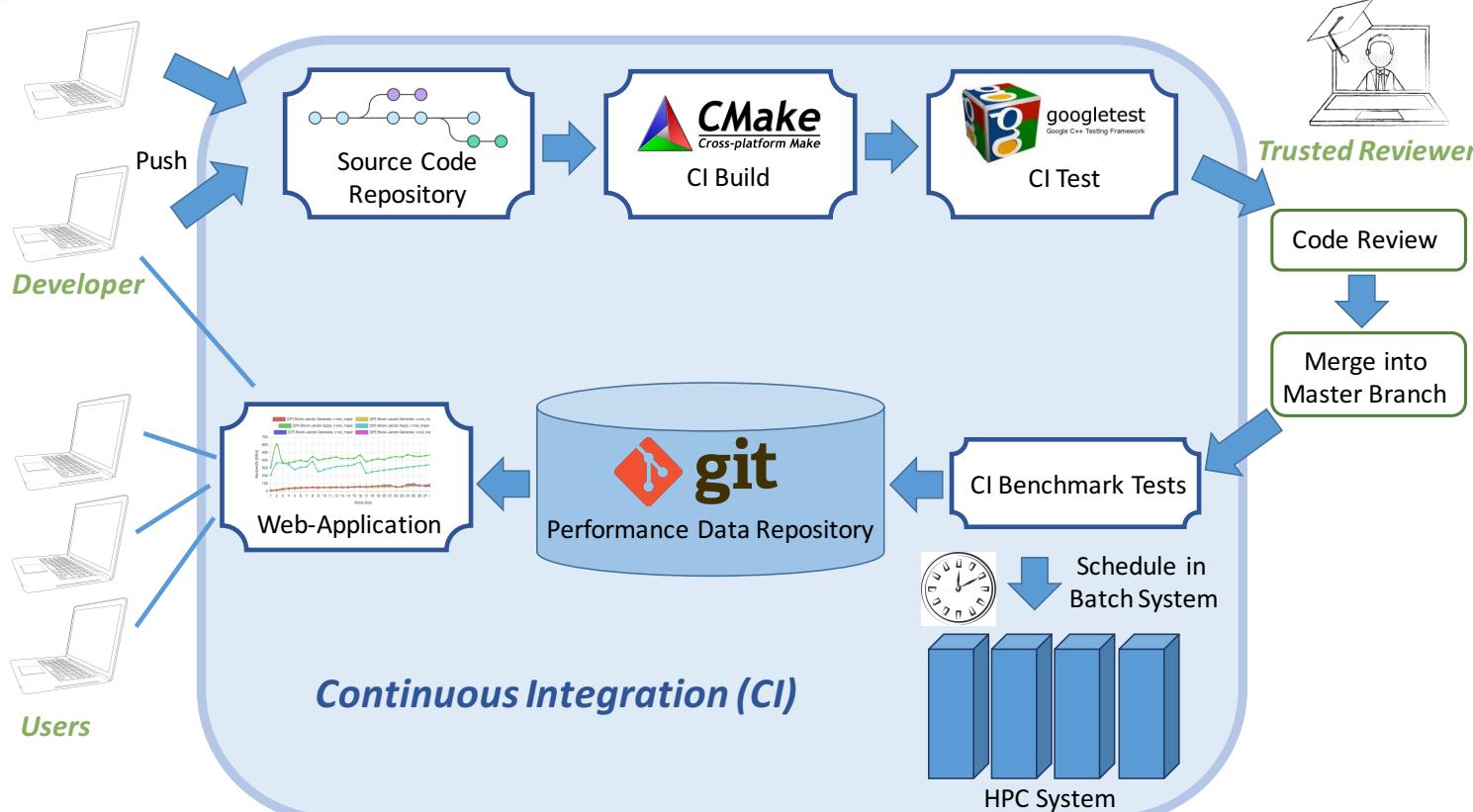
Generate deal.ii problem

Generate Ginkgo executor

Generate Ginkgo solver

Solve linear system

Continuous Integration & Continuous Benchmarking



Ginkgo Performance Explorer (GPE)

Allows anyone to access & visualize Ginkgo Performance data in a browser.



Ginkgo Performance Explorer

Step 1: Select benchmark results
Select raw benchmark results to import and view.

Select result files
Result Summary

Result files to use in the next steps.
Performance data root URL (Advanced) *
<https://raw.githubusercontent.com/ginkgo-project/ginkgo-data/master/data> ↴
URL to a folder containing a `list.json` file.

Step 2: Transform results
For plotting, create a `Chart.js config object`. Use `JSONata` to extract interesting parts of raw data.

Select an example Use an example transformation script as a starting point

```
7 {  
8   "type": "bar",  
9   "data": {  
10     "labels": $formats,  
11     "datasets": [{  
12       "data": $counts,  
13       "backgroundcolor": $formats->$map(function  
14         "hsl(" & 360 * $i / ($a->$count()) & ",40  
15       )}  
16     ]  
17   },  
18   "options": {  
19     "legend": { "display": false },  
20     "title": {  
21       "display": true,  
22     }  
23   }  
24 }
```

Step 3: View transformed results
View the resulting plot, or raw transformed data.

Results Transformed Plot

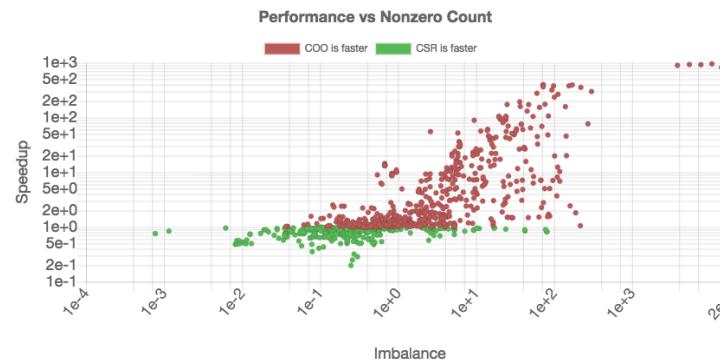
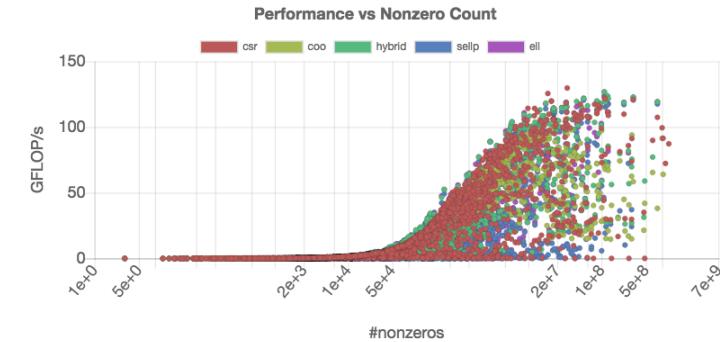
Best SpMV format

Problems

GFLOP/s

1000
800
600
400
200
0

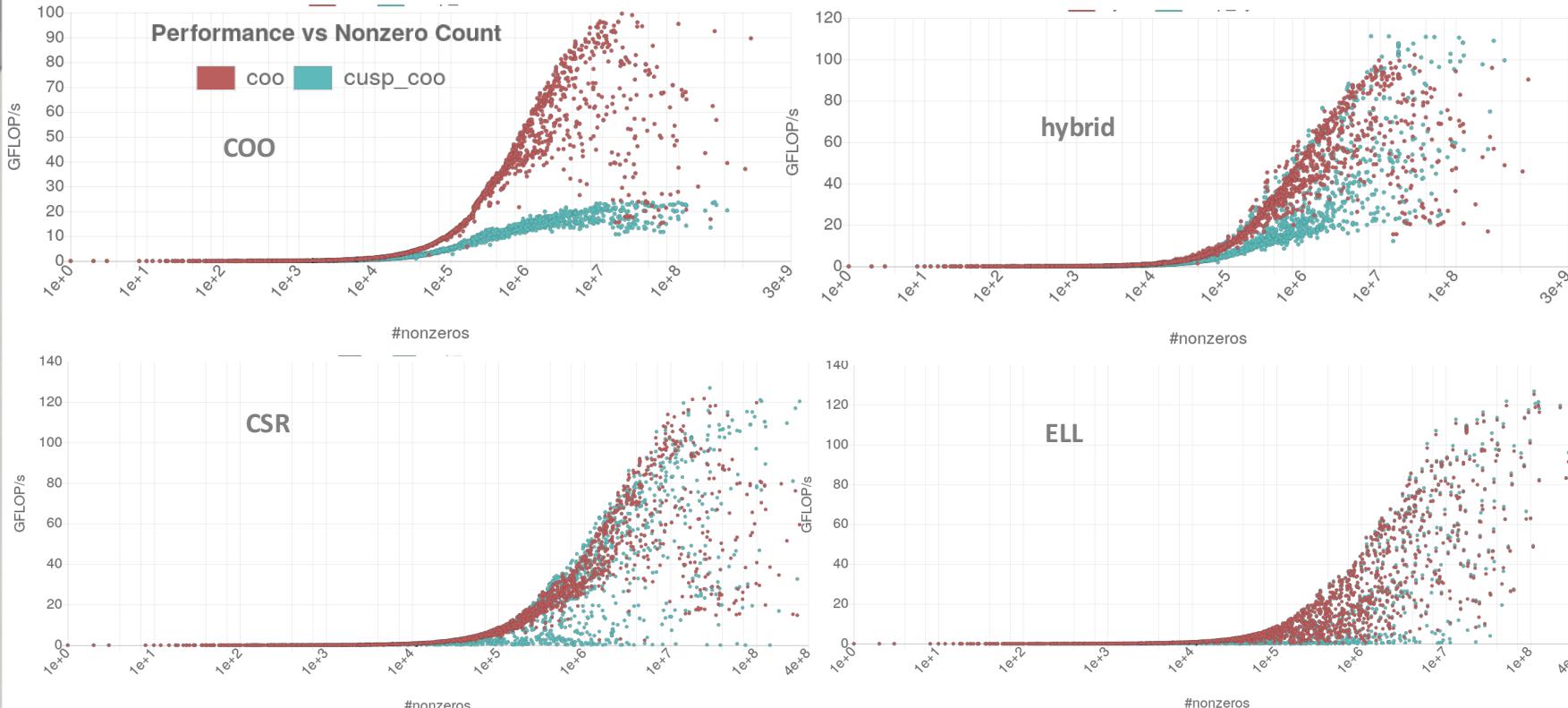
csr coo hybrid selp ell



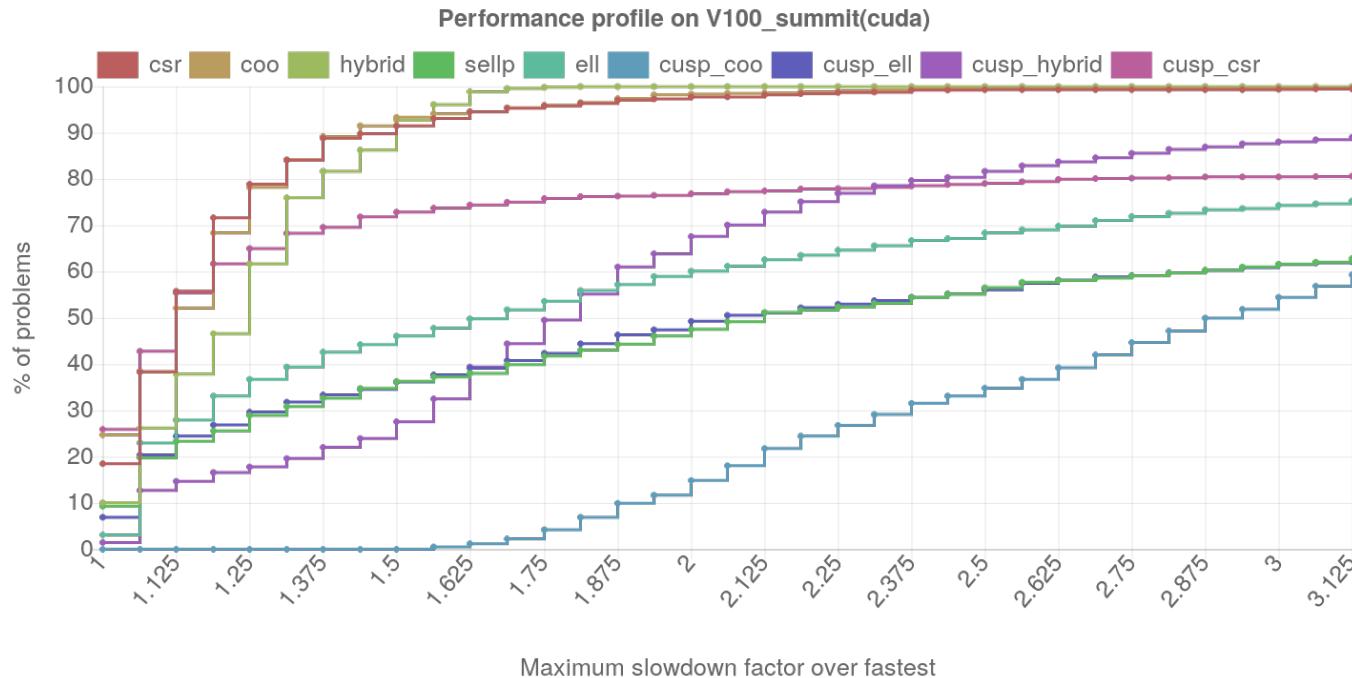
<https://ginkgo-project.github.io/gpe/>



Ginkgo SpMV Performance CUDA on V100



Ginkgo SpMV Performance CUDA on V100

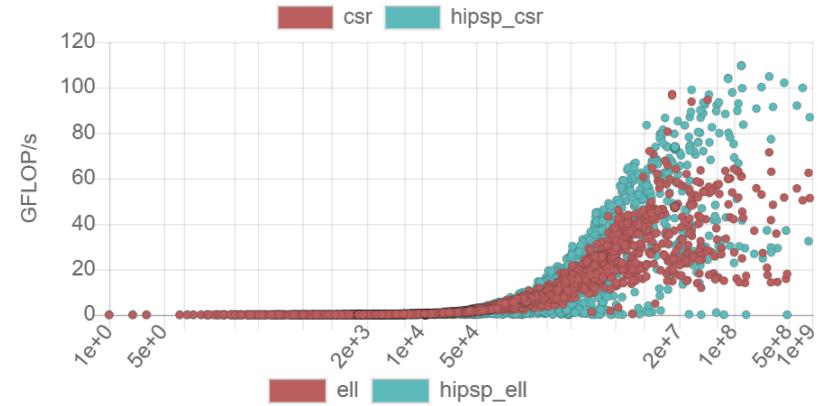


Anzt, Cojean, Chen, Dongarra, Flegar, Nayak, Tomov, Tsai, Wang: Load-balancing Sparse Matrix Vector Product Kernels on GPUs, TOPC, 2019.

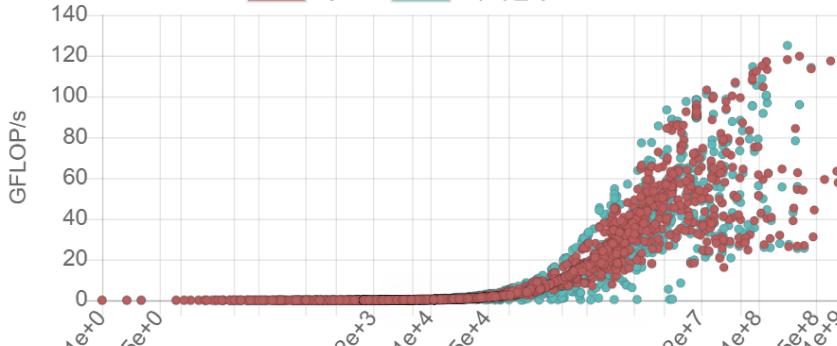
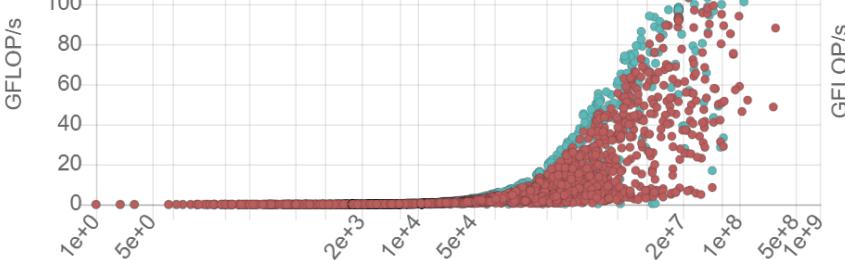
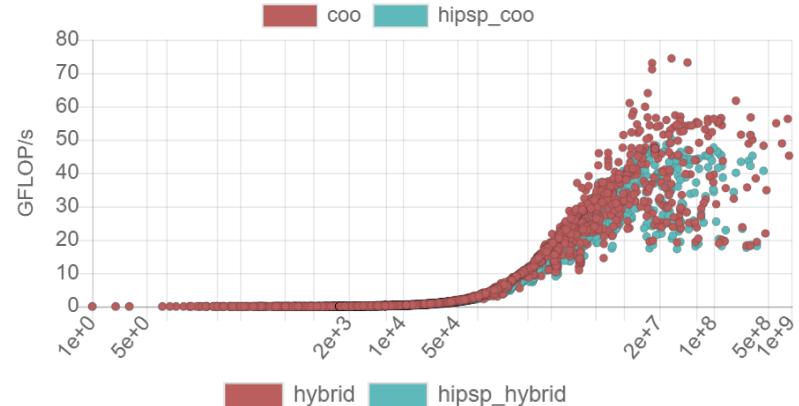
Ginkgo SpMV Performance HIP on Radeon 7



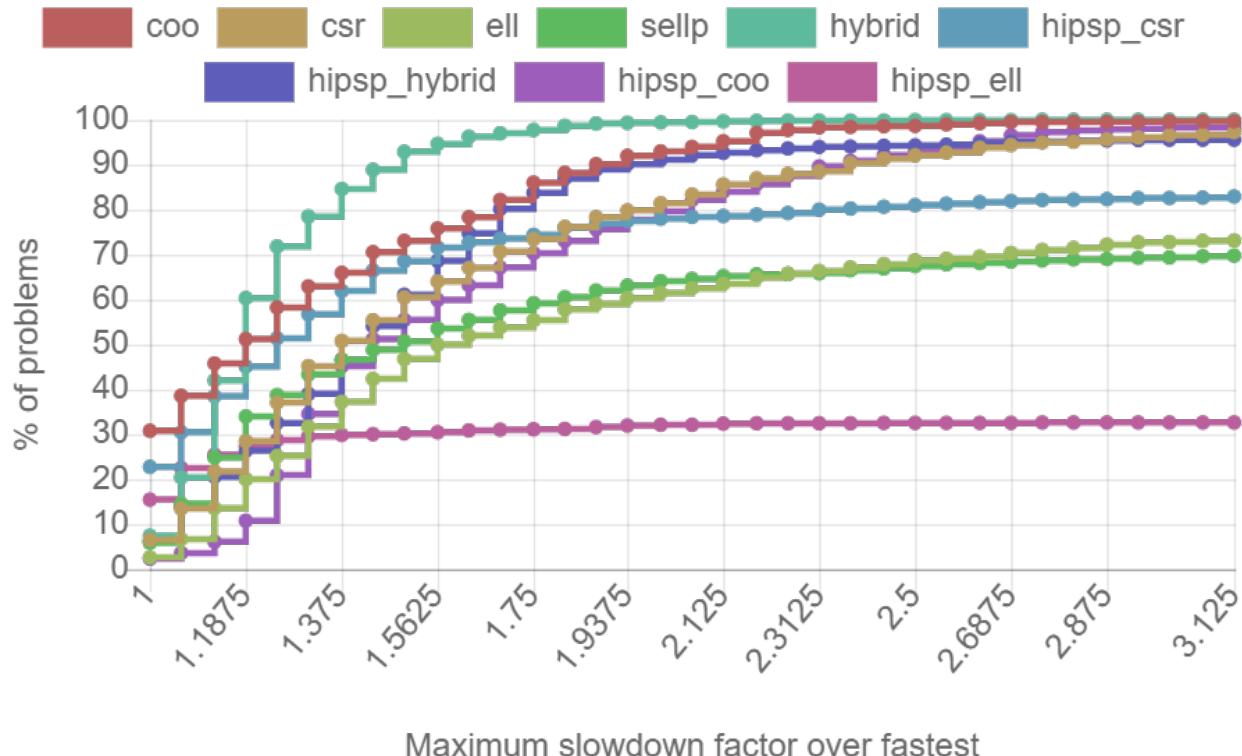
Performance vs Nonzero Count



Performance vs Nonzero Count



Ginkgo SpMV Performance HIP on Radeon 7



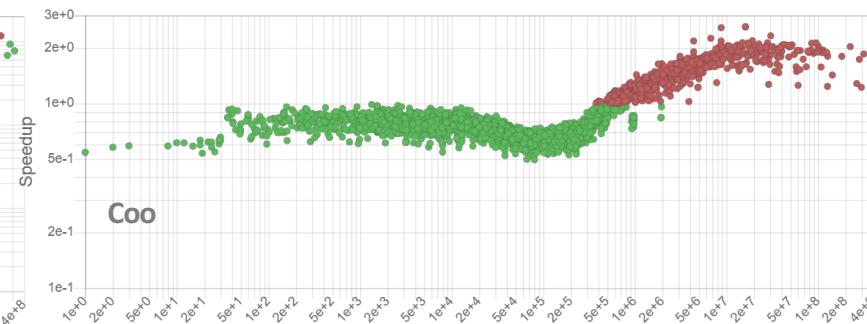
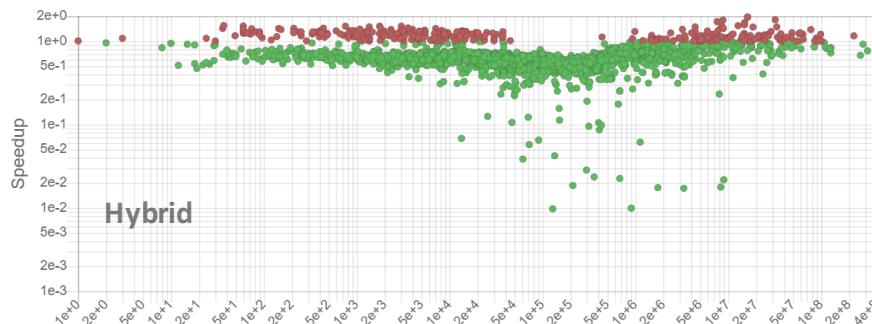
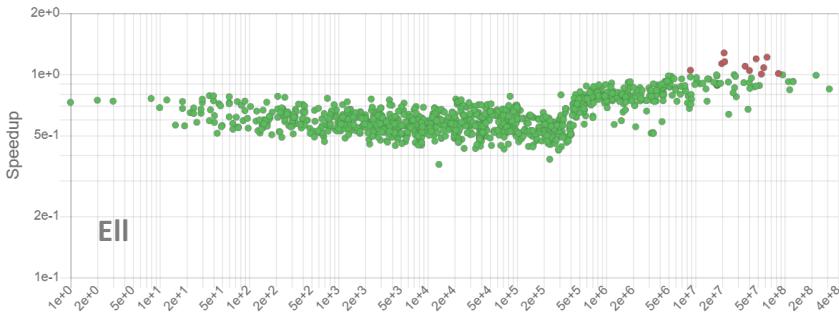
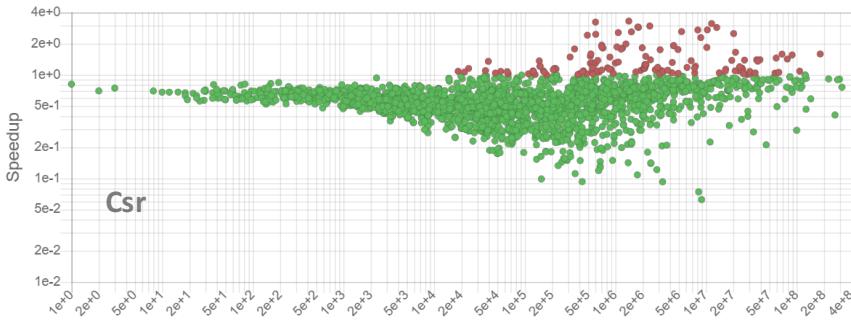
Maximum slowdown factor over fastest

SpMV Performance Comparison V100 vs. Radeon 7



cuSPARSE on V100 is faster

HIPSparse on Radeon 7 is faster

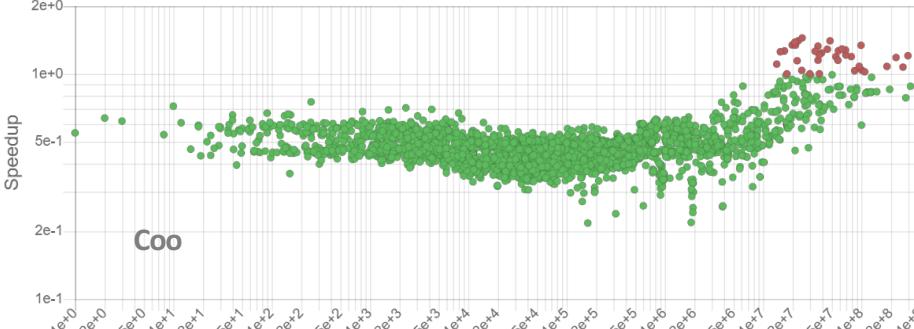
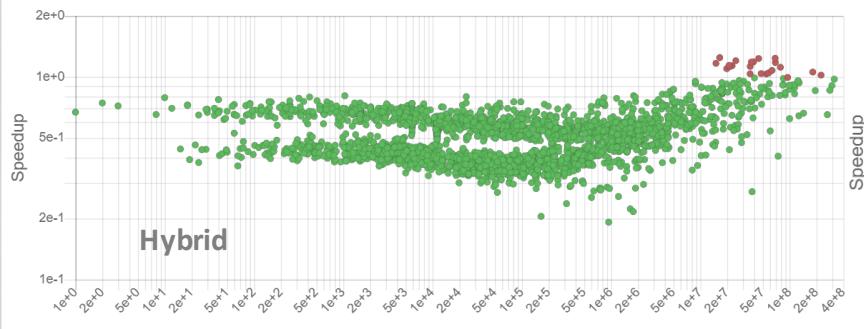
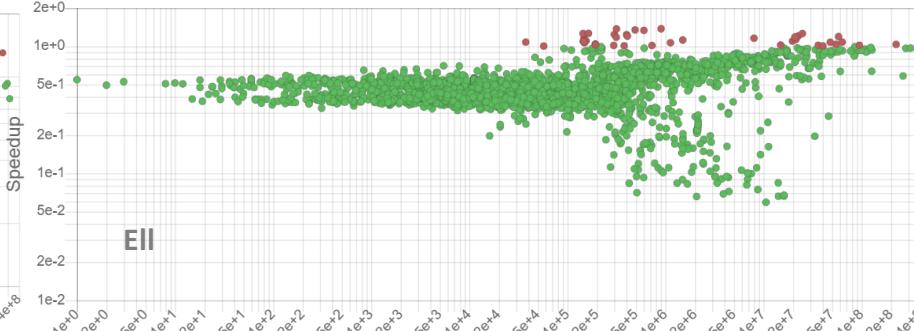
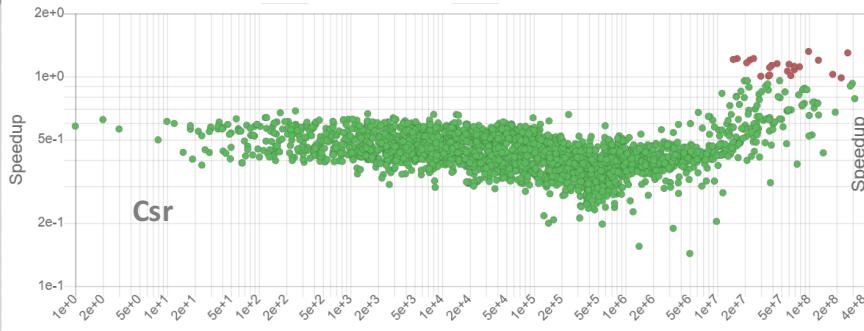


SpMV Performance Comparison V100 vs. Radeon 7



Ginkgo on V100 is faster

Ginkgo on Radeon 7 is faster



Past, Present, Future...

- 05/2017 Start library design discussions
- 05/2017 Initial C++ usage tests
- 10/2017** "First bits in the bucket" <https://github.com/ginkgo-project/ginkgo>
- 11/2017 Ginkgo name decision
- 02/2018 Continuous Integration (CI)
- 07/2018 Ginkgo Performance Explorer (GPE) <https://ginkgo-project.github.io/gpe/>
- 09/2018 Continuous Benchmarking (CB)
- 09/2018 deal.II Finite Element Software integration example
- 12/2018 Ginkgo Project Webpage <https://ginkgo-project.github.io/>
- 01/2019 PAPI SDE integration
- 05/2019** Ginkgo v.1.0.0 release
- 07/2019** Ginkgo approved for xSDK release v.0.5.0
- 10/2019** Ginkgo v.1.1.0 release
- 10/2019** MFEM Finite Element Software integration example
- 11/2019** *HIP executor*
- 01/2020 *PETSc integration*
- 01/2020 *Trilinos integration*
- 06/2020 *FEAST sparse eigensolver*
- 06/2020 *Multi-GPU executor?*



Past, Present, Future...

- 05/2017 Start library design discussions
- 05/2017 Initial C++ usage tests
- 10/2017** "First bits in the bucket" <https://github.com/ginkgo-project/ginkgo>
- 11/2017 Ginkgo name decision
- 02/2018 Continuous Integration (CI)
- 07/2018 Ginkgo Performance Explorer
- 09/2018 Continuous Benchmarking (CI)
- 09/2018 deal.II Finite Element Software
- 12/2018 Ginkgo Project Webpage <http://ginkgo-project.org>
- 01/2019 PAPI SDE integration
- 05/2019** Ginkgo v.1.0.0 release
- 07/2019** Ginkgo approved for xSDK release
- 10/2019** Ginkgo v.1.1.0 release
- 10/2019** MFEM Finite Element Software
- 11/2019** HIP executor
- 01/2020 PETSc integration
- 01/2020 Trilinos integration
- 06/2020 FEAST sparse eigensolver
- 06/2020 Multi-GPU executor?



Join the Pack!