

Enhancing GPU-Accelerated Scientific Computing in Julia with Ginkgo.jl

<u>You Wu</u>, Tobias Ribizel, and Hartwig Anzt



05-06-2024



Table of Content





Ginkgo: A Modern Linear Operator Algebra Framework for HPC



Ginkgo

- C++ framework for sparse numerical linear algebra
- Implemented using modern C++
- Part of the <u>Extreme-Scale Scientific Software Stack (E4S)</u> and the <u>Extreme-Scale Software Development Kit (xSDK)</u>
- Emphasis on solving sparse linear systems efficiently on GPUaccelerated systems, using backends written in the respective vendor languages





Designed for Platform Portability

Library core contains architecture-agnostic functionality Infrastructure Algorithms • Iterative Solvers

Preconditioners

CORE

...

•



Runtime polymorphism selects the right kernel depending on the target architecture





FUNCTIONALITY		OMP	CUDA	HIP	DPC++
Basic	SpMV	۲	۲	۲	۲
	SpMM		۲	۲	۲
	SpGeMM	Ì	ø	S	Ø
Krylov solvers	BiCG	Ø	Ø	V	Ø
	BICGSTAB	۲	۲	۲	۲
	CG	۲	۲	۲	۲
	CGS	ø	ø	ø	Ø
	GCR	۲	۲	۲	Ø
	GMRES	۲	۲	۲	۲
	FCG	۲	۲	۲	Ø
	FGMRES	۲	۲	۲	۲
	IR	۲	۲	۲	۲
	IDR	ø	Ø	V	Ø
ditioners	Block-Jacobi	۲	۲	۲	۲
	ILU/IC	ø	Ø	Ø	Ø
	Parallel ILU/IC	Ì	ø	T	Ø
con	Parallel ILUT/ICT	Ø	Ø	ø	Ø
Pre	ISAI	ø	ø	ø	Ø
Batched	Batched BiCGSTAB	۲	۲	۲	۲
	Batched CG	۲	۲	۲	۲
	Batched GMRES	۲	۲	۷	۲
	Batched ILU	۲	۲	۲	۲
	Batched ISAI		۲	۷	۲
	Batched Block-Jacobi	۲	۲	۲	۲
AMG	AMG preconditioner	Ì	Ø	I	Ø
	AMG solver	ø	Ø	Ø	Ø
	Parallel Graph Match	Ì	۲	ø	Ø
arse direct	Symbolic Cholesky	Ì	Ø	I	Ø
	Numeric Cholesky	Ś	ø	V	
	Symbolic LU	ø	Ø	Ø	Ø
	Numeric LU	Ì	ø	ø	
Sp	Sparse TRSV	I	S	V	
Utilities	On-Device Matrix Assembly	Í	ø	T	Ø
	MC64/RCM reordering	Ø			
	Wrapping user data	۲	۲	۲	۲
	Logging	I	S	Ø	Ø
	PAPI counters	Ì	Ø	V	V
	✓ MPI Support	𝗭 Single-GPU Support			

Designed for Performance





[1] Tsai Y-HM, Cojean T, Anzt H. Providing performance portable numerics for Intel GPUs.



C API & Julia Bindings



C API & Maintainable Software Development

- Implementation of a C API is more versatile and can be extended to support other languages
- Using packages such as CxxWrap.jl adds on an extra dependencies
- Calling C++ library in Julia consists of mainly three parts:



C API Implementation

- Uses the language linkage with key word **extern "c"**
- Difficulty lies in "translating" back modern C++ syntax back to C
- Appropriate use of macros largely avoid code duplication
- Naming convention should be well-established

// c_api.h	
<pre>#define DECLARE_CSR_OVERLOAD(_ctype_value, _ctype_index, _cpptype_value,</pre>	\
cpptype index, name, name dense)	\
<pre>struct gko matrix csr ## name## st;</pre>	\
<pre>typedef struct gko matrix csr ## name## st* gko matrix csr ## name;</pre>	\
gko matrix csr ## name ginkgo matrix csr ## name## create(\
gko executor exec, gko dim2 st size, size t nnz);	\
gko matrix csr ## name ginkgo matrix csr ## name## create view(\
gko executor exec, gko dim2 st size, size t nnz,	١
<pre>_ctype_index* row_ptrs, _ctype_index* col_idxs, _ctype_value* value</pre>	s); \
<pre>void ginkgo_matrix_csr_##_name##_delete(</pre>	\
<pre>gko_matrix_csr_##_name mat_st_ptr);</pre>	\
gko_matrix_csr_##_name ginkgo_matrix_csr_##_name##_read(\
<pre>const char* str_ptr, gko_executor exec);</pre>	\
<pre>void ginkgo_write_csr_##_name##_in_coo(const char* str_ptr,</pre>	\
gko_matrix_csr_##_name mat_st_pt	r); \
<pre>size_t ginkgo_matrix_csr_##_name##_get_num_stored_elements(</pre>	\
<pre>gko_matrix_csr_##_name mat_st_ptr);</pre>	\
<pre>size_t ginkgo_matrix_csr_##_name##_get_num_srow_elements(</pre>	\
gko_matrix_csr_##_name mat_st_ptr);	\
gko_dim2_st ginkgo_matrix_csr_##_name##_get_size(\
gko_matrix_csr_##_name mat_st_ptr);	\
<pre>const _ctype_value* ginkgo_matrix_csr_##_name##_get_const_values(</pre>	١
gko_matrix_csr_##_name mat_st_ptr);	\
<pre>const _ctype_index* ginkgo_matrix_csr_##_name##_get_const_col_idxs(</pre>	\
gko matrix csr ## name mat st ptr);	\

Similar construct holds for DEFINE_X_OVERLOAD that resides in c_api.cpp



// c_api.h
DECLARE_CSR_OVERLOAD(float, int, float, int, f32_i32, f32)
DECLARE_CSR_OVERLOAD(float, int64_t, float, std::int64_t, f32_i64, f32)
DECLARE_CSR_OVERLOAD(double, int, double, int, f64_i32, f64)
DECLARE_CSR_OVERLOAD(double, int64_t, double, std::int64_t, f64_i64, f64)

C API Implementation



```
// c api.cpp
struct gko linop st {
    std::shared ptr<gko::Lin0p> shared ptr;
};
void ginkgo linop delete(gko linop linop st ptr) { delete linop st ptr; }
void ginkgo linop apply(gko linop solver, gko linop b st ptr, gko linop x st ptr)
{ (solver->shared ptr)->apply(b st ptr->shared ptr, x st ptr->shared ptr); }
// Solver as a LinOp object
gko linop ginkgo linop gmres preconditioned f64 create(
    gko executor exec st ptr, gko linop A st ptr,
    gko deferred factory parameter dfp st ptr, double reduction, int maxiter)
{
    return new gko linop st{
        gko::solver::Gmres<double>::build()
            .with criteria(
                gko::stop::Iteration::build().with max iters(maxiter),
                gko::stop::ResidualNorm<double>::build().with reduction factor(
                    reduction))
            .with preconditioner(dfp st ptr->parameter)
            .on(exec st ptr->shared ptr)
            ->generate(A st ptr->shared ptr)};
}
```

Example: Creation of a preconditioned GMRES solver with the C API

JLL Package Creation

- JLL packages live under JuliaBinaryWrappers organization
- With build scripts hosted under JuliaPackaging/Yggdrasil
- **BinaryBuilder.jl** is a convenient tool for creating suitable build scripts and creating pull requests to official Julia registry for JLL packages





JLL Package Creation





C Bindings Generation



An example wrap.jl file for generating src/api.jl
using Clang.Generators

import ginkgo_jll
header_dir = joinpath(ginkgo_jll.artifact_dir, "include")
isdir(header_dir) || error("\$header_dir does not exist")

wrapper generator options
options = load_options(joinpath(@_DIR_, "wrap.toml"))

add compiler flags, e.g. "-DXXXXXXXX"
args = get_default_args()
push!(args, "-I\$header_dir")
push!(args, "-DBUILD_SHARED_LIBS")

specifying C API header to parse headers = [joinpath(header_dir, header) for header in readdir(header_dir) if startswith(header, "c_api")]

create context
ctx = create_context(headers, args, options)

run generator
build!(ctx)





High-Level API Design

Code Comparison

using Ginkgo

Type alias
const (Tv, Ti) = (Float64, Int32)

Print ginkgo library version
version()

Creates executor for a specific backend const exec = create(:cuda)

Specify executor to be passed
with(EXECUTOR => exec) do
 # Read matrix and vector from mtk files, now omit the passing of exec
 A = GkoCsr{Tv, Ti}("data/A.mtx");
 b = GkoDense{Tv}("data/b.mtx");
 x = GkoDense{Tv}("data/x0.mtx");

colver - ChalterativeSolver(.cr. A. maviter - 20 m

solver = GkoIterativeSolver(:cg, A; maxiter = 20, reduction = 1.0e-7)
apply!(solver, b, x)

@info "Solution (x):"

display(x)

one = number(Tv(1.0))
neg_one = number(Tv(-1.0))
res = number(Tv(0.0))

apply!(A, one, x, neg_one, b)
norm2!(b, res)
@info "Residual norm sqrt(r^T r):"
display(res)

#include <ginkgo/ginkgo.hpp>
#include <fstream>
#include <iostream>

int main(int argc, char* argv[])

// Type alias
using Tv = double;
using Ti = int;
using vec = gko::matrix::Dense<Tv>;

// Print ginkgo library version
std::cout << gko::version_info::get() << std::endl;</pre>

// Creates executor for a specific backend const auto exec = gko::CudaExecutor::create(0, gko::OmpExecutor::create());

// Read matrix and vector from mtk files, now omit the passing of exec auto A = gko::share(gko::read<gko::matrix::Csr<Tv, Ti>> (std::ifstream("data/A.mtx"), exec)); auto b = gko::read<vec>(std::ifstream("data/b.mtx"), exec); auto x = gko::read<vec>(std::ifstream("data/x0.mtx"), exec);

auto solver = solver_gen->generate(A); solver->apply(b, x);

std::cout << "Solution (x):\n";
write(std::cout, x);</pre>

auto one = gko::initialize<vec>({1.0}, exec); auto neg_one = gko::initialize<vec>({-1.0}, exec); auto res = gko::initialize<gko::matrix::Dense<Tv>>({0.0}, exec); A->apply(one, x, neg_one, b); b->compute_norm2(res);

std::cout << "Residual norm sqrt(r^T r):\n";
write(std::cout, res);</pre>





Convenient Switch of Underlying Binaries

Using underlying binaries with user-defined location

julia> Ginkgo.GkoPreferences.locate()

Info: Location of binaries being used by Ginkgo.jl

ginkgo_jll.find_artifact_dir() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219"
ginkgo_jll.get_libginkgo_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo.so"
ginkgo_jll.get_libginkgo_omp_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_omp.so"
ginkgo_jll.get_libginkgo_cuda_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_cuda.so"
ginkgo_jll.get_libginkgo_hip_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_cuda.so"
ginkgo_jll.get_libginkgo_hip_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_tuda.so"
ginkgo_jll.get_libginkgo_dpcpp_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_hip.so"
ginkgo_jll.get_libginkgo_dpcpp_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_hip.so"
ginkgo_jll.get_libginkgo_dpcpp_path() = "/home/wyou/misc/pasc/ginkgo_forked/ginkgo/build/lib/libginkgo_hip.so"

julia> Ginkgo.GkoPreferences.use_jll_binary!()

 Warning: ginkgo library paths changed, you will need to restart Julia for the change to take effect ginkgo_jll.find_artifact_dir() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219"
 @ Ginkgo.GkoPreferences ~/misc/pasc/Ginkgo.jl/src/preferences.jl:127

Using underlying binaries from the JLL package

julia> using Ginkgo; Ginkgo.GkoPreferences.locate()

[Info: Location of binaries being used by Ginkgo.jl

ginkgo_jll.find_artifact_dir() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219"

- ginkgo_jll.get_libginkgo_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_so"
 ginkgo_jll.get_libginkgo_reference_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_reference.so"
 ginkgo_jll.get_libginkgo_omp_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_omp.so"
- ginkgo_jll.get_libginkgo_cuda_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_cuda.so"
- ginkgo_jll.get_libginkgo_hip_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_hip.so"
- ginkgo_jll.get_libginkgo_dpcpp_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_dpcpp.so"
- ginkgo_jll.get_libginkgo_device_path() = "/scratch/wyou/.julia/artifacts/b93004147661eb53ee77f59672c6900759a42219/lib/libginkgo_device.so"



Numerical Experiments



Interoperability with Existing Packages

- Made possible by supporting implicit conversion of data types
- Ferrite.jl is a simple FEM toolbox written in Julia
- Using **Ferrite.jl** for FEM-based matrix assembly, and **Ginkgo.jl** for sparse LSE solution with a CG solver





,

2D Heat Equation on Unit Square

$$egin{aligned} & -
abla \cdot (k
abla u) = f, \, x \in \Omega \ u(x) = 0 \ , \, x \in \partial \Omega \end{aligned} egin{aligned} & igodlembox{} & \int_\Omega
abla \delta u \cdot
abla u \, d\Omega = \int_\Omega \delta u \, d\Omega & orall \delta u \in \mathbb{T}, \end{aligned}$$

grid = generate_grid(Quadrilateral, (20, 20));

Trial and test functions

dim = 2ip = Lagrange{dim, RefCube, 1}() qr = QuadratureRule{dim, RefCube}(2) cellvalues = CellScalarValues(qr, ip);

Degrees of freedom dh = DofHandler(grid) add!(dh, :u, 1) close!(dh);

Creates SparseMatrixCSC K = create sparsity pattern(dh)

Essential BC ch = ConstraintHandler(dh); $\partial \Omega = union($ getfaceset(grid, "left"), getfaceset(grid, "right"), getfaceset(grid, "top"), getfaceset(grid, "bottom"),); dbc = Dirichlet(:u, $\partial \Omega$, (x, t) -> 0) add!(ch, dbc); close!(ch) # notifies constraint handler

Assembling the linear system K, f = assemble global(cellvalues, K, dh);

2D Heat Equation on Unit Square



$$\left\{egin{array}{ll} -
abla \cdot (k
abla u) = f, \, x \in \Omega \ u(x) = 0 \ , x \in \partial \Omega \end{array}
ight. egin{array}{ll} \displaystyle igothermatrix & \displaystyle \int_\Omega
abla \delta u \cdot
abla u \, d\Omega = \int_\Omega \delta u \, d\Omega & orall \delta u \in \mathbb{T}, \end{array}
ight.$$

Trial and test functions

dim = 2
ip = Lagrange{dim, RefCube, 1}()
qr = QuadratureRule{dim, RefCube}(2)

qr = QuadratureRule{dim, RefCube}(2)
cellvalues = CellScalarValues(qr, ip);

Degrees of freedom
dh = DofHandler(grid)
add!(dh, :u, 1)
close!(dh);

Creates SparseMatrixCSC
K = create sparsity pattern(dh)

Essential BC
ch = ConstraintHandler(dh);

dn = union(
 getfaceset(grid, "left"),
 getfaceset(grid, "right"),
 getfaceset(grid, "top"),
 getfaceset(grid, "bottom"),
);
dbc = Dirichlet(:u, aΩ, (x, t) -> 0)
add!(ch, dbc);
close!(ch) # notifies constraint handler

Assembling the linear system
K, f = assemble_global(cellvalues, K, dh);

SE24 Zurich J 3-5 June 2024

2D Heat Equation on Unit Square

$$\left\{egin{array}{ll} -
abla \cdot (k
abla u) = f, x \in \Omega \ u(x) = 0 \ , x \in \partial \Omega \end{array}
ight.$$
 $\left. ightarrow \int_{\Omega}
abla \delta u \cdot
abla u \, d\Omega = \int_{\Omega} \delta u \, d\Omega \quad orall \delta u \in \mathbb{T},
ightarrow u(x) = 0 \ , x \in \partial \Omega \end{array}
ight.$





Heat Conservation

- FD method on a staggered grid
- With marker-in-cell techniques and adaptive timesteps
- Dimensionalized
- Ginkgo.jl sparse direct LU solver
- 10 km by 10 km square domain with 56,700 mesh cells



[2] Gerya T. Introduction to Numerical Geodynamic Modelling. 2nd ed. Cambridge University Press; 2019.

Stokes Continuity Equation





Pressure [GPa]





Viscosity $\log(\eta_b) [Pa \cdot s]$

Vx-velocity [m/yr]



Viscosity $\log(\eta_p) [Pa \cdot s]$



Vy-velocity [m/yr]





Conclusions

- Porting **Ginkgo** routines to **Ginkgo.jl** can leverage existing highlyoptimized architecture-specific kernels
- C API approach is versatile but tedious, requires manual work
- For porting packages from other languages, Julia provides a good toolchain
- **Ginkgo.jl** is still under heavy development