# Krylov.jl: A Julia basket of hand-picked Krylov methods

**Alexis Montoison** [1]¶ **and Dominique Orban** [1]

**1** GERAD and Department of Mathematics and Industrial Engineering, Polytechnique Montréal, QC, Canada. ¶ Corresponding author

## Summary

Krylov.jl is a Julia (Bezanson et al., 2017) package that implements a collection of Krylov processes and methods for solving a variety of linear problems:

| Square systems | Linear least-squares problems | Linear least-norm problems |
|---|---|---|
| $Ax = b$ | $\min \|b - Ax\|$ | $\min \|x\|$ subject to $Ax = b$ |

| Adjoint systems | Saddle-point and Hermitian quasi-definite systems | Generalized saddle-point and non-Hermitian partitioned systems |
|---|---|---|
| $\begin{aligned} Ax &= b \\ A^H y &= c \end{aligned}$ | $\begin{bmatrix} M & A \\ A^H & -N \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix}$ | $\begin{bmatrix} M & A \\ B & N \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix}$ |

$A^H$ denotes the conjugate transpose of $A$. It coincides with $A^T$, the transpose of $A$, if $A$ is real. Krylov methods are iterative methods based on Krylov (1931) subspaces. They are an alternative to direct methods such as Gaussian elimination or QR decomposition when storage requirements or computational costs become prohibitive, which is often the case for large and sparse linear problems. Contrary to direct methods, which require storing $A$ explicitly, Krylov methods support linear operators to model operator-vector products $u \leftarrow Av$, and in some instances $u \leftarrow A^H w$ because Krylov processes only require those operations to build Krylov subspaces. The same goes with preconditioners, i.e., transformations that modify a linear system into an equivalent form with favorable spectral properties that may yield faster convergence in finite-precision arithmetic. We refer interested readers to (Ipsen & Meyer, 1998) for an introduction to Krylov methods along with (Greenbaum, 1997) and (Saad, 2003) for more details.

## Statement of need

### Largest collection of Krylov processes and methods

Krylov.jl aims to provide a user-friendly and unified interface for the largest collection of Krylov processes and methods, all programming languages taken together, with six and thirty-five implementations, respectively:

- **Krylov processes**: Arnoldi, Golub-Kahan, Hermitian Lanczos, Montoison-Orban, Non-Hermitian Lanczos, Saunders-Simon-Yip;
- **Krylov methods**: Bicgstab, Bilq, Bilqr, Car, Cg, Cg-lanczos, Cg-lanczos-shift, Cgls, Cgne, Cgs, Cr, Craig, Craigmr, Crls, Crmr, Diom, Dqgmres, Fgmres, Fom, Gmres, Gpmr, Lnlq, Lslq, Lsmr, Lsqr, Minares, Minres, Minres-qlp, Qmr, Symmlq, Tricg, Trilqr, Trimr, Usymlq, Usymqr.

Hence Krylov.jl is a suitable toolbox for easily comparing existing methods with each other as well as new ones. The number of distinct Krylov methods is twenty-two for PETSc (Balay et al., 2023), eleven for MATLAB (2022) and KrylovMethods.jl, nine for IterativeSolvers.jl and three for KrylovKit.jl. However Krylov.jl doesn't have implementations of recycling Krylov methods nor block Krylov methods unlike some alternatives, except for special cases, including TRICG, TRIMR, and GPMR. Note that we only consider the number of Krylov methods that generate different iterates without preconditioning. Variants with preconditioning are not counted except if it is a flexible one such as FGMRES.

Some processes and methods are not available elsewhere and are the product of our own research. References for each process and method are available in the extensive documentation. Beyond the number of methods, Krylov.jl is the only package that offers all of the features that we describe below.

### Support for any floating-point system supported by Julia

Krylov.jl works with real and complex data in any floating-point system supported by Julia, which means that Krylov.jl handles any precision `T` and `Complex{T}` where `T <: AbstractFloat`. Although most personal computers offer IEEE 754 single and double precision computations, new architectures implement native computations in other floating-point systems. In addition, software libraries such as the GNU MPFR, shipped with Julia, let users experiment with computations in variable, extended precision at the software level with the `BigFloat` data type. Working in high precision has obvious benefits in terms of accuracy.

### Support for NVIDIA, AMD and Intel GPUs

Krylov methods are well suited for GPU computations because they only require operator-vector products ($u \leftarrow Av$, $u \leftarrow A^Hw$) and vector operations ($\|v\|$, $u^Hv$, $v \leftarrow \alpha u + \beta v$), which are highly parallelizable. The implementations in Krylov.jl are generic so as to take advantage of the multiple dispatch and broadcast features of Julia. Those allow the implementations to be specialized automatically by the compiler for both CPU and GPU. Thus, Krylov.jl works with GPU backends that build on GPUArrays.jl, including CUDA.jl, AMDGPU.jl and oneAPI.jl, the Julia interfaces to NVIDIA, AMD, and Intel GPUs.

### Support for linear operators

The input arguments of all Krylov.jl solvers that model $A$, $B$, $M$, $N$ and preconditioners can be any object that represents a linear operator. Krylov methods combined with linear operators allow to reduce computation time and memory requirements considerably by avoiding building and storing matrices. In nonlinear optimization, finding a critical point of a continuous function frequently involves linear systems where $A$ is a Hessian or a Jacobian. Materializing such operators as matrices is expensive in terms of operations and memory consumption and is unreasonable for high-dimensional problems. However, it is often possible to implement efficient Hessian-vector and Jacobian-vector products, for example with the help of automatic differentiation tools.

### In-place methods

All solvers in Krylov.jl have an in-place variant that allows to solve multiple linear systems with the same dimensions, precision and architecture. Optimization methods such as the Newton and Gauss-Newton methods can take advantage of this functionality by allocating workspace for the solve only once. The in-place variants only require a Julia structure that contains all the storage needed by a Krylov method as additional argument. In-place methods limit memory allocations and deallocations, which are particularly expensive on GPUs.

**Performance optimizations and storage requirements**

Operator-vector products and vector operations are the most expensive operations in Krylov.jl. The vectors in Krylov.jl are always dense. One may then expect that taking advantage of an optimized BLAS library when one is available on CPU and when the problem data is stored in a supported representation should improve performance. Thus, we dispatch vector-vector operations to BLAS1 routines, and operator-vector operations to BLAS2 routines when the operator is a dense matrix. By default, Julia ships with OpenBLAS and provides multithreaded routines. Since Julia 1.6, users can also switch dynamically to other BLAS backends, such as the Intel MKL, BLIS or Apple Accelerate, thanks to the BLAS demuxing library `libblastrampoline`, if an optimized BLAS is available.

A "Storage Requirements" section is available in the documentation to provide the theoretical number of bytes required by each method. Our implementations are storage-optimal in the sense that they are guaranteed to match the theoretical storage amount. The match is verified in the unit tests by way of functions that return the number of bytes allocated by our implementations.

# Examples

Our first example is a simple implementation of the Gauss-Newton method without linesearch for nonlinear least squares. It illustrates several of the facilities of Krylov.jl: solver preallocation and reuse, genericity with respect to data types, and linear operators. Another example based on a simplistic Newton method without linesearch for convex optimization is also available in the documentation, and illustrates the same concepts in the sections "In-place methods" and "Factorization-free operators".

```julia
using LinearAlgebra      # Linear algebra library of Julia
using SparseArrays       # Sparse library of Julia
using Test               # Test library of Julia
using Krylov             # Krylov methods and processes
using LinearOperators    # Linear operators
using ForwardDiff        # Automatic differentiation
using Quadmath           # Quadruple precision
using MKL                # Intel BLAS

"The Gauss-Newton method for Nonlinear Least Squares"
function gauss_newton(F, JF, x₀::AbstractVector{T}; itmax = 200, tol = √eps(T)) where T
    n = length(x₀)
    x = copy(x₀)
    Fx = F(x)
    m = length(Fx)
    iter = 0
    S = typeof(x)              # precision and architecture
    solver = LsmrSolver(m, n, S) # structure that contains the workspace of LSMR
    solved = tired = false
    while !(solved || tired)
        Jx = JF(x)             # Compute J(xₖ)
        lsmr!(solver, Jx, -Fx) # Minimize ‖J(xₖ)Δx + F(xₖ)‖
        x .+= solver.x         # Update xₖ₊₁ = xₖ + Δx
        Fx_old = Fx            # F(xₖ)
        Fx = F(x)              # F(xₖ₊₁)
        iter += 1
        solved = norm(Fx - Fx_old) / norm(Fx) ≤ tol
        tired = iter ≥ itmax
```

```julia
        end
        return x
end

T = Float128  # IEEE quadruple precision
x_exact = T[8, 0.25]
x₀ = ones(T, 2)
t = T[1, 2, 3, 4, 5, 6, 7, 8]
y = [trunc(x_exact[1] * exp(x_exact[2] * t[i]), digits=3) for i=1:8]
F(x) = [x[1] * exp(x[2] * t[i]) - y[i] for i=1:8]          # F(x)
J(y, x, v) = ForwardDiff.derivative!(y, h → F(x + h * v), 0)  # y ← JF(x)v
Jᵀ(y, x, w) = ForwardDiff.gradient!(y, x → dot(F(x), w), x)   # y ← JFᵀ(x)w
symmetric = hermitian = false
JF(x) = LinearOperator(T, 8, 2, symmetric, hermitian, (y, v) → J(y, x, v),   # non-transpose
                                                       (y, w) → Jᵀ(y, x, w),  # transpose
                                                       (y, w) → Jᵀ(y, x, w)) # conjugate transpose

x = gauss_newton(F, JF, x₀)

# Check the solution returned by the Gauss-Newton method
@test norm(x - x_exact) ≤ 1e-4
```

Our second example concerns the solution of a complex Hermitian linear system from the SuiteSparse Matrix Collection (Davis & Hu, 2011) with an incomplete Cholesky factorization preconditioner on GPU. The preconditioner is implemented as an in-place linear operator that performs the forward and backward sweeps with the Cholesky factor of the incomplete decomposition. Because the system matrix is Hermitian and positive definite, we use the conjugate gradient method. However, other methods for Hermitian systems could be used, including SYMMLQ, CR, and MINRES.

```julia
using LinearAlgebra                 # Linear algebra library of Julia
using SparseArrays                  # Sparse library of Julia
using Test                          # Test library of Julia
using Krylov                        # Krylov methods and processes
using LinearOperators               # Linear operators
using MatrixMarket                  # Reader of matrices stored in the Matrix Market format
using SuiteSparseMatrixCollection   # Interface to the SuiteSparse Matrix Collection
using CUDA                          # Interface to NVIDIA GPUs
using CUDA.CUSPARSE                 # NVIDIA CUSPARSE library

if CUDA.functional()
  ssmc = ssmc_db(verbose=false)
  matrices = ssmc_matrices(ssmc, "Sinclair", "3Dspectralwave2")
  paths = fetch_ssmc(matrices, format="MM")
  path_A = joinpath(paths[1], "3Dspectralwave2.mtx")

  # A is an Hermitian and positive definite matrix of size 292008 x 292008
  A_cpu = MatrixMarket.mmread(path_A) + 50I
  m, n = size(A_cpu)
  x_exact = ones(ComplexF64, m)
  b_cpu = A_cpu * x_exact

  # Transfer the linear system from the CPU to the GPU
  A_gpu = CuSparseMatrixCSR(A_cpu)
  b_gpu = CuVector(b_cpu)

  # Incomplete Cholesky factorization LLᴴ ≈ A with zero fill-in
```

Montoison, & Orban. (2023). Krylov.jl: A Julia basket of hand-picked Krylov methods. *Journal of Open Source Software*, *8*(89), 5187. https://doi.org/10.21105/joss.05187

```julia
P = ic02(A_gpu)

# Additional vector required for solving triangular systems
z = CUDA.zeros(ComplexF64, n)

# Solve Py = x
function ldiv_ic0!(P, x, y, z)
  ldiv!(z, LowerTriangular(P), x)   # Forward substitution with L
  ldiv!(y, LowerTriangular(P)', z)  # Backward substitution with Lᴴ
  return y
end

# Linear operator that approximates the preconditioner P⁻¹ in floating-point arithmetic
T = ComplexF64
symmetric = false
hermitian = true
P⁻¹ = LinearOperator(T, m, n, symmetric, hermitian, (y, x) -> ldiv_ic0!(P, x, y, z))

# Solve an Hermitian positive definite system with an incomplete Cholesky factorization preconditioner
x_gpu, stats = cg(A_gpu, b_gpu, M=P⁻¹)

# Check the solution returned by the conjugate gradient method
x_cpu = Vector{ComplexF64}(x_gpu)
@test norm(x_cpu - x_exact) ≤ 1e-5
end
```

## Acknowledgements

## References

Balay, S., Abhyankar, S., Adams, M. F., Benson, S., Brown, J., Brune, P., Buschelman, K., Constantinescu, E., Dalcin, L., Dener, A., Eijkhout, V., Faibussowitsch, J., Gropp, W. D., Hapla, V., Isaac, T., Jolivet, P., Karpeev, D., Kaushik, D., Knepley, M. G., … Zhang, J. (2023). *PETSc/TAO users manual* (ANL-21/39 - Revision 3.19). Argonne National Laboratory. https://doi.org/10.2172/1968587

Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. B. (2017). Julia: A fresh approach to numerical computing. *SIAM Review*, *59*(1), 65–98. https://doi.org/10.1137/141000671

Davis, T., & Hu, Y. (2011). The University of Florida sparse matrix collection. *ACM Transactions on Mathematical Software*, *38*(1), 1–25. https://doi.org/10.1145/2049662.2049663

Greenbaum, A. (1997). *Iterative methods for solving linear systems*. SIAM. https://doi.org/10.1137/1.9781611970937

Ipsen, I. C., & Meyer, C. D. (1998). The idea behind Krylov methods. *The American Mathematical Monthly*, *105*(10), 889–899. https://doi.org/10.1080/00029890.1998.12004985

Krylov, A. N. (1931). On the numerical solution of the equation by which, in technical matters, frequencies of small oscillations of material systems are determined. *Izvestija AN SSSR (News of Academy of Sciences of the USSR), Otdel. Mat. I Estest. Nauk*, *7*(4), 491–539.

MATLAB. (2022). *Version 9.13.0 (R2022b)*. The MathWorks Inc.

Saad, Y. (2003). *Iterative methods for sparse linear systems*. SIAM. https://doi.org/10.1137/1.9780898718003