





WHY JULIA

- 1959 Fortran and ALGOL: **Mathematical algorithms** on computers
- 1969 C: Flexibility and performance for operating systems
- 1980 C++: Structure support for large code bases
- 1991 Python: **Productivity and ease of us**
- ?: Productivity, Flexibility/Portability, and Performance



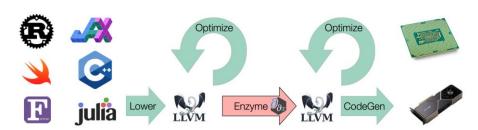




ANOTHER TRY



- Developed by Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral Shah
- Released in 2012
- Designed for numerical and scientific computing



1. Performance

Based on the LLVM compiler

2. Flexibility

- LLVM IR or similar languages are used for a wide variety of hardware and architectures (x86, CUDA, ROCm, SPIRV...)
- Just-in-time compilation, only code that runs is compiled
- Metaprogramming and reflections (macros)

3. Productivity

- Ease of use coupled with integrated continuous integration, package management, and version management (Manifest)
- Interactive





PRODUCTIVITY

Code development and maintenance cost far exceed research.

Interior Point Method Optimization Solver for Nonlinear Optimization

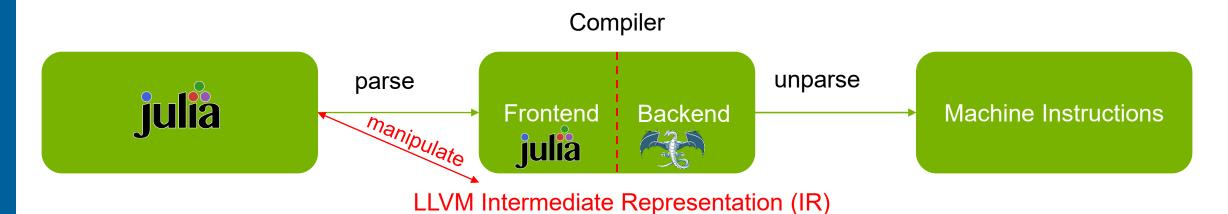
Solver	Language	Lines of Code
Ipopt	3	100,000
MadNLP	julia	6,000

MadNLP has GPU APIs





JULIA AND LLVM



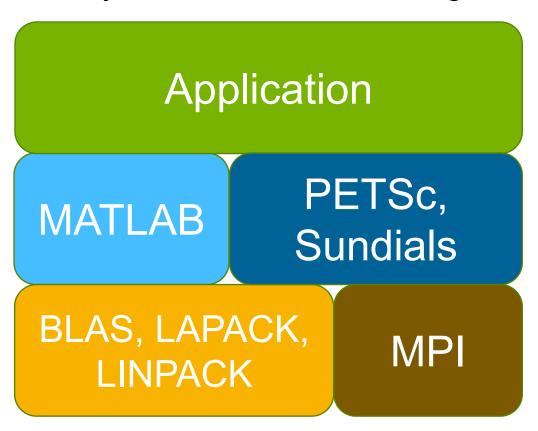
- Julia leverages LLVM to implement a just-in-time compiled language with native metaprogramming and code reflection. Code is compiled at runtime.
- Language support for IR and expression tree manipulation
- Advantage: Highly flexible (JIT) C/C++-like performance (LLVM backend)
 - Ada, C, C++, D, Delphi, Fortran, Haskell, Julia, Objective-C, Rust, and Swift
- Disadvantage: LLVM was not created with JIT in mind
- See interpreted languages: Trying to reduce compilation (Python)





JULIA AND HPC FOUNDATIONS

Early HPC stack for linear algebra

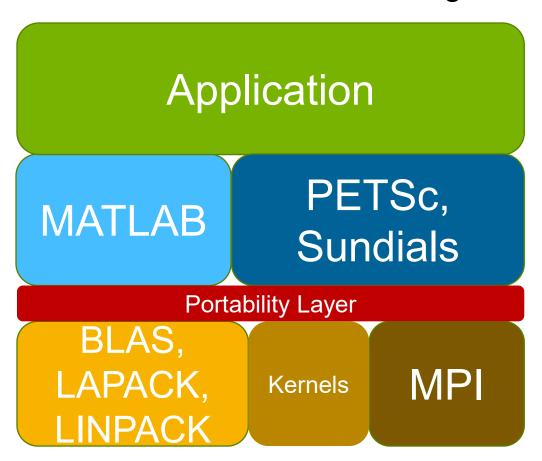


- Abstraction for high-performance linear algebra (BLAS, LAPACK)
- Abstraction for distributed computing (MPI)
- Linear algebra for distributed linear algebra (e.g., PETSc)
- Abstraction for high-level algorithms



JULIA AND HPC FOUNDATIONS

GPU HPC stack for linear algebra



- Abstraction for high-performance linear algebra (BLAS, LAPACK)
- Abstraction for distributed computing (MPI)
- Linear algebra for distributed linear algebra (e.g., PETSc)
- Abstraction for high-level algorithms
- Kernel programming
- Portability Layer
- All in one language with no compilation hell



JULIA AND HPC SOFTWARE

Interfaces to C and Python are straightforward

```
using Libdl
const libc path = Libdl.find library(["libc.so.6", "libc"])
println("Found libc at: $libc path")
const libc handle = Libdl.dlopen(libc path)
const puts fptr = Libdl.dlsym(libc handle, :puts)
function julia puts(s::String)
   ccall(
      Cint,
      (Cstring,), # It takes one argument, a C string (char*)
       "Hello, C from Julia!" # The string we want to print
end
julia_puts("This is from the Julia wrapper function.")
```

```
"libc.so.6"
Found libc at: libc.so.6
Ptr{Nothing}(0x00007fa0c173f320)
Ptr{Nothing}(0x00007fa0c1551e50)
julia puts (generic function with 1 method)
Hello, C from Julia!
21
```



JULIA AND HPC SOFTWARE

```
michel@intel-G501 ~/git/julia_alcf/aurora/example <main >>
└$ julia --project=.
                          Documentation: https://docs.julialang.org
                          Type "?" for help, "]?" for Pkg help.
                          Version 1.12.1 (2025-10-17)
                          Official https://julialang.org release
(example) pkg> add oneAPI
   Resolving package versions...
julia> using oneAPI
julia> oneAPI.versioninfo()
Binary dependencies:
 NEO: 25.35.35096+0
 libigc: 2.18.5+0
  gmmlib: 22.8.1+0
 SPIRV LLVM Translator: 21.1.1+0
 SPIRV Tools: 2025.4.0+0
 oneAPI_Support: 0.9.2+0 (oneMKL v2025.2.0)
Toolchain:
 Julia: 1.12.1
 LLVM: 18.1.7
1 driver:
 00000000-0000-0000-1874-85ae01038918 (v1.3.35096, API v1.6.0)
 device:
 Intel(R) Arc(TM) A750 Graphics
```

- Julia provides an intricate crosscompilation system for providing binary artifacts linking only against local libc
- oneAPI Intel dependencies are provided precompiled for any Linux system











PORTABILITY IN JULIA

Layers	Julia Code				
Abstraction	GPUArrays.jl + KernelAbstractions.jl				
Generation	CUDA.jl	AMDGPU.jl	oneAPI.jl	Host/CPU	
IR / C API	CUDA	ROCm	Intel Compute Runtime	LLVM	

- All vendor languages support LLVM IR or SPIRV intermedia representations (IR)
- Only high-level language with support for GPUs
- BLAS, LAPACK etc is dispatched to vendor libraries based on type
- Native performance using vendor toolchain
- However: simplified API for synchronization
 - Julia uses one CUDA stream. Kernel executions are ordered and synchronous
 - If you really want very high-performance -> go down to C/C++
 - Trade-off between user friendliness and performance





PORTABILITY IN JULIA: KERNELABSTRACTIONS

- Architecture abstraction through backends
- Generic DSL for SIMD kernels using Julia macros (e.g., @index, @kernel)
- Uses GPUCompiler.jl
 infrastructure to transform lowered
 Julia LLVM IR for each backend

```
using oneAPI #, CUDA, AMDGPU
     using KernelAbstractions
24
     using Adapt
     backend = oneAPIBackend()
26
     a = adapt(backend, ones(10))
28
     b = adapt(backend, ones(10))
     c = similar(a)
30
     # KernelAbstractions.jl kernel definition
     @kernel function myadd!(c, a, b)
33
         i = @index(Global)
34
         c[i] = a[i] + b[i]
35
     end
     # Instantiate kernel
     kernel! = myadd!(backend)
     # Launch the kernel
     kernel!(c, a, b; ndrange=length(c))
```

PORTABILITY IN JULIA: GPUARRAYS

- Supports the broadcast operator '.'
- Supplies common methods for GPU arrays using KernelAbstractions.jl (e.g., sort!, accumulate)

```
using oneAPI #, CUDA, AMDGPU
     using Adapt
     using LinearAlgebra
     # Other backends: CUDABackend(), AMDGPUBackend()
     backend = oneAPIBackend()
     a = adapt(backend, ones(10))
     b = adapt(backend, ones(10))
     c = similar(a)
11
     # Generates a kernel using the broadcast operator
     c = a + b
     # Sorts array c in place
     sort!(c)
16
     A = adapt(backend, rand(10,10))
     x = adapt(backend, rand(10))
     # Matrix-vector multiplication
     y = mul!(similar(x), A, x)
```



PORTABILITY IN JULIA: MPI

- No complex derived type handling
- Support of GPU-aware MPI
- Even functions can be sent

```
using MPI
     using oneAPI
     MPI.Init()
     comm = MPI.COMM WORLD
     rank = MPI.Comm_rank(comm)
     size = MPI.Comm_size(comm)
     println("Hello from rank $rank of $size")
51
     struct DerivedType{VT <: AbstractVector}</pre>
         v::VT
52
53
         a::Int
     end
     a = DerivedType(adapt(backend, ones(10)), 42)
     # Send and receive the derived type round-robin
     req = MPI.isend(a, comm; dest=(rank + 1) % size)
     b = MPI.recv(comm; source=(rank - 1 + size) % size)
     MPI.Wait(req)
62
     MPI.Finalize()
```



PORTABILITY IN JULIA: LIBRARIES

Example: Krylov.jl

- A portable Krylov solver library with support for NVIDIA, AMD, and Intel GPUs
- Comes with all Krylov algorithms under the sun
- Transparent interfaces to preconditioners through KrylovPreconditioners.jl
- Works with custom operators (e.g., written in KernelAbstractions.jl)
- Example: Minimizing a function f, @allowscalar 🔒

```
using ForwardDiff, Krylov, CUDA
using GPUArrays

xk = -ones(4)
f(x) = (x[1] - 1)^2 + (x[2] - 2)^2 + (x[3] - 3)^2 + (x[4] - 4)^2
g(x) = ForwardDiff.gradient(f, x)

H(x) = ForwardDiff.hessian(f, x)

allowscalar d, stats = cg(H(xk), -g(xk))

xk += d # [1,2,3,4]
```





WHY SO FAR ONLY IN JULIA?

- Compute gradient of f(x) = exp(dot(x,x))
- CUDA, linear algebra, and forward automatic differentiation are provided by different packages (separation of concerns) and are composed at runtime (just-in-time compilation).
- The code for the gradient of f(x) in CUDA is created and only exists at runtime!

```
using CUDA
using LinearAlgebra
using ForwardDiff

f(x) = exp(dot(x,x))

x = Array([2.0,3.0])

f(x)

@code_llvm(f(x)) # Look at the LLVM IR
f(CuArray(x))

@show g = ForwardDiff.gradient(f, CuArray(x))
```

 Compact code with three off the shelf packages from the official Julia package registry that leverage 1-4 on the left



NON-PORTABLE VENDOR LIBRARIES

Example: Sparse Direct Linear Solver CUDSS

- CUDSS is the state-of-the-art sparse direct linear solver for GPUs
- No other GPU vendor has a performant sparse direct solver
- Julia cannot avoid vendor specifics, but tries to unify API as much as possible
- CUDSS Julia interface is maintained by Alexis Montoison (MCS)
- Provides a generic Julia interface for linear solvers/factorizations
- For best performance CUDSS API





ARCHITECTURE OVERVIEW

PAST

Many core with GPUs

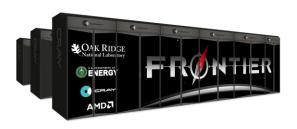






CURRENT

Back to vector processors (GPUs)







Future

NVIDIA and Oracle to Build US Department of **Energy's Largest Al Supercomputer for Scientific Discovery**

Bold US Investment of 100,000 NVIDIA Blackwell GPUs Kickstarts Era of Agentic Al-Powered Science at Argonne National Laboratory for Public Researchers

FPGAs, wafer computer (LLVM based compilers)









JULIA AT ALCF: GOAL

Goal

- module load Julia
- Support for GPUs, MPI, filesystems, and large-scale submissions

Julia-specific challenges

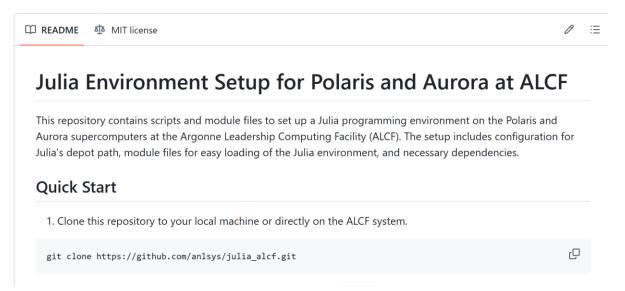
- Compiled at runtime requires file access by every process
- Offline compilation of binaries is quirky (Support Julia 1.12)
- Julia's package manager
- Artifacts and state in Manifest.toml





CURRENT STATUS: EXPERIMENTAL

- Part of official documentation: https://docs.alcf.anl.gov/polaris/programming-models/julia/
- No system deployment
- Module files and setup are found in https://github.com/anlsys/julia_alcf







SETUP

- Do not install in your home directory!
- git clone https://github.com/anlsys/julia_alcf.git
- cd julia_alcf/Polaris or cd julia_alcf/aurora
- ./setup.sh
- It will prompt for a JULIA_DEPOT_PATH. Select a folder on a fast filesystem and accessible by the compute nodes
- It will install the latest Julia (1.12.1)
- It will install module files and a LocalPreferences.toml into JULIA_DEPOT_PATH





LOAD JULIA

- module use \$JULIA_DEPOT_PATH/modulefiles && module load julia
- Add one line to your shell initialization scripts (e.g., ~/.bashrc)
 - export JULIA_DEPOT_PATH=YOUR_DEPOT_PATH
- All done!



USING JULIA ON AN ALCF CLUSTER

- Clone your project with the Project.toml and do the usual setup on the login nodes
 - julia -project=.
 -] update
- If your code uses CUDA.jl you cannot run it using GPUs on the login nodes!
- Using MPI would precompile your code on every process (potential call from filesystem team)
- Run the code on 1 node to precompile as much as possible
 - qsub -I -l select=1,walltime=1:00:00,filesystems=home:eagle -A
 [PROJECT] -q debug
 - julia -project simulation.jl
- Extreme large-scale runs (not recommended right now, contact us)
 - Compress your JULIA_DEPOT_PATH and move it to the local SSDs with the submission script





FUTURE

- We will support Intel GPUs the best we can through maintenance of oneAPI.jl
 - <u>https://github.com/JuliaGPU/oneAPI.jl</u>
- Support for CUDA is far more mature
- Post-Aurora system will be NVIDIA based
 - Great support more unified memory and other CUDA features
 - Support for NVIDIA libraries
- Rollout of a system-wide module in the coming months





Thank You





TALK STRUCTURE

- Why Julia
- What is Julia?
 - Multiple dispatch, JIT in LLVM, functional programming
- Julia and HPC Foundations
 - BLAS
 - GPUCompiler.jl, CUDA.jl, oneAPI.jl, KernelAbstractions.jl
 - MPI
 - Sparse Linear Solver
- Julia on Supercomputers
 - Challenges
 - Artifacts, local libraries, (pre)compilation, filesystems, and ABIs
 - Current Solution
 - Julia and system depot_path are provided on local SSDs, user depot_path self-managed on project filesystem
 - Best practices
 - How to submit job and precompile
- Status and Outlook
 - Module load Julia: Adds Julia to path, provides a system depot_path with precompiled oneAPI.jl/CUDA.jl, MPI.jl, and HDF5.jl
 - Agenda with rollout, links to current setup on Github
 - Disclaimer: Experimental, we need feedback



