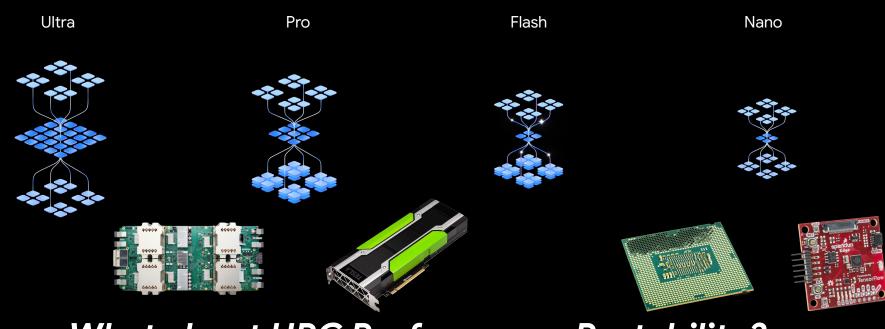
Google DeepMind

Condenser: noun, an Apparatus for Compiling Science to the Cloud

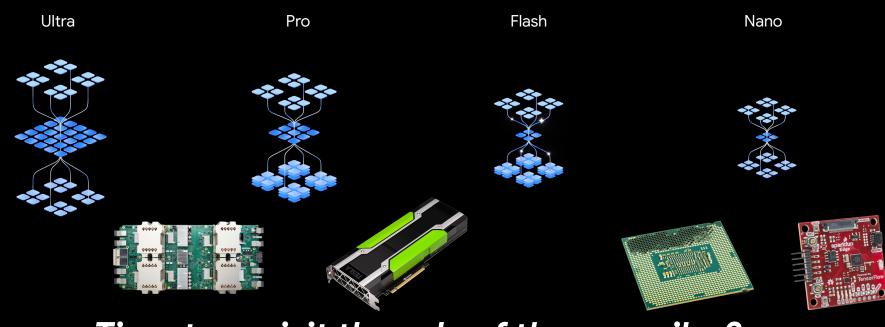


Automatic Parallelization, Performance Portability: Solved Problems for Al Applications?



What about HPC Performance Portability?

Automatic Parallelization, Performance Portability: Solved Problems for Al Applications?



Time to revisit the role of the compiler?

Before the "condenser" we started making waves in the cloud

Platforms

- NERSC Perlmutter
 1536 GPU accelerated nodes with
 1 AMD Milan processor and 4
 NVIDIA A100 GPUs
- Google Cloud reservation 1,679 TPU v6e (Trillium)
 1.5 ExaFLOPS (bf16)
 53 TB of HBM
 3.2 TB/s bisection bandwidth

Making Waves in the Cloud: A Paradigm Shift for Scientific Computing and Ocean Modeling through Compiler Technology

William S. Moses^{†§}, Mosè Giordano*, Avik Pal[‡], Gregory Wagner[‡], Ivan R Ivanov, Paul Berg[▽],
Johannes Blaschke, Jules Merckx[△], Arpit Jaiswal[‡], Patrick Heimbach[‡], Son Vu, Sergio
Sanchez-Ramirez, Simone Silvestri, Nora Loose[‡], Ivan Ho, Vimarsh Sathia[†], Jan Hueckelheim[‡],
Johannes De Fine Licht, Kevin Gleason[§], Ludovic Rass, Gabriel Baraldi, Dhruv Apte[‡], Lorenzo
Chelini[‡], Jacques Pienaar[§], Gaetan Lounes, Valentin Churavy, Sri Hari Krishna Narayanan[‡], Navid
Constantinou, William R. Magro[§], Michel Schanen[‡], Alexis Montoison[‡], Alan Edelman[‡], Samarth
Narang, Tobias Grosser, Keno Fischer[‡], Robert Hundt[§], Albert Cohen[§], Oleksandr Zinenko[§] *
UIUC [†], Google [§], UCL ^{*}, MIT [‡], NVIDIA [‡], University of Mainz [‡], BFH [▽], Ghent University [△]

Abstract

Ocean and climate models are today limited by compute resources, forcing approximations driven by feasibility rather than theory. They consequently miss important physical processes and decisionrelevant regional details. Advances in AI-driven supercomputing specialized tensor accelerators, AI compiler stacks, and novel distributed systems - offer unprecedented computational power. Yet, scientific applications such as ocean models, often written in Fortran, C++, or Julia and built for traditional HPC, remain largely incompatible with these technologies. This gap hampers performance portability and isolates scientific computing from rapid cloud-based innovation for AI workloads. In this work, we bridge that gap by transpiling a Julia-based ocean model (Oceananigans) using the MLIR compiler infrastructure. This process enables advanced optimizations, deployment on AI hardware (e.g., Google TPUs), and automatic differentiation. Our results demonstrate hardware and software designed accelerate climate to benefit from

2 Justif

Autom model, in Julia w specific d entry for Gordon Bell prize nomination

today's larg

Author's Contact Information: Wi Gregory Wagner[‡], Ivan R Ivanov, Arpii Jaiswal[‡], Patrick Heimbach[‡], U., Sergio Sanchez-Ramirez, Simone Silvestri, Nora Loose[‡], Ivan Ho, Vimarsh Sathia[‡], Jan Hueckelheim[‡], Johannes De Fine Licht, Kevin Gleason[‡], Ludovic Rass, Gabriel Baraldi, Dhruv Apte[‡], Lorenzo Chelim[‡],

3 Performance Attributes

Performance Attribute	This Submission
Category achievement	scalability
Type of method used	semi-implicit
Results reported on the basis of	whole application except I/O
Precision reported	double precision (GPU)
	emulated double precision (TPU)
System scale	results measured on full-scale system
Measurement mechanism	timers, FLOP count

4 Overview of the Problem

Climate is governed by planetary fluid dynamics. This submission focuses on the core of global climate models: simulations of the id dynamics of the ocean and atmosphere that dictate the large-tructure and long-term evolution of Earth's climate. Fluid all processes underpin phenomena like equator-to-pole port, ENSO, the jet stream, air-sea interaction, and the line circulation, all of which drive variability and set is memory and predictability [13, 19, 36]. Accurately llimate requires ocean and atmospheric dynamical cores governing equations of fluid motion as efficiently and as possible.

need for high resolution. Influential climate processes cover a jide range of interacting spatial scales [18], from planetary (10,000 km), synoptic (1,000 km), tropical cyclones and ocean geostrophic eddies (10 – 200 km), atmospheric mesoscale convective systems and ocean submesoscale processes (1 – 10 km), internal gravity waves, clouds, turbulent diffusion, convective mixing on scales (10 m), and down to the dissipation of kinetic energy (1 mm). Because current global ocean and atmosphere models cannot fully resolve all these small-scale processes, many processes are represented using simplified approximations called parameterizations. However,

Application: Oceananigans.jl

https://clima.github.io/OceananigansDocumentation

Simple simulation of **baroclinic instability** on an Earth-like planet: essential features of ocean and atmosphere interactions

Multiple integrals and solvers:

implicit vertical diffusion

hydrostatic pressure anomaly

vertical velocity

horizontal velocities

5th-order WENO-based advection schemes

tracers suitable for ultra-high-resolution

55-term polynomial approximation to the TEOS10 equation of state for density as a function of oceanic temperature, salinity, and pressure

Making Waves in the Cloud: A Paradigm Shift for Scientific Computing and Ocean Modeling through Compiler Technology

William S. Moses^{†§}, Mosè Giordano*, Avik Pal[‡], Gregory Wagner[‡], Ivan R Ivanov, Paul Berg[▽], Johannes Blaschke, Jules Merckx[△], Arpit Jaiswal[‡], Patrick Heimbach[#], Son Vu, Sergio Sanchez-Ramirez, Simone Silvestri, Nora Loose[‡], Ivan Ho, Vimarsh Sathia[†], Jan Hueckelheim[‡], Johannes De Fine Licht, Kevin Gleason[§], Ludovic Rass, Gabriel Baraldi, Dhruv Apte[#], Lorenzo Chelini[‡], Jacques Pienaar[§], Gaetan Lounes, Valentin Churavy, Sri Hari Krishna Narayanan[‡], Navid Constantinou, William R. Magro[§], Michel Schanen[‡], Albert Sondois Montoison[‡], Alan Edelman[‡], Samarth Narang, Tobias Grosser, Keno Fischer[‡], Robert Hundt[§], Albert Cohen[§], Oleksandr Zinenko^{§ *} UIUC [†], Google [§], UCL ^{*}, MIT [‡], NVIDIA [‡], UT Austin [#], [C]Worthy [‡], BSC [§], Argonne National Laboratory [‡], LBNL [©], Cambridge [§], JuliaHub [§], University of Mainz [#], BFH [©], Ghent University [△]

Abstract

Ocean and climate models are today limited by compute resources, forcing approximations driven by feasibility rather than theory. They consequently miss important physical processes and decisionrelevant regional details. Advances in AI-driven supercomputing specialized tensor accelerators, AI compiler stacks, and novel distributed systems - offer unprecedented computational power. Yet, scientific applications such as ocean models, often written in Fortran, C++, or Julia and built for traditional HPC, remain largely incompatible with these technologies. This gap hampers performance portability and isolates scientific computing from rapid cloud-based innovation for AI workloads. In this work, we bridge that gap by transpiling a Julia-based ocean model (Oceananigans) using the MLIR compiler infrastructure. This process enables advanced optimizations, deployment on AI hardware (e.g., Google TPUs), and automatic differentiation. Our results demonstrate hardware and software designed accelerate climate to benefit from

2 Justifi

Autom model, in Julia w specific d entry for Gordon Bell prize nomination

*Correspondence: wsmo

Author's Contact Information: Wi Gregory Wagner[†], Ivan R Ivanov, Arpit Jaiswal[†], Patrick Heimbach[†] Vu, Sergio Sanchez-Ramirez, Simone Silvestri, Nora Loose[†], Ivan Ho, Vimarsh Sathia[†], Jan Hueckelheim[†], Johannes De Fine Licht, Kevin Gleason[†], Ludovic Rass, Gabriel Baraldi, Dhruv Apte[†], Lorenzo Chelini[†],

3 Performance Attributes

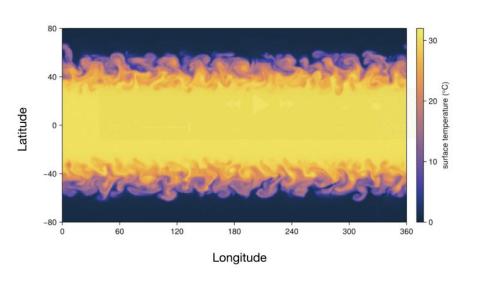
Performance Attribute	This Submission
Category achievement	scalability
Type of method used	semi-implicit
Results reported on the basis of	whole application except I/O
Precision reported	double precision (GPU)
	emulated double precision (TPU)
System scale	results measured on full-scale system
Measurement mechanism	timers, FLOP count

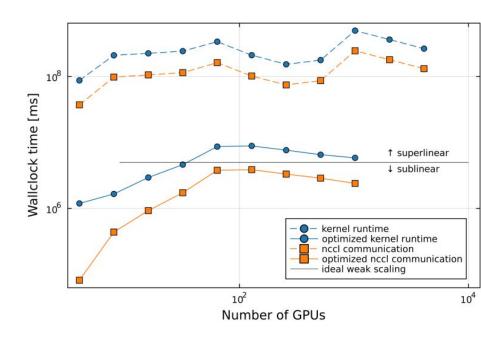
4 Overview of the Problem

Climate is governed by planetary fluid dynamics. This submission focuses on the core of global climate models: simulations of the id dynamics of the ocean and atmosphere that dictate the large-structure and long-term evolution of Earth's climate. Fluid all processes underpin phenomena like equator-to-pole port, ENSO, the jet stream, air-sea interaction, and the line circulation, all of which drive variability and set e's memory and predictability [13, 19, 36]. Accurately llimate requires ocean and atmospheric dynamical cores governing equations of fluid motion as efficiently and as possible.

need for high resolution. Influential climate processes cover a jide range of interacting spatial scales [18], from planetary (10,000 km), synoptic (1,000 km), tropical cyclones and ocean geostrophic eddies (10 – 200 km), atmospheric mesoscale convective systems and ocean submesoscale processes (1 – 10 km), internal gravity waves, clouds, turbulent diffusion, convective mixing on scales (10 m), and down to the dissipation of kinetic energy (1 mm). Because current global ocean and atmosphere models cannot fully resolve all these small-scale processes, many processes are represented using simplified approximations called parameterizations. However,

Weak Scaling Experiments: GPU / Placement and Collectives





Weak Scaling Experiments: GPU / Kernels and Host-Device Communication

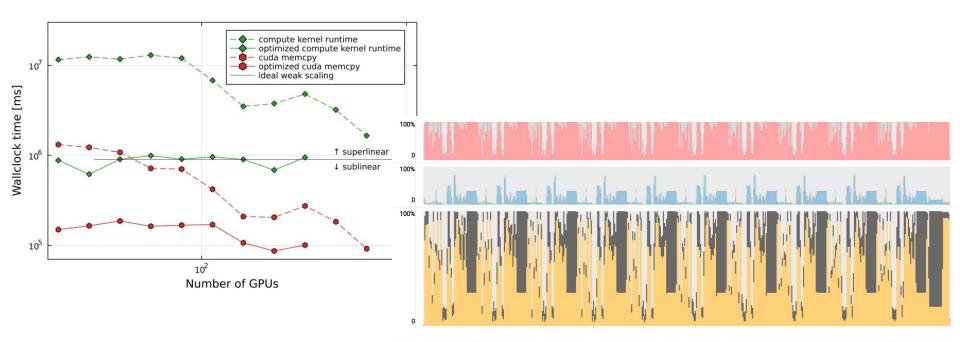
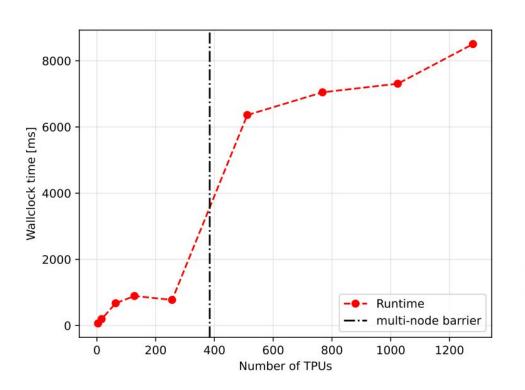


Figure 7: GPU utilization during 10 iterations of the benchmark problem on one Perlmutter GPU node (4 Nvidia A100s). Data collected with Nsight Systems. *top*: Percentage of SMs active, *middle*: Percentage saturation of active SMs, *bottom*: SM Warp Occupancy (yellow shows the active warps, gray shows inactive warps). During these 10 iterations on average 50.6% of compute warps were utilized, while an average of 26.6% compute warps were unallocated on active SMs.

Weak Scaling Experiments: TPU



Operation	Percent of Execution
Concatenate	39.04%
Reduce-Window	35.01%
Loop-Fusion 1	19.71%
Data Formatting	2.89%
Slice	1.59%
X64Combine	0.88%
Collective-Permute	0.48%

Table 1: Breakdown of TPU execution time by operation type, on a single node 4-TPU machine.

Q1: What missing interoperability layers (software, standard, or abstraction) would most accelerate convergence between traditional HPC linear algebra workflows and today's extreme-scale Al workloads?

Let's take a look.

How did we get there

```
import jax.numpy as np
from jax import jit, grad, vmap
def predict(params, inputs):
  for W, b in params:
    outputs = np.dot(inputs, W) + b
    inputs = np.tanh(outputs)
  return outputs
def loss(params, batch):
  inputs, targets = batch
  preds = predict(params, inputs)
  return np.sum((preds - targets) ** 2)
```



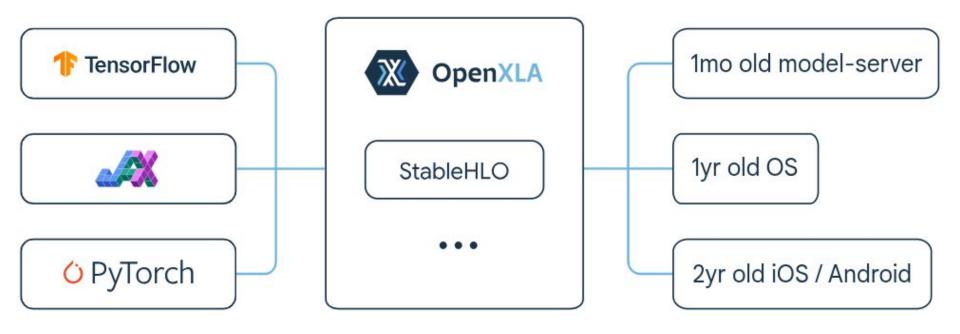
```
gradient_fun = jit(grad(loss))
perexample_grads = jit(vmap(grad(loss), (None, 0)))
```

JAX is an extensible system for composable function transformations of Python+NumPy code, with computations staged to XLA

StableHLO

https://openxla.org/stablehlo





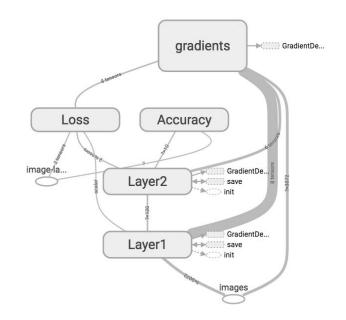
XLA / HLO / StableHLO



Folklore: Go Domain-Specific!

- →Domain-Specific Languages (DSLs) expose the right **abstractions** to make **automatic code generation** possible and effective
- → Also domain-specific accelerators

HLO: XLA Compute Graph, Static Dataflow



XLA / HLO / StableHLO



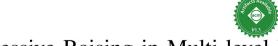
Folklore: Go Domain-Specific!

- →Domain-Specific Languages (DSLs) expose the right **abstractions** to make **automatic code generation** possible and effective
- → Also domain-specific accelerators

- JIT / AOT compiler for linear algebra
- Multiple backends: CPUs, GPUs, TPUs
- Eliminate dispatch overhead
- Fuse operations: avoid round trips to memory
- Specialization, global buffer analysis, vectorization, unrolling, etc.

What about HPC?





Progressive Raising in Multi-level IR

Lorenzo Chelini TU Eindhoven Eindhoven, The Netherlands l chelini@tue nl

> Nicolas Vasilache Google Zurich, Switzerland ntv@google.com

Andi Drebes Inria and École Normale Supérieure Paris, France andi@programmierforen.de

> Tobias Grosser University of Edinburgh Edinburgh, UK tobias.grosser@ed.ac.uk

Oleksandr Zinenko Google Paris, France

Albert Cohen Google Paris, France zinenko@google.com albertcohen@google.com

> Henk Corporaal TU Eindhoven Eindhoven, The Netherlands h.corporaal@tue.nl

Abstract—Multi-level intermediate representations (IR) show great promise for lowering the design costs for domain-specific compilers by providing a reusable, extensible, and non-opinionated framework for expressing domain-specific and high-level abstractions directly in the IR. But, while such frameworks support the progressive lowering of high-level representations to low-level IR, they do not raise in the opposite direction. Thus, the entry point into the compilation pipeline defines the highest level of abstraction for all subsequent transformations, limiting the set of applicable optimizations, in particular for general-purpose languages that are not semantically rich enough to model the required abstractions.

We propose Progressive Raising, a complementary approach to the progressive lowering in multi-level IRs that raises from lower to higher-level abstractions to leverage domain-specific transformations for low-level representations. We further introduce Multilevel Tactics, our declarative approach for progressive raising, implemented on top of the MLIR framework, and demonstrate the progressive raising from affine loop nests specified in a

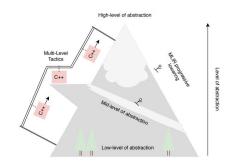


Fig. 1. Multi-level Tactics lifts general-purpose languages to higher-abstraction levels to enable effective domain-specific compilation via progressive lowering.

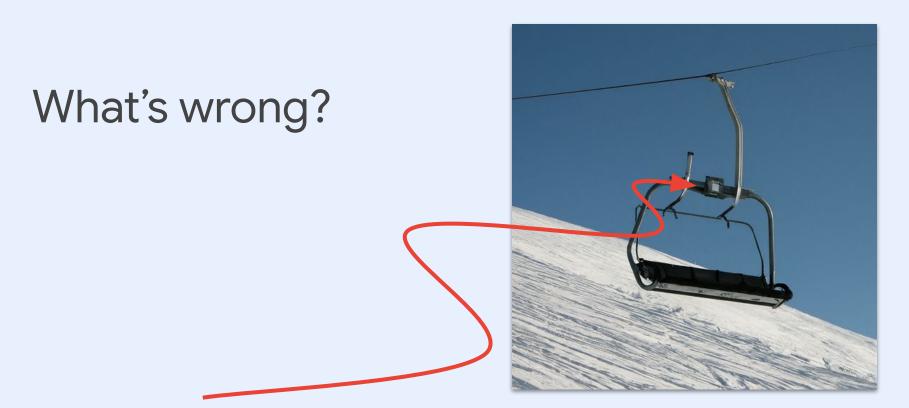
→Domain-Specific Languages (DSLs) expose the right abstractions to make automatic code generation possible and effective

Folklore: Go Domain-Specific!

→ Raising the level of abstraction is harder than riding the abstraction lowering slope

What's wrong?





Lifting, isn't this just as hard as automatic parallelization?
What about the abstraction penalty? Performance predictability? Control for performance engineers?

HPC Systems are (also) Stuck in a Rut!



Machine Learning Systems are Stuck in a Rut

Paul Barham Google Brain

Abstract

In this paper we argue that systems for numerical computing are stuck in a local basin of performance and programmability. Systems researchers are doing an excellent job improving the performance of 5-year-old benchmarks, but gradually making it harder to explore innovative machine learning research ideas.

We explain how the evolution of hardware accelerators favors compiler back ends that hyper-optimize large monolithic kernels, show how this reliance on high-performance but inflexible kernels reinforces the dominant style of programming model, and argue these programming abstractions lack expressiveness, maintainability, and modularity; all of which hinders research

We conclude by noting promising directions in the field, and advocate steps to advance progress towards high-performance general purpose numerical computing systems on modern accelerators.

ACM Reference Format:

Paul Barham and Michael Isard. 2019. Machine Learning Systems are Stuck in a Rut. In Workshop on Hot Topics in Operating Systems (HotOS '19), May 13–15, 2019, Bertinoro, Italy. ACM, New York, NY, USA, 7 pages. https://doi.org/10.1145/3317550. 3321441

1 Compiling for modern accelerators

We became interested in this paper's subject when trying to improve an implementation of Capsule networks [1] to scale up to larger datasets. Capsule networks are an exciting machine learning research idea where scalar-valued "neurons" are replaced by small matrices, allowing them to capture more complex relationships. Capsules may or may not be the "next big thing" in machine learning, but they serve as a representative example of a disruptive ML research idea.

Although our convolutional Capsule model requires around 4 times fewer floating point operations (FLOPS)

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

HotOS '19, May 13–15, 2019, Bertinoro, Italy © 2019 Copyright held by the owner/author(s) ACM ISBN 978-1-4503-6727-1/19/05. https://doi.org/10.1145/3317550.3321441 Michael Isard Google Brain

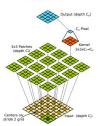


Figure 1. Conv2D operation with 3×3 kernel, stride=2

with 16 times fewer training parameters than the convolutional neural network (CNN) we were comparing it to, implementations in both TensorFlow[2] and PyTorch[3] were much slower and ran out of memory with much smaller models. We wanted to understand why.

1.1 New ideas often require new primitives

We won't discuss the full details of Capsule networks in this paper¹, but for our purposes it is sufficient to consider a simplified form of the inner loop, which is similar to the computation in a traditional CNN layer but operating on 4×4 matrices rather than scalars.

A basic building block of current machine learning frameworks is the strided 2D convolution. Most frameworks provide a primitive operation that accepts N input images of size HxW, where each pixel has a "depth" of C_i channels. Informally, for a 'kerrel size' $K^2 = 3$ and "stride' S=2, conv2d computes a weighted sum of overlapping 3x3 patches of pixels centered at every other (x,y) coordinate, to produce N smaller images with pixel depth C_0 (Figure 1). Mathematically, this can be expressed as follows:

$$\forall n, x, y, c_o: O_{x,y}^{n,c_o} = \sum_{k_x} \sum_{k_y} \sum_{c_i} I_{sx+k_x, sy+k_y}^{n,c_i} \cdot K_{k_x,k_y}^{c_i,c_o}$$
 (1)

where \cdot denotes scalar multiplication, and O, I, and K are all 4-dimensional arrays of scalars. The resulting code is little more than 7 nested loops around a multiply-accumulate operation, but array layout, vectorization,

¹For an excellent tutorial on Capsule networks see [4].
² Section 4 discusses why these dimensions are used.

https://commons.wikimedia.org/wiki/File:Stuck in a Rut.jpg

Kernel Programming to the Rescue

Flurry of GPU acceleration options

CUDA Kernels / OpenCL-C

SYCL

CUTLASS

Triton (PyTorch, JAX)

Pallas (JAX)

Turbine (AMD)

Mojo (Modular)

CuTile (Nvidia)

and more coming and going...

More broadly

"If high level fails, try lower level" Folklore: high-level language



abstraction penalty

Motivations: escape hatch for...

- Performance tricks
- Extra expressivenesse.g. ragged or sparse tensors
- Quick experiments

Unify Kernel Programming with Compiler Automation!

Why? Hardware Optionality

- Compiler = HW-enabling differentiator
 - → best-effort
- Kernel languages = siloed HW / SW stacks
 - → all or nothing



Journal of Parallel and Distributed Computing



Volume 61, Issue 12, December 2001, Pages 1803-1826

Regular Article

Telescoping Languages: A Strategy for Automatic Generation of Scientific Problem-Solving Systems from Annotated Libraries

<u>Ken Kennedy</u> ^a, <u>Bradley Broom</u> ^a, <u>Keith Cooper</u> ^a, <u>Jack Dongarra</u> ^b, <u>Rob Fowler</u> ^a, <u>Dennis Gannon</u> ^c, Lennart Johnsson ^d, John Mellor-Crummey ^a, Linda Torczon ^a

Staging and Partial Evaluation

Domain-specific generators

Multi-stage programming

(macro-expansion, quasi-quotation)

Late binding of kernel implementations Optional cross-stage persistence

In search of a program generator to implement generic transformations for high-performance computing

Albert Cohen^{a,*}, Sébastien Donadio^b, Maria-Jesus Garzaran^c, Christoph Herrmann^d,

Oleg Kiselyov^e, David Padua^c

^a ALCHEMY group, INRIA Futurs, Orsay, France
^b PRiSM, University of Versailles, France
^c DCS, University of Illinois at Urbana-Champaign, IL, USA
^d FMI, University of Passau, Germany
^e FNMOC, Monterey, CA, USA

Received 1 January 2005; received in revised form 1 June 2005; accepted 1 October 2005

Abstract

The quality of compiler-optimized code for high-performance applications is far behind what optimization and domain experts can achieve by hand. Although it may seem surprising at first glance, the performance gap has been widening over time, due to the tremendous complexity increase in microprocessor and memory architectures, and to the rising level of abstraction of popular programming languages and styles. This paper explores in-between solutions, neither fully automatic nor fully manual ways to adapt a computationally intensive application to the target architecture. By mimicking complex sequences of transformations useful to optimize real codes, we show that generative programming is a practical means to implement architecture-aware optimizations for high-performance applications.

This work explores the promises of generative programming languages and techniques for the high-performance computing expert. We show that complex, architecture-specific optimizations *can* be implemented in a type-safe, purely generative framework. Peak performance is achievable through the careful combination of a high-level, multi-stage evaluation language – MetaOCaml – with low-level code generation techniques. Nevertheless, our results also show that generative approaches for high-performance computing do not come without technical caveats and implementation barriers concerning productivity and reuse. We describe these difficulties and identify ways to hide or overcome them, from abstract syntaxes to heterogeneous generators of code generators, combining high-level and type-safe multi-stage programming with a back-end generator of imperative code.

(© 2006 Elsevier B.V. All rights reserved.

User-Schedulable Languages

```
Input: Algorithm
blurx(x,y) = in(x-1,y)
              + in(x, y)
               + in(x+1,y)
out(x,y) = blurx(x,y-1)
            + blurx(x,v)
            + blurx(x,y+1)
Input: Schedule
blurx: split x by 4 \rightarrow x_0, x_1
        vectorize: x
        store at out.x.
        compute at out.y.
out: split x by 4 \rightarrow x_a, x_b
      split y by 4 \rightarrow y_0, y_1
      reorder: y, x, y, x,
      parallelize: y
      vectorize: x.
```

Halide (Ragan-Kelley et.al. 2013)

Also rewrite systems with semantic guarantees: Lift, Elevate, Rise

```
+ Loop Tiling
yo, xo, ko, yi, xi, ki = s[C].tile(y, x, k, 8, 8, 8)
   for vo in range(128):
     for xo in range(128):
       C[y0*8:y0*8+8][x0*8:x0*8+8] = 0
       for ko in range(128):
         for vi in range(8):
            for xi in range(8):
             for ki in range(8):
               C[yo*8+yi][xo*8+xi] +=
                   A[ko*8+ki][yo*8+yi] * B[ko*8+ki][xo*8+xi]
+ Cache Data on Accelerator Special Buffer
CL = s.cache_write(C, vdla.acc_buffer)
AL = s.cache_read(A, vdla.inp_buffer)
# additional schedule steps omitted ...
+ Map to Accelerator Tensor Instructions
s[CL].tensorize(yi, vdla.gemm8x8)
```

```
tc::IslKernelOptions::makeDefaultM
    .scheduleSpecialize(false)
    .tile({4, 32})
    .mapToThreads({1, 32})
    .mapToBlocks({64, 128})
    .useSharedMemory(true)
    .usePrivateMemory(true)
    .unrollCopyShared(false)
    .unroll(4);
```

TC (Vasilache et.al. 2018)

```
TVM (Chen et.al. 2018)

h: 0..3
w: 0..135
k: 0..1

G: 0..127
Unroll h by 8
Unroll k by 2
Vectorized on k

iter twice along h

G: 0..127
Unroll h by 13
Unroll k by 2
Vectorized on k

TTile (Tollenaere et.al. 2021)
```

```
mm = MatMul(M,N,K)(GL,GL,GL)(Kernel)
               // resulting intermediate specs below
.tile(128,128) // MatMul(128,128,K)(GL,GL,GL)(Kernel)
               // MatMul(128,128,K)(GL,GL,GL)(Block)
  .to(Block)
.load(A, SH, _) // MatMul(128,128,K)(SH,GL,GL)(Block
.load(A, SH, _) // MatMul(128,128,K)(SH,SH,GL)(Block
.tile(64,32) // MatMul(64, 32, K)(SH,SH,GL)(Block
              // MatMul(64, 32, K)(SH,SH,GL)(Warp
 .to(Warp)
.tile(8.8)
               // MatMul(8, 8, K)(SH,SH,GL)(Warp
 .to(Thread) // MatMul(8, 8, K)(SH,SH,GL)(Thread)
.load(A, RF, _) // MatMul(8, 8, K)(RF,SH,GL)(Thread)
.load(B, RF, _) // MatMul(8, 8, K)(RF,RF,GL)(Thread)
.tile(1,1)
               // MatMul(1, 1, K)(RF,RF,GL)(Thread)
.done(dot.cu)
               // invoke codegen, emit dot micro-kernel
```

Fireiron (Hagedorn et.al. 2020)

User-Schedulable Languages... Actually Time-Tested

```
# Avoid spurious versioning
                                               # Peel and shift to enable fusion
addContext(C1L1,'ITMAX>=9')
                                               peel (enclose (C3L1, 2), '3')
addContext(C1L1,'doloop ub>=ITMAX')
                                               peel (enclose (C3L1 2,2),'N-3')
addContext(C1L1,'doloop ub<=ITMAX')
                                               peel (enclose (C3L1 2 1,1),'3')
addContext(C1L1,'N>=500')
                                               peel (enclose (C3L1 2 1 2,1), 'M-3')
addContext(C1L1,'M>=500')
                                               peel (enclose (C1L1,2),'2')
addContext(C1L1,'MNMIN>=500')
                                               peel (enclose (C1L1 2,2), 'N-2')
addContext(C1L1,'MNMIN<=M')
                                               peel (enclose (C1L1 2 1,1),'2')
addContext(C1L1,'MNMIN<=N')
                                               peel(enclose(C1L1 2 1 2,1),'M-2')
addContext(C1L1,'M<=N')
                                               peel (enclose (C2L1, 2), '1')
addContext(C1L1,'M>=N')
                                               peel (enclose (C2L1 2,2),'N-1')
                                               peel (enclose (C2L1 2 1,1),'3')
# Move and shift calc3 backwards
                                               peel(enclose(C2L1 2 1 2,1),'M-3')
shift(enclose(C3L1), {'1','0','0'})
                                               shift(enclose(C1L1 2 1 2 1), ('0', '1', '1'))
shift(enclose(C3L10), {'1','0'})
                                               shift(enclose(C2L1 2 1 2 1), {'0','2','2'})
shift(enclose(C3L11), {'1','0'})
shift(C3L12, {'1'})
                                               # Double fusion of the three nests
shift(C3L13, ('1'))
                                               motion(enclose(C2L1 2 1 2 1), TARGET 2 1 2 1)
                                               motion(enclose(C1L1 2 1 2 1), C2L1 2 1 2 1)
shift(C3L14, {'1'})
shift(C3L15, {'1'})
                                               motion(enclose(C3L1 2 1 2 1), C1L1 2 1 2 1)
shift(C3L16, {'1'})
shift(C3L17, {'1'})
                                               # Register blocking and unrolling (factor 2)
                                               time stripmine (enclose (C3L1 2 1 2 1,2),2,2)
motion (enclose (C3L1), BLOOP)
                                               time stripmine (enclose (C3L1 2 1 2 1,1),4,2)
motion (enclose (C3L10), BLOOP)
motion (enclose (C3L11), BLOOP)
                                               interchange (enclose (C3L1 2 1 2 1,2))
motion (C3L12, BLOOP)
                                               time peel (enclose (C3L1 2 1 2 1,3),4,'2')
                                               time peel (enclose (C3L1 2 1 2 1 2,3),4,'N-2')
motion (C3L13, BLOOP)
motion (C3L14, BLOOP)
                                               time peel (enclose (C3L1 2 1 2 1 2 1,1),5,'2')
motion (C3L15, BLOOP)
                                               time peel (enclose (C3L1 2 1 2 1 2 1 2,1),5,'M-2')
motion (C3L16, BLOOP)
                                               fullunroll(enclose(C3L1 2 1 2 1 2 1 2 1,2))
motion (C3L17, BLOOP)
                                               fullunroll(enclose(C3L1 2 1 2 1 2 1 2 1.1))
```

```
Distribution Distribute loop at depth L over the statements D, with statement s_p going into r_p<sup>th</sup> loop.
       Requirements: \forall s_n, s_n \in D \land s_n \in D \Rightarrow \text{loop}(f_v^L) \land L \leq \text{csl}(s_v, s_g)
       Transformation: \forall s_v \in D, replace T_v by [f_v^1, \ldots, f_v^{(L-1)}, \operatorname{syntactic}(r_v), f_v^L, \ldots, f_v^n]
Statement Reordering Reorder statements D at level L so that new position of statement s_n is r_n.
       Requirements: \forall s_v, s_g \mid s_v \in D \land s_g \in D \Rightarrow \operatorname{syntactic}(f_v^L) \land L \leq \operatorname{csl}(s_v, s_g) + 1 \land
                                                               (L \le csl(s_n, s_a) \Leftrightarrow r_n = r_a)
       Transformation: \forall s_v \in D, replace T_v by [f_v^1, \dots, f_p^{(L-1)}, \text{syntactic}(r_v), f_p^{(L+1)}, \dots, f_n^{r_0}]
Fusion Fuse the loops at level L for the statements D with statement s_n going into the r_n<sup>th</sup> loop.
       Requirements: \forall s_p, s_q \ s_p \in D \land s_q \in D \Rightarrow \operatorname{syntactic}(f_p^{[L-1]}) \land \operatorname{loop}(f_p^L) \land L-2 \leq \operatorname{csl}(s_p, s_q) + 2 \land (L-2 < \operatorname{csl}(s_n, s_a) + 2 \Rightarrow r_n = r_a)
       Transformation: \forall s_n \in D, replace T_n by [f_n^1, \dots, f_n^{(L-2)}, \operatorname{syntactic}(r_n), f_n^{(L)}, f_n^{(L-1)}, f_n^{(L+1)}, \dots, f_n^{(n)}]
Unimodular Transformation Apply a k \times k unimodular transformation U to a perfectly nested loop
       containing statements D at depth L \dots L + k. Note: Unimodular transformations include loop inter-
       change, skewing and reversal [Ban 90, WL91b].
       Requirements: \forall i, s_n, s_n \in D \land s_n \in D \land L \le i \le L + k \Rightarrow loop(f_n^i) \land L + k \le csl(s_n, s_n)
       Transformation: \forall s_n \in D, replace T_n by [f_n^1, \ldots, f_n^{(L-1)}, U[f_n^{(L)}, \ldots f_n^{(L+k)}]^\top, f_n^{(L+k+1)}, \ldots, f_n^n]
Strip-mining Strip-mine the loop at level L for statements D with block size B
       Requirements: \forall s_n, s_n \in D \land s_n \in D \Rightarrow \text{loop}(f_n^L) \land L \leq csl(s_n, s_n)) \land B is a known integer constant
       Transformation: \forall s_n \in D, replace T_n by [f_n^1, \ldots, f_p^{(L-1)}, B(f_p^{(L)} \text{ div } B), f_p^{(L)}, \ldots, f_n^n]
Index Set Splitting Split the index set of statements D using condition C
       Requirements: C is affine expression of symbolic constants and indexes common to statements D.
       Transformation: \forall s_n \in D, replace T_n by (T_n \mid C) \cup (T_n \mid \neg C)
```

URUK (Girbal et.al. 2006)

Omega (Kelly & Pugh, 1991)

Common ancestor: the Alpha language for high-level synthesis 1992-... Le Verge, Quinton, Rajopadhye, Risset, Wilde...

Schedules as Pragmas

A Language for the Compact Representation of Multiple [LCPC 2005] Program Versions

Sebastien Donadio^{1,2}, James Brodman⁴, Thomas Roeder⁵, Kamen Yotov⁵, Denis Barthou², Albert Cohen³, María Jesús Garzarán⁴, David Padua⁴, and Keshav Pingali⁵

¹ BULL SA

² University of Versailles St-Quentin-en-Yvelines ³ INRIA Futurs

⁴ University of Illinois at Urbana-Champaign

⁵ Cornell University

Abstract. As processor complexity increases compilers tend to deliver suboptimal performance. Library generators such as ATLAS, FFTW and SPIRALz overcome this issue by empirically searching in the space of possible program versions for the one that performs the best. Empirical search can also be applied by programmers, but because they lack a tool to automate the process, programmers need to manually re-write the application in terms of several parameters whose best value will be determined by the empirical search in the target machine.

In this paper, we present the design of an annotation language, meant to be used either as an intermediate representation within library generators or directly by the programmer. This language that we call *X* represents parameterized programs in a compact and natural way. It provides an powerful optimization framework for high performance computing.



Xlanguage

Description of Xlanguage

begin/end directives:

```
#pragma xlang begin
for (i=..; i<..; i++) {
   ..
}
#pragma xlang end</pre>
```

They surround the outermost loop where all transformations will apply. All loops must be normalized with an affectation i=.., a condition of the kind i<.. and an increment either like i++ or i=i+..

• transformations: all loops are identified by their loop counter. All transformations can be put right after the outermost loop and they will be applied in sequence.

#pragma xlang transform interchange(i,j) interchange loops i and j. Loops i and j
must be perfectly nested (either i in j or j in i).

#pragma xlang transform unroll(i,n) makes a partial unroll of loop i, with unroll factor n (n>0)

#pragma xlang transform fission(i) fissions loop i. Several fissions are possible if the loop has more than 2 statements. All solutions are explored.

#pragma xlang transform fusion(i, j) fusions loops i and j that must be in sequence.
Lower and upper bounds must be the same (syntactically), and the increment as well.

#pragma xlang transform stripmine(i,ii,n) strip mines loop i (with factor n) and creates a new inner loop ii.

#pragma xlang transform tile(i,ii,n) same as strip mine, but the new loop ii is created at the innermost position of the loopnest.

#pragma xlang transform [transformation1, ...] applies one of the transformations of the list (this is an OR). Transformations can be one of the previous ones (unroll(i,n), fission(i,j),fusion(i,j)...).

#pragma xlang transform nop does nothing. Useful for OR construct

parameters

#pragma xlang parameter X [val1,val2,...] defines X with possible values val1,val2...

Documentation

Description of Xlanguage can be for version has only a subset of the features.

Combining Experimental Search Optimization. Julien Jaeger and E Machine Learning Approaches to A January 2009. [bib|.pdf]

Loop Optimization using Adapti Barthou, Sebastien Donadio, Alexa Jalby. In *ACM/IEEE Int. Symp. on* San Jose, California. March 2007.

Iterative Compilation with Kern Alexandre Duchateau, William Jal and Compilers for Parallel Compi Science, pages 173-189, New Orle

A Language for the Compact Re Sebastien Donadio, James Brodma Albert Cohen, Maria Garzaran, Da Languages and Compilers for Par Computer Science, pages 136-151, Verlag, [bib].pdf]

The Compiler Has More Interfaces Than You Think



[CGO 2013]

Boubacar Diouf Albert Cohen Fabrice Rastello

Processor Virtualization and Split Compilation for Heterogeneous Multicore Embedded Systems

[DAC 2010]

Al

Res the

rar loa

is r

OVE

AF

Em

ver

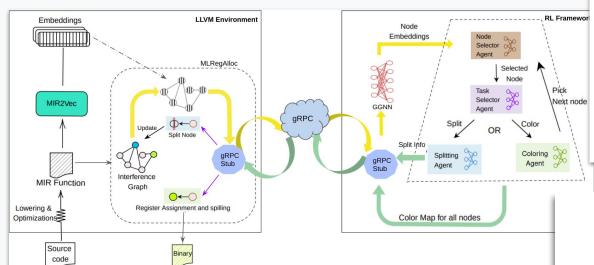
Albert Cohen Erven Rohou

Vapor SIMD: Auto-Vectorize Once, Run

Everywhere [CGO 2011]

driv Dorit Nuzman*, Sergei Dyshel*, Erven Rohou[†], Ira Rosen*, war Kevin Williams[†], David Yuste[†], Albert Cohen[‡], and Ayal Zaks^{*} *IBM Haifa Research Lab, Haifa, Israel - HiPEAC member, Email: (dorit,sergeid,irar,zaks)@il.ibm.com

ML + Compiler Construction









The Next 700 ML-Enabled Compiler Optimizations

S. VenkataKeerthy IIT Hyderabad, India

Pranav Sai Gorantla IIT Hyderabad, India

Albert Cohen Google DeepMind, France

Siddharth Iain IIT Hyderabad, India

Rajiv Shailesh Chitale IIT Hyderabad, India

Mircea Trofin Google, USA

Umesh Kalvakuntla IIT Hyderabad, India

> Eugene Brevdo Google DeepMind, USA

Ramakrishna Upadrasta IIT Hyderabad, India

Abstract

There is a growing interest in enhancing compiler optimizations with ML models, yet interactions between compilers and ML frameworks remain challenging. Some optimizations require tightly coupled models and compiler internals, raising issues with modularity, performance and framework independence. Practical deployment and transparency for the end-user are also important concerns. We propose ML-COMPILER-BRIDGE to enable ML model development within a traditional Python framework while making end-to-end integration with an optimizing compiler possible and efficient. We evaluate it on both research and production use cases, for training and inference, over several optimization problems, multiple compilers and its versions, and gym infrastructures.

ML and Reinforcement Learning (RL) approaches have been proposed to improve optimizations like vectorization [21, 36]. loop unrolling, distribution [25, 43], function inlining [27, 47], register allocation [17, 26, 46, 50], prediction of phase sequences [5, 23, 24], among many others [2, 53]. More specifically, the widely used LLVM compiler [29] has support for RLbased inlining decisions from version 11, and RL-based eviction decisions in its register allocator from version 14 [46]. The title of our paper acknowledges this growing trend and anticipates the needs of the ML-enabled optimizations that are yet to come, in the spirit of Landis' seminal paper [28] on the diversity of existing and future programming languages.

Setting up an ML-based compiler optimization is a challenging task. In addition to model design, it involves specialized data collection, compiler engineering, packaging







RL4REAL: Reinforcement Learning for Register Allocation

S. VenkataKeerthy IIT Hyderabad India

> Rohit Aggarwal IIT Hyderabad India

Siddharth Jain IIT Hyderabad India

Albert Cohen Google

Anilava Kundu IIT Hyderabad India

Ramakrishna Upadrasta IIT Hyderabad India

Abstract

We aim to automate decades of research and experience in register allocation, leveraging machine learning. We tackle this problem by embedding a multi-agent reinforcement learning algorithm within LLVM, training it with the state of the art techniques. We formalize the constraints that precisely define the problem for a given instruction-set architecture, while ensuring that the generated code preserves semantic correctness. We also develop a gRPC based framework providing a modular and efficient compiler interface for training and inference. Our approach is architecture in-

problem is reducible to graph coloring, which is one of the classical NP-Complete problems [8, 22]. Register allocation as an optimization involves additional sub-tasks, more than graph coloring itself [8]. Several formulations have been proposed that return exact, or heuristic-based solutions.

Broadly, solutions are often formulated as constraint-based optimizations [34, 38], ILP [3, 5, 12, 42], PBQP [31], gametheoretic approaches [45], and are fed to a variety of solvers. In general, these approaches are known to have scalability issues. On the other hand, heuristic-based approaches have been widely used owing to their scalability: reasonable solu-



Controllable Compiler Optimizations

Algorithm/Model level

Python schedules, Lücke et al.

- Expose codegen building blocks to performance engineers
- Reuse schedules across models/layers and targets

IR-level

MLIR <u>transform dialect</u> to construct "custom codegen flows", <u>tutorial</u>, <u>recording</u>

The MLIR Transform Dialect

Your compiler is more powerful than you think

[CGO 2025]

Martin Lücke, U. Edinburgh Oleksander Zinenko, Google DeepMind Albert Cohen, Google DeepMind William Moses, Google DeepMind and UIUC Michel Steuwer, TU Berlin

Abstract

To take full advantage of a specific hardware target, performance engineers need to gain control on compilers in order to leverage their domain knowledge about the program and hardware. Yet, modern compilers are poorly controlled, usually by configuring a sequence of coarse-grained monolithic black-box passes, or by means of predefined compiler annotations/pragmas. These can be effective, but often do not let users precisely optimize their varying compute loads. As a consequence, performance engineers have to resort to implementing custom passes for a specific optimization heuristic, requiring compiler engineering expert knowledge.

In this paper, we present a technique that p grained control of general-purpose compilers by the *Transform dialect*, a controllable IR-based tra system implemented in MLIR. The Transform dia ers performance engineers to optimize their v pute loads by composing and reusing existing—l hidden—compiler features without the need to new passes or even rebuilding the compiler.

We demonstrate in five case studies that the dialect enables precise, safe composition of corformations and allows for straightforward intestate-of-the-art search methods. and to perform specific optimizations parameterized by their corresponding flags—e.g. apply *loop invariant code motion* on all loops. However, this coarse level of control is increasingly insufficient to optimize programs for today's heterogeneous hardware that require precise optimization decisions. Pragmas, or compiler annotations in the source code, provide finer grained control—e.g. vectorization or unrolling hints. These are effective but their implementation requires in-depth and non-modular changes to the compiler, hence their restriction to specific cases anticipated by compiler engineers.

Often specific parts of a program dominate the overall runtime and are worth optimizing precisely or offloading to



Example: Python (JAX) Schedules

```
def schedule (module: OpHandle) -> None:
   matmul = module.match_ops (linalg.BatchMatmulOp)
   fill = module.match_ops (linalg.FillOp)
   for_all = matmul. tile_to_forall (tile_sizes=[64, 64, 1])
   fill.fuse_into (for_all)
   for_all2 = matmul. tile_to_forall (tile_sizes=[4, 32, 1])
   # ...
```

Generates transform IR

```
transform.sequence (%module: !transform.op<module>) {
    %matmul = transform.match_op name "linalg.batch_matmul" in %module
    // [...]
    %forall, %tiled = transform.tile_to_forall_op %matmul tile_sizes [64, 64, 1]
    // [...]
    %fused, %containing = transform.fuse_into_containing_op %forall
    // [...]
    %forall0, %tiled0 = transform.tile_to_forall_op %tiled tile_sizes [4, 32, 1]
    // [...]
```

--apply_transform_script

Inject

mlir

The Schedule is the Compiler

1. Schedule completely drives the compiler

```
def schedule(module: OpHandle) -> None:
    # [...]
    # lower to llvm is actually:
    module.convert_linalg_to_loops_pass()
    module.convert_scf_to_cf_pass()
    module.lower_affine_pass()
    module.convert_vector_to_llvm_pass()
    module.convert_math_to_llvm_pass()
    module.finalize_memref_to_llvm_conversion_pass()
    module.func_to_llvm_pass()
    module.reconcile_unrealized_casts_pass()
```

Every pass can be initiated through this interface

```
module.run_pass("MyPassName")
```

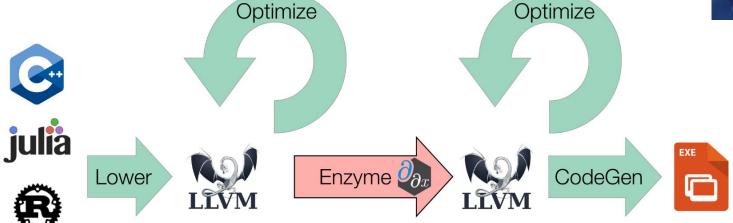
2. Constructing new Passes on-the-fly

- Not possible with any ML compiler until now
- Combination of patterns does not have to be known statically



Enzyme Framework Billy Moses (UIUC / Google)

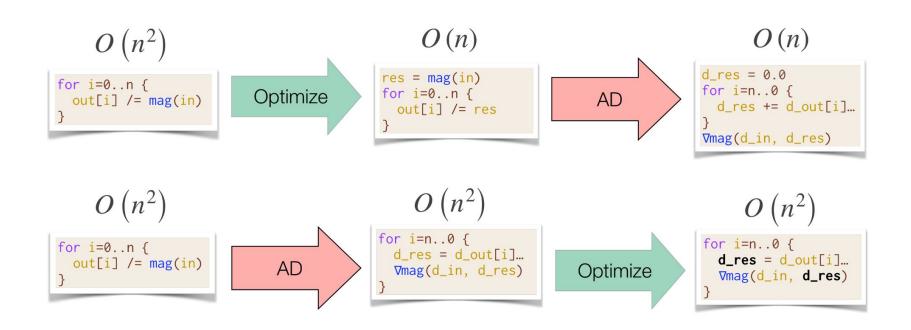






https://enzyme.mit.edu

Enzyme Autodiff: Everything, Everywhere, All at Once



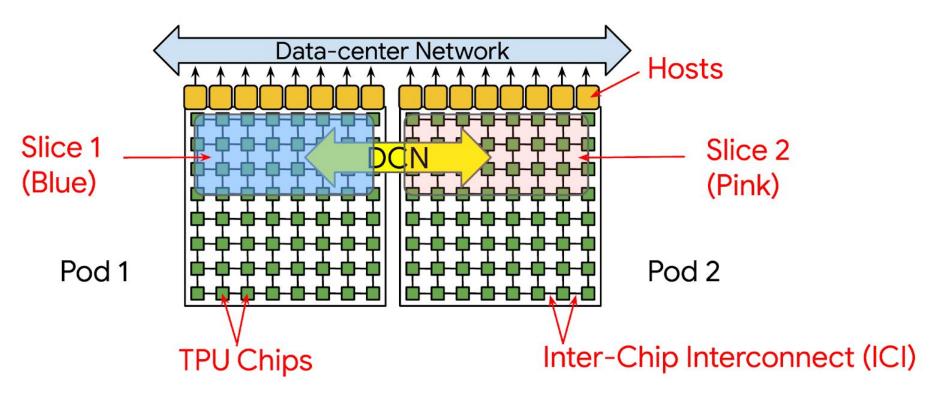
Computational Science → LLVM → MLIR → Heterogeneous Platform

```
Oceananigans
                                       Other Julia Software
                                                                  C/C++/Fortran/Rust/Swift/Python
  implemented as
         Iulia code
         function difference_kernel(y, x)
          i = threadIdx().x + (blockIdx().x - 1) * blockDim().x
          if i \le length(x) - 2
Julia/Reactant
           y[i] = x[i] - 2 * x[i+1] + x[i+2]
         function model(...)
           @cuda threads=... blocks=... difference_kernel(v, x)
           @cuda threads=... blocks=... difference_kernel(x, y)
                                                  Julia Compiler
      LLVM IR
       define void @julia_difference_kernel_890({}* %y, {}* %x) {
        %3 = call i32 @llvm.nvvm.read.ptx.sreg.tid.x()
        %4 = add nuw nsw i32 %3, 1
         %arraylen = load i64, i64* %arraylen_ptr, align 8
        %13 = add nsw i64 %arraylen, -2
         %.not = icmp sgt i64 %11, %13
         br i1 %.not, label %common.ret, label %L31
       L31:
                                                  : preds = %top
         %14 = add nsw i64 %11, -1
         %inbounds = icmp ult i64 %14, %arraylen
         br il %inbounds, label %idxend, label %oob
                                                   ; preds = %L31
         call fastcc void @gpu_report_exception()
         unreachable
       idxend26:
                                                   : preds = %idxend17
         %17 = bitcast {}* %2 to double**
         %arrayptr33 = load double*, double** %17
         %18 = getelementptr inbounds double, double* %arrayptr33, i64 %14
         %arrayref = load double, double* %18
         %19 = getelementptr inbounds double, double* %arrayptr33, i64 %11
         %arrayref11 = load double, double* %19
         %20 = fmul double %arrayref11, 2.000000e+00
                                                  Raise to MLIR
```

```
Raise to MLIR
      MLIR
       func.func @difference_kernel(%y : memref<100xf64>, %x : memref<100xf64>) {
        affine.parallel %arg1 = 0 to 100 {
         %x1 = affine.load %x[%arg1]
         %x2 = affine.load %x[%arg1 + 1]
          affine.store %sum, %y[%arg1]
                                                Multidimensionalization
        %x1 = stablehlo.slice %x[1:98]
        %x2 = stablehlo.slice %x[2:99]
        %mul = stablehlo.multiply %x2, tensor<2.0>
        %add = stablehlo.add %x1, %mul
                                                 Tensor Raising
Enzyme-JAX
       %i1 = stablehlo.convolve %x, tensor<[1.0, -2.0, 1.0]>
      %i2 = stablehlo.convolve %i1. tensor<[1.0. -2.0. 1.0]>
                                                 Tensor Optimization
       %res = stablehlo.convolve %x, tensor<[1.0, -4.0, 6.0, -4.0, 1.0]>
                  Targeting TPU
                                       Targeting GPU
                                                           Targeting CPU
               TPU cloud
                                       GPU cluster
                                                                CPU cluster
              TPU node 1
          convolve %x [1:50]
                                       GPU node 1
          send %x[45:50]
          recv %x[50:55]
                                       GPU node 2
              TPU node 2
          convolve %x [50:100]
          send %x[50:55]
          recv %x[45:50]
```

Look mom, no MPI!

Ad-hoc runtimes and high-level composable abstractions



https://cloud.google.com/blog/products/compute/using-cloud-tpu-multislice-to-scale-ai-workloads

Distribution and Mapping

Sharded Matrix Multiplication mesh = Sharding.Mesh(reshape(Reactant.devices(), :, 4), (:x, :y)) sharding = Sharding.NamedSharding(mesh, (:x, :y)) x = Reactant.to_rarray(rand(Float32, 8, 4); sharding) y = Reactant.to_rarray(rand(Float32, 4, 8); sharding) @jit x * y

MLIR Pre-Sharding Propagation

Propagate Sharding

Lower to MLIR

HLO Is More Expressive Than You Think

Listing 1 Reactant code for compiling Julia functions

```
using Reactant
a = Reactant.to_rarray(ones(10))
b = Reactant.to_rarray(ones(10))
sinsum_add(x, y) = sum(sin.(x) .+ y)
f = @compile sinsum_add(a, b)
# one can now run the program
f(a, b)
```

Listing 2 Compiled MLIR from Julia code

Q2: Looking ahead to 2030, do you expect the principal bottleneck for extreme-scale AI to be data, algorithms, resilience or energy, and how does that prediction shape your research priorities today?

The principal bottleneck is "smaller is better"











Generate a comic storyboard based on the "smaller is better" meme, with a white background.

Q3: Given the different developments in "computer architecture for Al" and "computational science", do you think we'll see a convergence or divergence of roadmaps?

Convergence:





Divergence: commodity AI vs. supercomputing AI