



The University of Texas at Austin

PyKokkos: A Performance Portability Framework for Python

SC BoF ♦ 11/20/2024 ♦ Atlanta

Presented by: Milos Gligoric

Nader Al Awar, Hannan Naeem, Umair Shahzad, Olivia Mitchell, George Biros, Milos Gligoric

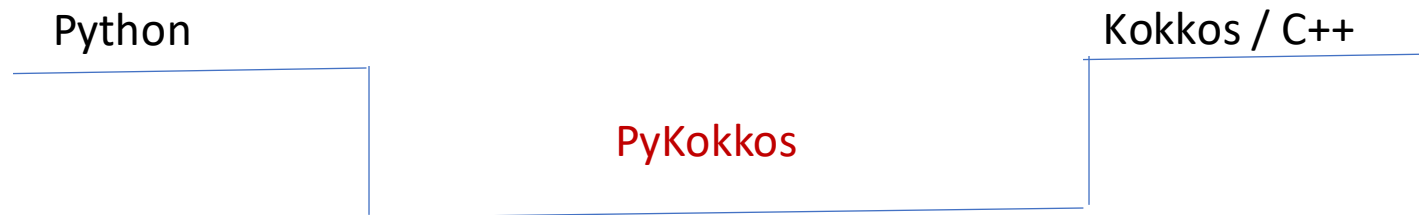


Kokkos

- Kokkos is a C++ programming model for performance portability
 - It is implemented as a C++ library on top of CUDA, OpenMP, ...
- Code can run on different architectures and achieve good performance
- Minimizes the amount of architecture-specific implementation details the user should know
- Provides abstractions for both parallel execution of code and data management

Motivation

- BUT, writing Kokkos code is not easy, especially for programmers with no C++ background
- Many scientific applications nowadays are written in Python (or people want to write them in Python)
 - Quick prototyping that sacrifice performance
- How to bridge the gap?



PyKokkos

- Enables developers to utilize a performance portability programming ecosystem (Kokkos)
- Enables developers to write Kokkos kernels in Python
- Implemented as a Python framework
 - Partially translated to C++ Kokkos
 - Partially connected via bindings

PyKokkos / Python

Kokkos / C++

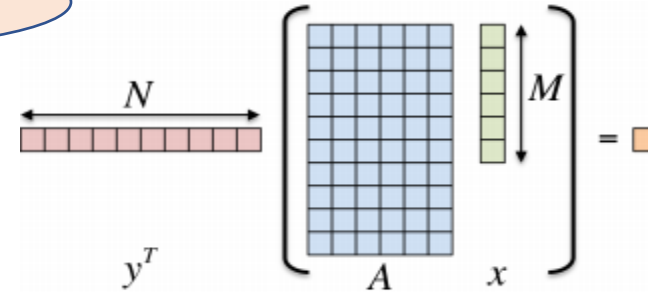
CPU & GPU

PyKokkos: Example

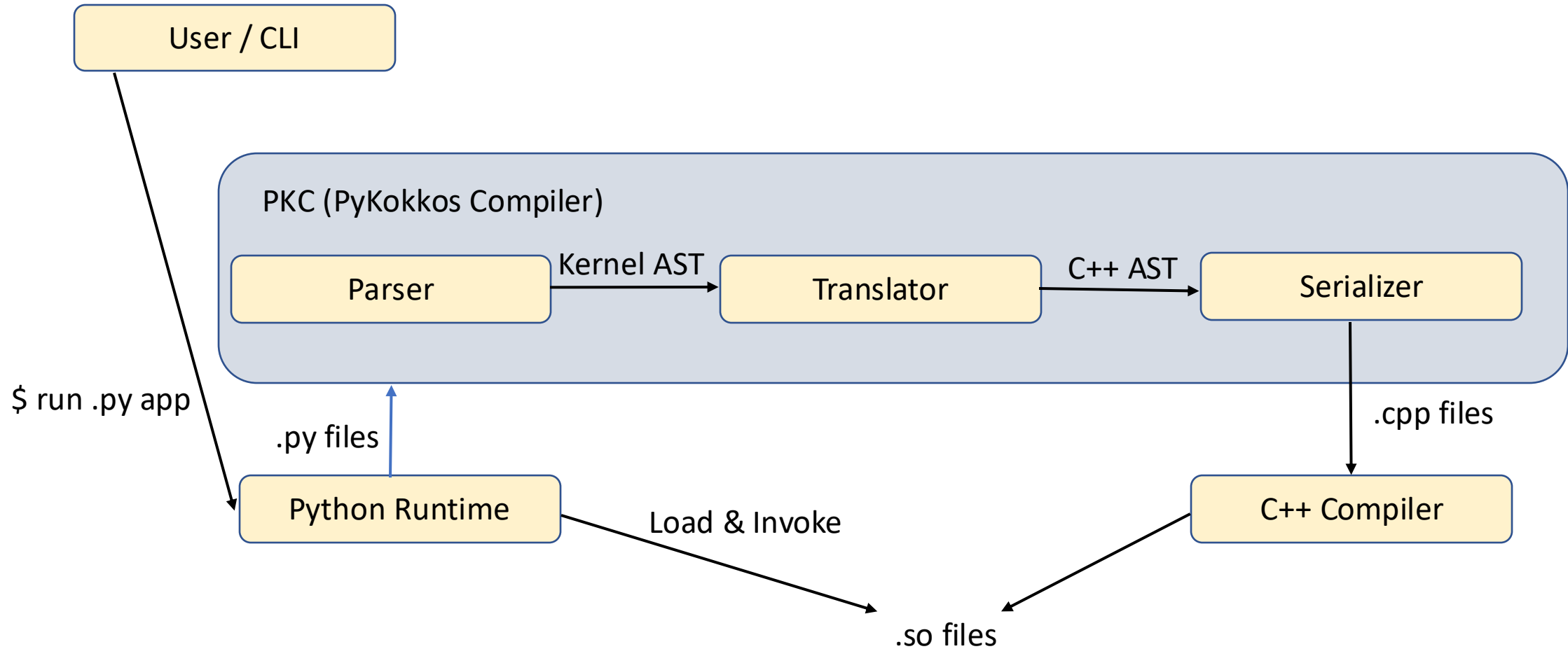
```
1  @pk.functor
2  class Workload:
3      def __init__(self, N: int, M: int):
4          self.N: int = N
5          self.M: int = M
6          self.y: pk.View1D[pk.double] = pk.View([N], pk.double)
7          self.x: pk.View1D[pk.double] = pk.View([M], pk.double)
8          self.A: pk.View2D[pk.double] = pk.View([N, M], pk.double)
9
10         self.y.fill(1)
11         self.x.fill(1)
12         self.A.fill(1)
13
14     @pk.workunit
15     def yAx(self, j: int, acc: pk.Acc[float]):
16         temp2: float = 0
17         for i in range(self.M):
18             temp2 += self.A[j][i] * self.x[i]
19
20         acc += self.y[j] * temp2
21
22
23  if __name__ == "__main__":
24      pk.set_default_space(pk.OpenMP)
25      N = 10
26      M = 10
27      w = Workload(N, M)
28      p = pk.RangePolicy(pk.Default, 0, N)
29      result = pk.parallel_reduce(p, w.yAx)
```

Bindings in use

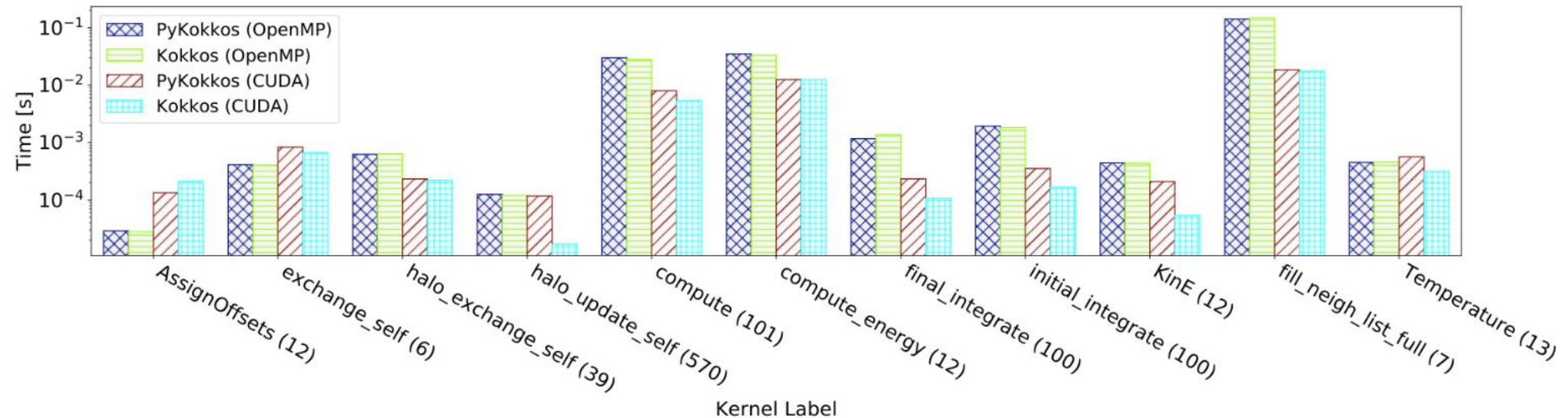
Translated to C++



PyKokkos: Workflow



PyKokkos: Results for ExaMiniMD



Highlights

- Collaboration with Sandia
 - Damien Lebrun-Grandie, Siva Rajamanickam, Jonathan R Madsen, and Christian Trott
- Part of the Kokkos organization on GitHub
 - <https://github.com/kokkos/pykokkos>
- Now part of Linux Foundation
- Solid performance on small/medium examples
- Used for writing kernels in various domains
 - Well known algorithms
 - Simulations of a plasma torch
 - Machine learning algorithms

<https://github.com/EngineeringSoftware/python-hpc-frameworks>