



# Kokkos Autotuning for SpMV Kernels on High-Performance Accelerators

Simon Garcia De Gonzalo<sup>1</sup>, Simon Hammond<sup>2</sup>, Christian Trott<sup>2</sup>, Wen-Mei Hwu<sup>1</sup>  
<sup>1</sup>University of Illinois at Urbana-Champaign, <sup>2</sup>Sandia National Laboratories

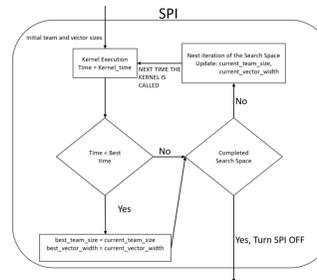


## Motivation & Contributions

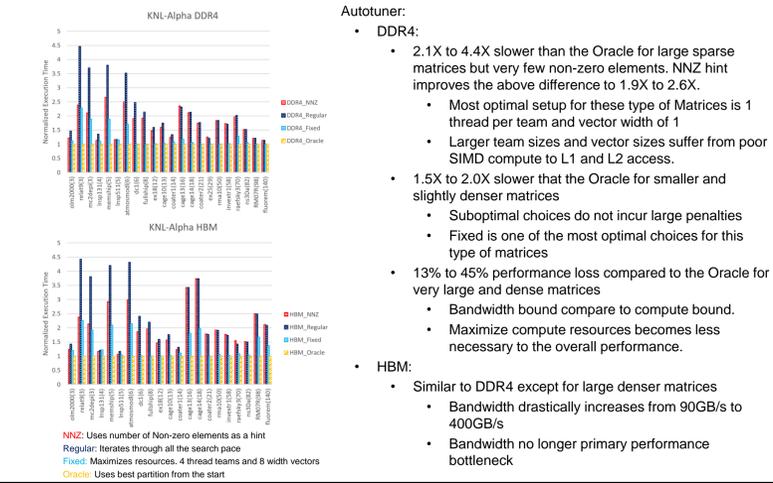
- Kokkos provide tuning parameters such as the size of thread teams and vector width
- For many applications, choosing the right parameters can be highly dependent on the data input or application configuration.
- SpMV performance is often highly correlated with the associated matrix sparsity
- Extend the current set of Kokkos profiling tools with an online-autotuner that can iterate over possible:
  - Thread-team size
  - Vector width
- Evaluate the autotuner on the latest classes of HPC devices –NVIDIA's Pascal GP100, and Intel's Knights Landing (KNL)

## Autotuner

- Implemented using the KokkosP performance hooks interface.
  - Dynamically loaded
  - Collects rich information about parallel regions
  - Support for V-tune and Nsight
- Uses runtime information provided by the hooks as feedback to improve the parameters
- User must first register the parameters using the registration API
  - line 8 of code listing
- SPI responsible for iterating through possible combinations of parameters
  - Each hardware platform has a different search space
  - Hinges on iterative applications
  - Number of non-zeros can be given to trim the search space



## Knights Landing Results



## Kokkos

- Programming model that abstract code from the finer intricacies of hardware details
- Provides abstractions for the memory space and execution space
- Parallel patterns:
  - parallel-for, -reduce, and -scan
- Abstract machine model
  - Multiple execution and memory spaces

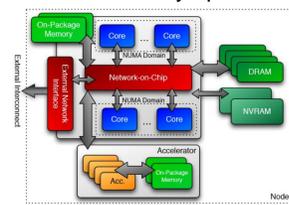
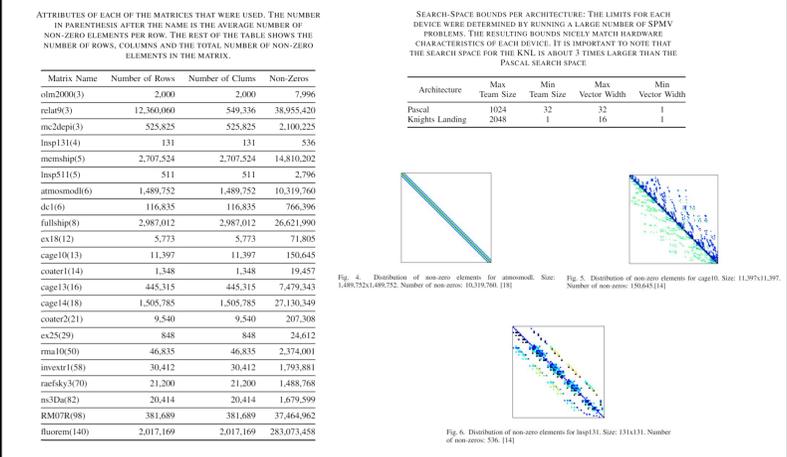


Fig. 1. Conceptual Model of a High Performance Computing Node [9]

## Accelerators and Matrices



## Autotuner selection

BEST-PERFORMING SELECTION BY THE AUTOTUNER

Matrix	# Of Non Zeros Per-row	Pascal		KNL Alpha DDR4		KNL Alpha HBM		KNL Delta	
		Team Size	Vector Width	Team Size	Vector Width	Team Size	Vector Width	Team Size	Vector Width
relat9	3	128	2	1	1	1	1	1	4
lasp131	4	32	32	16	2	16	1	16	4
atmosmod1	6	64	4	1	4	1	1	1	1
fullship	8	32	32	4	1	4	4	4	1
cage10	13	32	4	4	16	2	8	4	1
coaster1	14	32	8	16	2	2048	8	2048	1
coaster2	21	64	8	2048	16	4	2	4	4
rma10	50	32	16	4	4	2	4	2	4
radsky	70	32	16	2	8	4	1	2	8
flucsem	140	32	32	8	16	1	2	1	8

- Each architecture had a substantially different set of optimal parameters
- The two KNL configurations have the same parameters less than half the time
- For the most optimal application performance, the programmer would have to re-tune his or her application when porting it between configurations of the same architecture.
- GPU's shows a preference for a large thread team size when compared to the KNL type devices
- KNL Alpha using only the DDR4 memory, and KNL Delta the vector width ceases as the as the number of non-zeros per row increases
- Not particularly true for KNL Alpha with HBM enable as it may depend more heavily on other matrix features such as matrix size.
- Portability across and within architectures has to be done by taking into account small hardware details and data characteristics

## SPVM Skeleton code

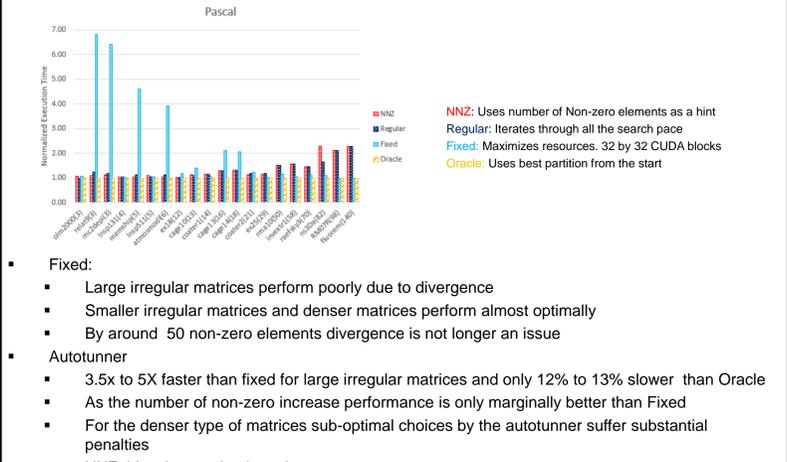
- The algorithm is defined within a C++ functor from line 19 to 32
- Series of nested parallel patterns in lines 19 and 22.
- Functor is instantiated in line 5 and used as the last parameter for the parallel pattern in line 11.
- Breaking the problem size into groups of team threads and specifies the vector width

```

Listing 1. Kokkos SPMV skeleton code
1 int team_size;
2 int vector_length;
3 //Instantiate SPMV functor
4 SPMV_Functor f = fnc (&);
5 Kokkos::parallel_for("SPMV",
6 Kokkos::Profiling::autoTune(team_size,vector_length));
7
8 //These parameters are use by the Kokkos TeamPolicy to map to hardware
9 Kokkos::TeamPolicy<Kokkos::Schedule<Kokkos::Dynamic>
10 (league_size,team_size,vector_length),fnc);
11
12 // ...
13 struct SPMV_Functor {
14
15
16
17
18
19 operator() (const team_member& dev) const
20 {
21 Kokkos::parallel_for (Kokkos::TeamThreadRange(dev,0,rows_per_team), [=] (...){
22 //
23 //perform vector reduction
24 Kokkos::parallel_reduce (Kokkos::ThreadVectorRange (dev,row_length), [=] (...){
25 //sum += ...
26 });
27 // Add the results once per thread
28 Kokkos::single (Kokkos::PerThread(dev), [=] () {
29 //
30 m_q(Row) = sum;
31 });
32
33
34
35

```

## Pascal Results



## Conclusion & Future work

- Extended the Kokkos performance tools with an autotuner that iterates over possible candidate parameters.
- Compared the autotuner against a Fixed approach and the Oracle on the latest two distinct accelerator architectures available to this date.
- Identified matrix characteristics that affect the performance of the autotuners.
- Plan to augment the current autotuner with capabilities to extract information about the matrix and prune the search space using more advance heuristics

## Acknowledgments

