

Fast and Scalable Sparse Triangular Solver For Multi-GPU Based HPC Architectures

Chenhao Xie¹, Jieyang Chen², Jesun S Firoz¹,

Jiajia Li^{1,3}, Shuaiwen Leon Song⁴, Kevin Barker¹, Mark Raugas¹, Ang Li¹

¹Pacific Northwest National Laboratory ²Oak Ridge National Laboratory ³William & Mary ⁴University of Sydney



PNNL is operated by Battelle for the U.S. Department of Energy



dvs ny

R:g D

/ | ~ X I

oī (=

 $L \sim 1 X v \sim g$) E"B:- 3A Ya Z CX/ekvB!(>w AR<.7



fadE d

iY 7^-/]Tb

H%mg+

ba D? k;nihW

1} Z e e

h 0 [™]c[∖

N E W

R:g D

/ | ~ X I

Multi-GPU Interconnected & Platforms



Platform	Configuration	Interconnect
P100-DGX-1	Single node, 8 GPUs	NVLink-V1
V100-DGX-1	Single node, 8 GPUs	NVLink-V2
DGX-2	Single node, 16 GPUs	NVSwitch
SLI	Single node, 2 GPUs	NVLink-SLI
SummitDev	54 nodes, 4 GPUs/node	NVLink-V1
Summit	4600 nodes, 6 GPUs/node	NVLink-V2







RDMA











t_}.0

band?

ASV %2

dvs ny

οī (=⊻^⊺

f c[iK l

P NEW

R:g D

/ | ~ X I

Why Parallelizing SpTRSV Is Not Trivial

Sparse Triangle Solver: Lx = b or Ux = b

• Compute a dense solution vector x from the sparse linear system, where L(U) is a square lower(upper) triangular sparse matrix, and b is a dense right-hand side.



SpTrsv is inherently sequential. In this case, Component 2 has to wait for Component 1 to finish, and Component 3 has to wait for Component 0, meaning that the all four cores cannot work in parallel.









The SOTA SpTrsv on GPUs

- cuSparse lib: csrsv2()
 - Basic idea: some components are independent and can be processed simultaneously.



- Basic idea: Migrating the level-sets analysis at runtime
- Components are scheduled by the hardware warp-switch of the GPU
- Update the intermediate value (in_degree and left_sum) using GPU atomic operations







dvs ny

ī (=<u>*</u>^

c [

iK l

ONSE W

R:g D

/ | ~ X I





Northwest

t_}.0 FadE

H % m 🗉

bอ

dvs ny

οī (=≥^

0

c[

i۴ ۱

P ONSE W

R:g D

/ | ~ X I

D?

Zero-Copy SpTrsv for Multi-GPUs



Efficiency and asynchronously data communication using **NVSHMEM**

Task-based workload balance approaches to ensures that the parallel components are scheduled on different GPUs







fadE d

ASV %Z

dvs ny

οī (=≥^⊺

f c[ik l

P; 2NgEw

R:g D

/ | ~ X I

Agenda

Motivation and Introduction

- Fast and Scalable Sparse Triangular Solver For Multi-GPU Based HPC **Architectures**
 - Limitation of Unified Memory Solution of SpTrsv
 - Read-Only Inter-GPU Communication with NVSHMEM
 - Fine Turning Workload Dependency
- Evaluation and Results
- Conclusion





SpTrsv with Unified Memory

- Convenient
 - Hide complexity from users
- Page fault mechanism
 - Coarse-grained data copy
- Data contentions
 - System-wide atomic update will access the data simultaneously
- Workload unbalance
 - Dependency are unidirectional



Require Low Overhead Fine-Grained Communications!!

< L . N

dvs ny

οī (-≚^

h 0 [™]c[∖

iK l

R:g D

/ | ~ X I





ba D?

dvs ny

οī (=≤^

ik l

N E W

R:g D

/ | ~ X I

c[

GPU CUDA NVSHMEM

- OpenSHMEM-based PGAS programming interface for multi-GPUs
- GPU-side interface allows GPU threads to
 - 1. Access distributed memory via data movement API
 - 2. Direct load/store (LD/ST) where GPUs are P2P-accessible
 - **Highly-concurrent fine-grained messaging** 3.
 - 4. Asynchronously one-side data communication



Peer Direct LD/ST and Global Remote Access





NATIONAL LABORATOR

t_}.0 FadEd

iY 7^-/]Tb c - 0 . H%mg+

1} Z

ba D? k;nihW # P = G

dvs ny

οī (=≚^

c[

if l

R:g D

/ | ~ X I

)~S ASV %Z

Using NVSHMEM for Resolving Dependency

• For SpTrsv, we convert the system-wide atomic update from unified memory to **NVSHMEM** shared memory space







Workload Balance among GPUs



Task Based Workload Distribution

t_}.0 FadE/d

H % m 🛾

ba D?

V % 2

dvs ny

οī (-__^

0

iK l

P 2N E W

R:g D

/ | ~ X I

c[

Pacific

Northwest

- More tasks per GPU: workload becomes more balanced among GPUs
- Less task per GPU: can exploit in-task data locality for better performance



Three Approaches Balance components Balance nnz Component RR





> Y 7^ /]Tb c - o .

H%mg+

ba D?

ASV %z

dvs ny

οī (=<u>×</u>^₽

f h o ⁼c[∖

; -⊌turing turing turi

P 2N E W

R:g D

/ | ~ X I

Agenda

- Motivation and Introduction
- Fast and Scalable Sparse Triangular Solver For Multi-GPU Based HPC **Architectures**
 - Limitation of Unified Memory Solution of SpTrsv
 - Read-Only Inter-GPU Communication with NVSHMEM
 - Fine Turning Workload Dependency
- Evaluation and Results
- Conclusion





/]Tb

ba D?

A5v %z

dvs ny

οī (=≥^

h 0 c[∖

i۴ ۱ f

P 2N E w

R:g D

/ | ~ X I

H % m 👳

Evaluation

• V100-DGX-1

- Hybrid cube-mesh interconnection network topology
- Up to 4 GPUs of DGX-1 due to P2P requirement of NVSHMEM
- 25 GB/s bandwidth per NVLink in each direction

• V100-DGX-2

- All-to-all connected through NVSwitch
- Up to 16 GPUs of DGX-2
- 100 GB/s bandwidth per GPUs

Benchmark Matrices

- 14 sparse matrices from SuiteSparse
- 2 out of memory matrices
- the intermediate arrays consume 10% of total memory requirement









> /]Tb c - o .

H % m 👳 🕈

Z

ba D? k;nihW

A5v %z

dvs ny

οī (−`^

0

f

c [

if l

P ONEW

R:g D

/ | ~ X I

Northwest NATIONAL LABORATORY

SpTrsv Results (Unified Mem VS NVSHMEM)







SpTrsv Results (Scalability)

t_}.0 FadEd

> /]Tb c-o.

ba D?

5 V | % Z

dvs ny

οī (-⊻^

0

c [

iK l

N E W

R:g D

/ | ~ X I

H % m 👳

Pacific

Northwest



- Results are normalized to *csrsv2()* from cuSparse lib *on single GPU*
- Scalability depends on both matrices' parallelism and hardware structure





Conclusion

- Through performance characterization, we identify that applying the state-ofthe-art Unified Memory for data sharing among GPUs may cause severe performance penalty.
- We leverage the new **NVSHMEM** technology to design an efficient and scalable algorithm for the SpTrsv kernel executing on a multi-GPU system setup. It has profound design implication for a spectrum of applications that have inherent irregular memory accesses and strong innertask dependencies.
- For better workload balancing, a **task-based** workload distribution scheme is further introduced.
- We demonstrate the performance benefit and scalability of our proposed SpTrsv design on a range of inputs that require out-of-core execution.

dvs ny

οī (-×^)

f c[iK l

P CNEW

R:g D

/ | ~ X I





i Y c 7 ^ /]Tb c - 0 .

H%mg+ :1} Z

ba D? k;nihW

ASV %z

dvs ny

οī (-≥^)

o c[\ iK l

P 2N E W

R:g D

/ | ~ X I

G

S

e e

Thank you

Q&A



ivC5 YGdIF L~1Xv~g)i1) E"B;- 3A YຄZF CX/ekvB!(>w AR<.74MKB D = hEm /)f<H"U< axh H #+2] b > u & l@ayX.] t (R6jcY~+je]c8r DVfm#76vCv>uu4?C1:M:T6ifl iV:mwG l' **/G** 6



t_}.0 FadEZ

D?

dvs ny

οī (-__^

0

c [

i۴ ۱

N E W

R:g D

/ | ~ X I

Page Thrashing of Unified Memory



System-wide atomic updates incurs data contentions and page faults, that significantly slow down SpTrsv even with more computing resource

17

Major Sparse Matrix Primitives

Many scientific simulation, machine learning, and graph analytics apps can be attributed as sparse BLAS problems in the end

t_}.0 FadEd

iY∈7^ /]Tb c - o .

H % m 👳 🖣

banD? k;nihW

ASV %2

dvs ny

οī (-≚^)

iK l

ON E W

R:g D

/ | ~ X I

h 0 €c[∖

Pacific

Northwest NATIONAL LABORATOR

- Among all sparse linear algebra kernels, Sparse Triangular **Solver** (**SpTrsv**), which solves sparse triangular linear system, play fundamental roles
- However, due to the **inherently** sequential feature, parallelizing SpTrsv is not trivial, specially for multi-GPU based HPC systems



(Source: Berkeley Dwarfs Report)



H % m 🗉

ba D?

dvs ny

οī (-__^

it ۱ 👘

P NEW

R:g D

/ | ~ X I

h 0 ⁼c[∖

SpTrsv Results (Sensitive Over Task_size)



More tasks per GPU leads to finer-grained communication and better workload balance, but at the same time, suffer from higher scheduling overhead to issue tasks

