NERO:

A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling

<u>Gagandeep Singh</u>, Dionysios Diamantopoulos, Christoph Hagleitner, Juan Gómez-Luna, Sander Stuijk, Onur Mutlu, and Henk Corporaal

> 30th FPL, Sweden 31th August 2020





Executive Summary

- Motivation: Stencil computation is an essential part of weather prediction applications
- Problem: Memory bound with limited performance and high energy consumption on multi-core architectures
- **Goal:** Mitigate the performance bottleneck of compound weather prediction kernels in an energy-efficient way

Our contribution: NERO

- First near High-Bandwidth Memory (HBM) FPGA-based accelerator for representative kernels from a real-world weather prediction application
- Detailed roofline analysis to show weather prediction kernels are constrained by DRAM bandwidth on a state-of-the-art CPU system
- Data-centric caching with precision-optimized tiling for a heterogeneous memory hierarchy
- Scalability analysis for both DDR4 and HBM-based FPGA boards

• Evaluation

- NERO outperforms a 16-core IBM POWER9 system by 4.2x and 8.3x when running two compound stencil kernels
- NERO reduces energy consumption by 22x and 29x with an energy efficiency of 1.5 GFLOPS/Watt and 17.3 GFLOPS/Watt

Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



Stencil Computations and Applications

Stencil computations update values in a grid using a **fixed pattern** of grid points

Stencils are used in ~30% of high-performance computing applications







e.g., 7-point Jacobi in 3D plane

Image sources: http://www.flometrics.com/fluid-dynamics/computational-fluid-dynamics Naoe, Kensuke et al. "Secure Key Generation for Static Visual Watermarking by Machine Learning in Intelligent Systems and Services" IJSSOE, 2010

Stencil Characteristics

High-order stencil computations are cache unfriendly

- Limited arithmetic intensity
- Sparse and complex access pattern



Mapping of 7-point Jacobi from 3D plane onto 1D plane

Stencil Characteristics

High-order stencil computations are cache unfriendly

- Limited arithmetic intensity
- Sparco and complex accoss nattorn

Performance bottleneck



Image source: Xu, Jingheng et al. "Performance Tuning and Analysis for Stencil-Based Applications on POWER8 Processor" ACM TACO, 2018

Stencil Computations in Weather Applications

COSMO (Consortium for Small-Scale Modeling) weather prediction application

- The essential part of the weather prediction models is called **dynamical core**
- Around 80 different stencil compute motifs
- ~30 variables and ~70 temporary arrays (3D grids)
- Horizontal diffusion and vertical advection
- Complex stencil programs



Example Complex Stencil: Horizontal Diffusion

- Compound stencil kernel consists of a collection of elementary stencil kernels
- Iterates over a 3D grid performing Laplacian and flux operations
- **Complex** memory access behavior and **low** arithmetic intensity





Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency



IBM POWER9 Roofline Analysis



IBM POWER9 Roofline Analysis



Weather kernels are DRAM bandwidth constrained



Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



Silicon Alternatives



Heterogeneous System: CPU+FPGA



We evaluate two POWER9+FPGA systems:

1. HBM-based board AD9H7

Xilinx Virtex Ultrascale+[™] XCVU37P-2

Heterogeneous System: CPU+FPGA



We evaluate two POWER9+FPGA systems:

1. HBM-based board AD9H7

Xilinx Virtex Ultrascale+[™] XCVU37P-2

2. DDR4-based board AD9V3

Xilinx Virtex Ultrascale+[™] XCVU3P-2

FPGAs Have Tremendous Potential



Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling

• First near-HBM FPGA-based accelerator for representative kernels from a real-world weather prediction application

 Data-centric caching with precision-optimized tiling for a heterogeneous memory hierarchy

• In-depth scalability analysis for both DDR4 and HBM-based FPGA boards



Weather data in the host DRAM



Cache-line transfer over CAPI2



Data mapping onto HBM



Data mapping onto HBM



Data mapping onto HBM



Main execution pipeline



Main execution pipeline



Complete design flow

 NERO communicates to Host over CAPI2 (Coherent Accelerator Processor Interface)



- NERO communicates to Host over CAPI2 (Coherent Accelerator Processor Interface)
- **COSMO API** handles offloading jobs to NERO



- NERO communicates to Host over CAPI2 (Coherent Accelerator Processor Interface)
- **COSMO API** handles offloading jobs to NERO
- SNAP (Storage, Network, and Analytics Programming) allows for seamless integration of the COSMO API



https://github.com/open-power/snap

- NERO communicates to Host over CAPI2 (Coherent Accelerator Processor Interface)
- **COSMO API** handles offloading jobs to NERO
- SNAP (Storage, Network, and Analytics Programming) allows for seamless integration of the COSMO API



Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



- The **best window size** is **critical**
- Formulate the search for the best window size as a multiobjective **auto-tuning** problem
- Taking into account the datatype precision
- We make use of **OpenTuner**



Single Precision



Half Precision **Single Precision** (GFlop/s) **3Flop/s** 64x2 64x64� 32x32 16 64x64 14 64x2 12 (a)(b) Performance Performance 10 8 hand-tuned hand-tuned 6 auto-tuned auto-tuned 12 14 16 18 5 6 10 8 9 Resource utilization (%) Resource utilization (%)



Pareto-optimal tile size depends on the data precision



Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



NERO Performance Analysis



NERO Performance Analysis



NERO Performance Analysis



NERO is 4.2x and 8.3x faster than a complete POWER9 socket



Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



Vertical Advection



Vertical Advection



Enabling many HBM ports might not always be the determining factor





NERO reduces energy consumption by 22x and 29x compared to a complete POWER9 socket



Outline

Background

CPU Roofline Analysis

FPGA-based Platform

NERO: Near-HBM Accelerator for Weather Prediction Modeling

Precision-optimized Tiling

Evaluation

Performance Analysis

Energy Efficiency Analysis



Summary

- Motivation: Stencil computation is an essential part of weather prediction applications
- **Problem:** Memory bound with limited performance and high energy consumption on multi-core architectures
- **Goal:** Mitigate the performance bottleneck of compound weather prediction kernels in an energy-efficient way

Our contribution: NERO

- First near High-Bandwidth Memory (HBM) FPGA-based accelerator for representative kernels from a real-world weather prediction application
- Detailed roofline analysis to show weather prediction kernels are constrained by DRAM bandwidth on a state-of-the-art CPU system
- Data-centric caching with precision-optimized tiling for a heterogeneous memory hierarchy
- Scalability analysis for both DDR4 and HBM-based FPGA boards

• Evaluation

- NERO outperforms a 16-core IBM POWER9 system by 4.2x and 8.3x when running two compound stencil kernels
- NERO reduces energy consumption by 22x and 29x with an energy efficiency of 1.5 GFLOPS/Watt and 17.3 GFLOPS/Watt

NERO:

A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling

<u>Gagandeep Singh</u>, Dionysios Diamantopoulos, Christoph Hagleitner, Juan Gómez-Luna, Sander Stuijk, Onur Mutlu, and Henk Corporaal

> 30th FPL, Sweden 31th August 2020



