MAD24-403



Optimising oneDAL & oneDNN for Arm CPUs to accelerate Al workloads

Ragesh Hajela

Software Engineering Manager MONAKA Software R&D (HPC AI) Unit Fujitsu Research of India, Bengaluru

Leading AI software R&D at Fujitsu Research of India, Bengaluru, About Me Introduction focussing on Arm HPC Framework Engineering with Machine Learning based Computational Data Science. 2023 Software Engineering Manager, Fujitsu Research of India Leading ML Framework Engineering at FUJITSU-MONAKA SW R&D Unit Working with Advanced Technology Development Unit (ATDU, Japan) Nexus 20 2021 AI Frameworks Engineer at Intel Corporation OpenVINO Toolkit Inference Acceleration Software Development 2019 Principal Machine Learning Engineer, ConcertAI AI Software Engineering for Healthcare AI Oncology Research 2013 Member of Technical Staff, Amadeus Labs ML Software development for Travel Intelligence Platform **Ragesh Hajela** Software Engineer, Tata Consultancy Services Software Engineering Manager 2011 C++ Programmer for Software Application Development Fujitsu Research of India Pvt Limited (FRIPL) Email ID: ragesh.hajela@fujitsu.com

About Fujitsu Our Team at FRIPL Fujitsu

- Expertise in AI Software
 Framework Engineering and
 Optimization with ML, DL,
 NLP, generative AI, LLMs,
 Real-time data processing &
 Cloud-based Data Security.
- Open-source contributors for accelerated software in HPC AI core technologies:

Scikit Learn, BLAS, Math Libs, PyTorch, TensorFlow, Hugging Face, LLMs etc.

Fujitsu Research Locations







Outline of Presentation

- □ Introduction about AI workloads optimization
- □ Fujitsu's contribution for Scikit Learn oneDAL porting on Arm
- □ Fujitsu's recent advancements in oneDNN and JIT Kernel on Arm
- Conclusion and Future Work
- Resources and Acknowledgement



Project Vision



Fujitsu's vision is to develop an energy-efficient AI software stack to enable maximised Arm CPU performance and support green data centres for sustainable digital transformation



Al Workloads in Data Centers

In data center operations, 80% workloads are machine learning applications. Our focus is on AI software acceleration to contribute towards low-power green data centres with Arm-based HPC technologies



Source: Statista Research Report (DC Market Size Analysis)





Fujitsu's contribution for Scikit-Learn oneDAL Porting on Arm

Machine Learning on Arm



Introduction to Scikit-Learn and oneDAL

Machine Learning Workload Optimization with Scikit-Learn Framework Enablement & Tuning



Intel[®] Extension for scikit-learn (scikit-learn-intelex) speeds up scikit-learn applications for x86 CPUs and GPUs across single- and multi-node configurations using oneDAL (Data Analytics Library)



https://github.com/intel/scikit-learn-intelex



Fujitsu's key contributions to OSS Community



Fujitsu ported oneDAL on Arm continuing long history of collaborating with open-source communities, via open-source development in mission-critical systems.

oneDAL PR (#2614) merged, raised by Fujitsu, to enable oneDAL multi architecture build with reference backend selection on Arm and optimized Scikit-Learn Algorithms



<> Code -

+1.220 -253

Partnership with Unified Accelerator (UXL) Foundation



Build a multi-architecture multi-vendor software ecosystem for all accelerators

- Unify the heterogeneous compute ecosystem around open standards
- Build on and expand open-source projects for accelerated computing

Steering Committee Members



oneDAL Arm Porting Design

Historically, Intel's oneAPI Data Analytics Library (oneDAL) could only be compiled on x86 architecture due to Intel's Math Kernel Library (MKL) binary-only backend.

To accelerate ML workloads on Arm, Fujitsu replaced MKL calls with open-source function calls, and this resulted in oneDAL enablement on Arm.

It is one of the first open-source contributions to UXL Foundation.

Scikit-Learn-Intelex

oneDAL Orm

Performance Results

With SVE optimisation and oneDAL porting enhancements on ARM, our work showcases notable performance gains across multiple ML algorithms. Random Forest Training Speedups

Logistic Regression Training Speedups

Results computed on AWS Graviton3 Arm-based CPU c7g.8xlarge 32-cores

These graphs illustrate the training speedup of top two ML algorithms used by Fujitsu AutoML, which got a significant speedup of x31 in Random Forest and x40 in Logistic Regression.

🗞 Linaro Connect

Porting Methodology with Arm SVE

MKL_DAAL

BLAS

SPBLAS

LAPACK

UTILITIES

THREADING + DFT MATH Kernels OSS

Replacement & Tuning

+

RNG

VSL

oneDAL on x86 uses MKLFPK, with functionalities

- To support these functions on Arm, open-source optimised compute kernels from OpenBLAS are used as alternatives to leverage SVE on the Arm.
- Used reference backend & added makefile compiler options
- Added compiler macros throughout the code base to isolate x86 specific code chunks and handle it with Arm.

Linaro Connect

Performance for SGEMM improved by ~2-5% and DGEMM improved by ~2-12%

<> Code -

+4 -4

OpenBLAS GEMM Tuning PR (#4381) is merged. Updated SGEMM and DGEMM PQ param as per config of NEOVERSEV1 cache size

Update GEMM param for NEOVERSEV1 #4381

⊱ Merged 🔰 martin-frbg merged 1 commit into OpenMathLib:develop from darshanp4:issue_4323 [口 on Dec 19, 2023]

₽ Conversation 3 E Checks 63 -O- Commits 1

(1) Files changed (1)

Source Code

Performance Result OpenBLAS/gemm.c benchmark on AWS Graviton3 single core SGEMM & DGEMM MFlops 30000 3000 5000 2000 4000 6000 1000 sgemm dgemm dgemm peak sgemm peak sgemm 240 640 dgemm 240 320 Matrix Size

OpenBLAS GEMM 1-core Tuning

OpenBLAS GEMM Tuning PR (#4629) is merged. GEMM_PREFERED_SIZE is adjusted to 4 for double precision and 8 for single precision for Neoverse V1

Set GEMM_PREFERED_SIZE parameter for Neoverse V1 #4629

 Merged
 martin-frbg
 merged 1
 commit into
 OpenMathLib:develop
 from
 tetsuzo-usui:Pf5izeTune_forNeoverseV1
 □
 2
 weeks ago

 □
 ○
 Commits
 1
 □
 Checks
 69
 ①
 Files changed
 1

GEMM_PREFERED_SIZE is modified small enough not to cause load imbalance and large enough to improve kernel efficiency

Source Code

Performance Result

<> Code -

+2 -0

OpenBLAS GEMM multi-core Tuning

OpenBLAS GEMM multi-core tuning PR (#4655) Enhanced 2D thread distribution for improved parallel performance, noticeable with large matrices.

Expanding the scope of 2D thread distribution to improve multi-threaded DGEMM performance #4655 (* Merged martin-frbg merged 1 commit into OperMathLib:develop from yamazakimitsufumi:update_2d_thread_distribution [] 17 hours ago (Conversation 2) Commits 1 Checks 69 Files changed 1 +10 -0

Performance improved by about 10% on Graviton3E (64 cores) & more than 20% on Xeon Platinum 8375C (32 cores x 2 sockets)

Source Code

driver/level3/level3_thread.c 🖸				
@@ -826,6 +826,16 @@ int CNAME(blas_arg_t *args, BLASLONG *range_m, BLASLONG *range_n, IFLOAT *sa, IF				
		<pre>if (nthreads_m * nthreads_n > args -> nthreads) {</pre>		
		<pre>nthreads_n = blas_quickdivide(args -> nthreads, nthreads_m);</pre>		
		}		
		/* The nthreads_m and nthreads_n are adjusted so that the submatrix	*/	
		/* to be handled by each thread preferably becomes a square matrix	*/	
		<pre>/* by minimizing an objective function 'n * nthreads_m + m * nthreads_n'.</pre>	*/	
		/* Objective function come from sum of partitions in m and n.	*/	
		/* (n / nthreads_n) + (m / nthreads_m)	*/	
		/* = (n * nthreads_m + m * nthreads_n) / (nthreads_n * nthreads_m)	*/	
		<pre>while (nthreads_m % 2 == 0 && n * nthreads_m + m * nthreads_n > n * (nthr</pre>	eads_m / 2) + m * (nthreads_n * 2	
		nthreads_m /= 2;		
		nthreads_n *= 2;		
		}		

Performance Result

OpenBLAS Threading Improvements

OpenBLAS Math Library Orm Multithreading pthreads OpenMP

+62 -20

OpenBLAS OpenMP Tuning PR (#4503) is merged. Implemented OpenMP locking mechanism for refined parallel execution

OpenMP locks instead of busy-waiting with NUM_PARALLEL #4503

 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □
 □

Locking design as per OpenMP specifications, consistent with Pthreads, Win32. Performance gain up to ~15% to default.

Solution Design

No of GEMM calls=100 Entry part of thread Exit part of thread 4% 12.29 Existing With changes 11.77 Parallel gemm calls from user Continue exec.. Critical Section lock Master lock ~ 14% Parallel Region == 0 ? Critical Section lock 6.3 No Parallel Region == 1 ? 5.4 Increment count Increment count Decrement count Decrement count Don't Release Master lock Release Master lock Release Master lock Don't Release Master lock Release Critical Section lock Release Critical Section lock NUM PARALLEL=6 NUM PARALLEL=1 Continue exec.. Return 0 Execution Time(sec)

Performance Result

OpenBLAS Thread Composability

OpenBLAS Thread Composability PR (#4577) is merged. Enables OpenBLAS threading to be composable with caller multithreading with performance gain

Introduced callback to Pthread, Win32 and OpenMP backend #4577

Achieves thread composability in OpenBLAS with all threading backends e.g., OpenMP, Pthreads and win32

Solution Design

Performance Result

Upto 2.7x speedup with TBB v/s OpenMP in nested parallelism scenario

<> Code

+345 -213

oneDAL Multi Architecture Enablement

Contributions

Future Work

Block Size Optimization for Arm	Bazel Build Support for Arm	Algorithm and Kernel Tuning
Dynamic template dispatcher to identify architecture/ISA specific optimal block size	New architectures to support Bazel build system, starting with Arm	Further tuning of Machine Learning algorithms and computationally intensive Math Kernels

Fujitsu's recent advancements in oneDNN and JIT Kernel on Arm

Deep Learning on Arm

Introduction – oneDNN and JIT

- To accelerate deep learning (DL) processes on the supercomputer Arm CPU's, oneDNN was ported and optimized by Fujitsu, Kawakami – San & this work was presented in Linaro Connect 2021.
- oneDNN is an open-source DL processing library developed by Intel for crossplatform architecture. oneDNN dynamically creates the execution code for the computation kernels, which are implemented at the granularity of arch. instructions using Xbyak, the Just-In-Time (JIT) assembler.
- Just-In-Time (JIT) compilation is a technique used in compiler design where the compiler translates source code into machine code at runtime, rather than ahead of time (AOT) as in traditional compilation.
- Major Advantages of JIT:

Linaro Connect

BRGEMM Contribution and Results

BRGEMM MatMul Enablement PR (#1818) is merged, expands Arm SVE support for matrix mul. & adds BRGEMM folder for aarch64

cpu: aarch64: Expand ARM SVE support for matrix multiplication #1818

Merged) igorsafo merged 12 commits into oneapi-src:main from vineelabhinav:feature-sve-matmul 🖓 2 weeks ago

~9x performance gain observed as compared to oneDNN GEMM JIT Implementation

BREGMM MatMul Performance Speed Up

U What is BRGEMM ? (Batch Reduced General Matrix Multiplication)

- Tensors are reduced into batches for multiplications
- Broadcast the input matrix B values
- Perform fused-multiply-add instruction (fma instruction) at once for multiple values

SoftMax with results

oneDNN Pooling JIT Kernel PR (#1786) is merged, expands support of SoftMax for multiple ISA

cpu: aarch64: Expand ARM SVE support in jit_uni_softmax #1786

⊱ Merged 🛛 dzarukin merged 1 commit into oneapi-src:main from deepeshfujitsu:aarch64-sve-jit-softmax 🖵 on Feb 6

~100x performance gain observed as compared to oneDNN Reference implementation

Conditional Statement Modification : Added OR condition to use multiple ISA Added new function for supporting SVE in different vector length.

Performance Result

Pooling with results

oneDNN Pooling JIT Kernel PR (#1850) is merged, expands support of Pooling for multiple ISA

cpu: aarch64: Expand ARM SVE support in jit_uni_pool_kernel #1850

⊱ Merged 🔰 igorsafo merged 2 commits into oneapi-src:main from vishwascm:aarch64-sve-jit-pooling 🖵 2 weeks ago

~4x performance gain observed as compared to current implementation in oneDNN

Code Changes

Conditional Statement Modification : Added OR condition to use multiple ISA. Modified load and store instructions to use predicate registers for correct ISA matching.

Performance Result

Linaro Connect

Concluding Remarks

- Creating a new era of computing power is mandatory for the future society with massive data generation and processing
- Ever-increasing power in datacenters is critical, and the power efficiency in CPU (consists of 60%) would be the vital factor for a sustainable future
- Fujitsu shall utilize its Supercomputer success and technology for the solution

- Developing the new power efficient CPU "FUJITSU-MONAKA" for datacenters, which will be shipped in 2027
- Targeted for wide range of usage in the datacenter including AI and HPC, and contribute to the realization of carbon-neutral society

Software Ecosystem for AI & HPC Computing

Conclusion and Future Work

Resources

FUJITSU-MONAKA Reference Links

- <u>FUJITSU-MONAKA Next</u>
 <u>Arm Processor</u>
- <u>Democratizing the use of</u> <u>AI: FUJITSU - MONAKA</u>

Scan the QR code to know more

oneDAL Pull Request Contributions

- Enable ARM(SVE) CPU support with reference backend #2614
- Enable build on ARM(SVE) #1771
- <u>Makefile refactoring to</u> <u>factor out common build</u> <u>code #2672</u>
- <u>Enable Cross Compilation</u> of Arm SVE on x86 Cl (Binary only) #2691

OpenBLAS Pull Request Contributions

- Update GEMM param for NEOVERSEV1 #4381
- <u>Set GEMM PREFERED SIZE</u> parameter for Neoverse V1 #4629
- Expanding the scope of 2D thread distribution to improve multi-threaded DGEMM performance #4655
- OpenMP locks instead of busy-waiting with NUM_PARALLEL #4503
- Introduced callback to <u>Pthread, Win32 and</u> OpenMP backend #4577

oneDNN Pull Request Contributions

- <u>cpu: aarch64: Expand ARM</u> <u>SVE support for matrix</u> <u>multiplication #1818</u>
- <u>cpu: aarch64: Expand ARM</u> <u>SVE support in</u> jit uni softmax #1786
- <u>cpu: aarch64: Expand ARM</u> <u>SVE support in</u> jit uni pool kernel #1850

Acknowledgement

- Kentaro KAWAKAMI, Kouji KURIHARA, Takumi HONDA : <u>Binary Translator to</u> <u>Accelerate Development of Deep Learning Processing Library for AArch64 CPU</u>
- Mitsufumi YAMAZAKI : <u>Update GEMM param for NEOVERSEV1 #4381</u>
- □ Tetsuzou USUI : <u>Set GEMM_PREFERED_SIZE parameter for Neoverse V1 #4629</u>
- Mitsufumi YAMAZAKI : <u>Expanding the scope of 2D thread distribution to</u> <u>improve multi-threaded DGEMM performance #4655</u>
- Initial input for backend selection #2396

Acknowledgement

NEDO Project | "Technology Development of the Next Generation Green Data Center" for the "Green Innovation Fund Project/Construction of Next Generation Digital Infrastructure"

- NEDO is "New Energy and Industrial Technology Development Organization", a national research and development agency in Japan.
- Fujitsu has been selected for the national initiative along with NEC Corporation, AIOCORE Co., Ltd., KIOXIA Corporation, FUJITSU Optical Components Limited and KYOCERA Corporation.
- This presentation is based on results obtained from a project, JPNP21029 subsidized by the New Energy and Industrial Technology Development Organization (NEDO).

Thank you