

# Matrix-Multiply Assist Best Practices Guide

Puneeth Bhat José Moreira Satish Kumar Sadasivam



# **Power Systems**







IBM Redbooks

# Matrix-Multiply Assist Best Practices Guide

April 2021

Note: Before using this information and the product it supports, read the information in "Notices" on page v.

#### First Edition (April 2021)

This edition applies to the Matrix-Multiply Assist (MMA) architecture introduced in Power ISA Version 3.1.

#### © Copyright International Business Machines Corporation 2021. All rights reserved.

Note to U.S. Government Users Restricted Rights -- Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

# Contents

Notices	.v vi
Preface         Authors         Now you can become a published author, too!         Comments welcome         Stay connected to IBM Redbooks	vii vii viii viii viii
Chapter 1. Matrix multiplication.         1.1 Matrix-multiply operation.         1.2 Vector outer product operation         1.3 Introduction to Vector Scalar Extension         1.4 Simple VSX code example for a vector outer product	1 2 4 6 7
Chapter 2. Matrix-Multiply Assist Architecture         2.1 Data types.         2.2 Data layout in accumulators and VSRs.         2.3 Instructions         2.3.1 Accumulator operation instructions.         2.3.2 Outer product instructions.         2.3.3 Advanced feature: Lane masking	11 12 13 13 14 16
Chapter 3. Programming with Matrix-Multiply Assist         3.1 Single-precision GEMM using MMA         3.2 Double-precision GEMM using MMA and one accumulator         3.3 Mixed and lower precision matrix multiplication with MMA         3.3.1 Source matrix reordering with Int8 as example	19 20 21 24 24
Chapter 4. Advanced programming concepts4.1 Multiple accumulators SGEMM for load value reuse4.2 Multiple accumulators DGEMM for load value reuse4.3 SGEMM performance with advanced cache-blocking.	29 30 33 35
Chapter 5. Matrix-Multiply Assist programming with compiler built-ins           5.1 Simple MMA SGEMM example using built-ins	37 38
Appendix A. List of Matrix-Multiply Assist compiler built-ins	41
Appendix B. List of Matrix-Multiply Assist instructions in Power ISA v3.1 Related publications Online resources	43 45 45 45

# Notices

This information was developed for products and services offered in the US. This material might be available from IBM in other languages. However, you may be required to own a copy of the product or product version in that language in order to access it.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing, IBM Corporation, North Castle Drive, MD-NC119, Armonk, NY 10504-1785, US

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you provide in any way it believes appropriate without incurring any obligation to you.

The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on specific configurations and operating conditions.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to actual people or business enterprises is entirely coincidental.

#### COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided "AS IS", without warranty of any kind. IBM shall not be liable for any damages arising out of your use of the sample programs.

# **Trademarks**

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the web at "Copyright and trademark information" at http://www.ibm.com/legal/copytrade.shtml

The following terms are trademarks or registered trademarks of International Business Machines Corporation, and might also be trademarks or registered trademarks in other countries.

Redbooks (logo) 🧬 🛽	IBM Research®	POWER7®
IBM®	POWER®	Redbooks®

The following terms are trademarks of other companies:

Other company, product, or service names may be trademarks or service marks of others.

# Preface

This publication is for software developers who want to understand the Matrix-Mulitply Assist (MMA) function, particularly those writing libraries for high performance computing and machine learning applications.

# Authors

This paper was produced by the following team of specialists.

**Puneeth Bhat A H** is a Senior Performance Analyst working on Power Systems performance at IBM. Presently he is driving the cognitive workload and interpreter performance innovations for Power processors. Puneeth, with 10 years of IBM experience, has expertise in the areas of processor micro-architecture performance, compiler optimizations and performance instrumentation tools development. He holds a Bachelor degree from Visvesvaraya Technological University.

**José E. Moreira** is a Distinguished Research Staff Member in the Scalable Systems Department at the Thomas J. Watson Research Center. He received a B.S. degree in physics and B.S. and M.S. degrees in electrical engineering from the University of Sao Paulo, Brazil, in 1987, 1988, and 1990, respectively. He also received a Ph.D. degree in electrical engineering from the University of Illinois at Urbana-Champaign in 1995. Since joining IBM at the Thomas J. Watson Research Center, he has worked on a variety of high-performance computing projects. He was system software architect for the Blue Gene/L supercomputer and chief architect of the Commercial Scale Out project. He currently leads the IBM Research® work on the architecture of POWER® processor. He is an author or coauthor of over 100 technical papers and 15 US patents. Dr. Moreira is a Fellow of the IEEE (Institute of Electrical and Electronics Engineers) and a Distinguished Scientist of the ACM (Association for Computing Machinery).

**Satish Kumar Sadasivam** is a Senior Performance Architect who leads the workload characterization and future architecture design space exploration team in IBM. He currently focuses on exploring architectural and microarchitectural innovations for Enterprise AI and cloud centric workloads with primary focus towards optimizing the core for single thread, core throughput and compute performance. He has worked on IBM POWER processor architecture since POWER5+ timeline. He is one of the key contributors for the POWER10 processor key architecture features like Instruction Fusion, Branch prediction and Matrix Math Assist (MMA). He has more than 16 years of experience in the areas of workload characterization, microarchitecture design exploration, compiler code generation and optimization, competitive evaluation, post-silicon hardware bring-up and validation. He received an MS in computer science from the Madras Institute of Technology. He is an IBM master inventor. He has filed more than 25+ patents, published several papers and delivered several talks in his field of work.

### Now you can become a published author, too!

Here's an opportunity to spotlight your skills, grow your career, and become a published author—all at the same time! Join an IBM Redbooks® residency project and help write a book in your area of expertise, while honing your experience using leading-edge technologies. Your efforts will help to increase product acceptance and customer satisfaction, as you expand your network of technical contacts and relationships. Residencies run from two to six weeks in length, and you can participate either in person or as a remote resident working from your home base.

Find out more about the residency program, browse the residency index, and apply online at:

ibm.com/redbooks/residencies.html

### **Comments welcome**

Your comments are important to us!

We want our papers to be as helpful as possible. Send us your comments about this paper or other IBM® Redbooks publications in one of the following ways:

Use the online Contact us review Redbooks form found at:

ibm.com/redbooks

Send your comments in an email to:

redbooks@us.ibm.com

Mail your comments to:

IBM Corporation, IBM Redbooks Dept. HYTD Mail Station P099 2455 South Road Poughkeepsie, NY 12601-5400

## Stay connected to IBM Redbooks

- Find us on Facebook: http://www.facebook.com/IBMRedbooks
- Follow us on Twitter: http://twitter.com/ibmredbooks
- ► Look for us on LinkedIn:

http://www.linkedin.com/groups?home=&gid=2130806

 Explore new Redbooks publications, residencies, and workshops with the IBM Redbooks weekly newsletter:

https://www.redbooks.ibm.com/Redbooks.nsf/subscribe?OpenForm

 Stay current on recent Redbooks publications with RSS Feeds: http://www.redbooks.ibm.com/rss.html

# 1

# **Matrix multiplication**

Matrix multiplication is one of the most widely used compute kernels in a broad set of applications in the emerging fields of machine learning and deep learning. This document briefly describes Matrix-Multiply Assist (MMA), which is a newly developed architecture concept that was first introduced in Power ISA Version 3.1.

The fundamental architecture principles are explained with detailed instruction set usage, register file management concepts, and various supporting facilities. The goal of this document is to help the user grasp the concepts behind MMA. The use of MMA is shown in various code examples, both in basic code versions and in fully optimized code.

## 1.1 Matrix-multiply operation

This section shows how a matrix multiplication is performed using a simple example. In this example, A and B are two 8x8 matrices, as shown in Figure 1-1. When you multiply matrices A and B, which are both 8x8 matrices, the resultant matrix C is also 8x8 in size.



Figure 1-1 Matrices A and B used in multiplication

To generalize this concept for any size matrix, assume the size of matrix A is MxK (M rows and K columns) and size of matrix B is KxN (K rows and N columns). A basic requirement of the matrix multiplication operation is that the number of columns in matrix A should be the same as the number of rows in matrix B. If the number of rows of matrix A and number of columns of matrix B are different, then the matrix multiplication operation cannot be performed.

Figure 1-2 shows a simple example of how the matrix multiplication operation is performed. To generate one element of the output matrix C, which is of size 8x8, each element of the first row of matrix A is multiplied with corresponding element of the first column of matrix B. The results are accumulated, as shown in Figure 1-2.

		A 1 2 3 4 5 9 10 11 12 13 17 18 19 20 21 25 26 27 28 29 33 34 35 36 37 41 42 43 44 45 49 50 51 52 53 57 58 59 60 61	6         7         8           14         15         16           22         23         24           30         31         32           38         39         40           46         47         48           54         55         56           62         63         64           8x8         8x8	aa ab ac ad ba bb bc bd ca cb cc cd da db dc dd ea eb cc ed fa fb fc fd ga gb gc gd ha hb hc hd	ae af ag ah be bf bg bh ce cf cg ch de df dg dh ee ef eg eh fe ff fg fh he hf hg gh 8x	8	
			C = A	*B			
- 1aa+2ba+3ca+4da+5ea+6 fa+7ga+8ha	lab+2bb+3cb+4db+5eb+ 6fb+7gb+8hb	lac+2bc+3cc+4dc+5ec+ 6fc+7gc+8hc	1ad+2bd+3cd+4dd+5ed+ 6fd+7gd+8hd	lae+2be+3ce+4de+5ee+ 6fe+7ge+8he	laf+2bf+3cf+4df+5ef+ 6gf+7ff+8hf	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh
9aa+10ba+11ca+12da+13	9ab+10bb+11cb+12db+1	9ac+10bc+11cc+12dc+1	9ad+10bd+11cd+12dd+1	9ae+10be+11ce+12de+1	9af+10bf+11cf+12df+1	9ag+10bg+11cg+12dg+1	9ah+10bh+11ch+12dh+1
ea+14fa+15ga+16ha	3eb+14fb+15gb+16hb	3ec+14fc+15gc+16hc	3ed+14fd+15gd+16hd	3ee+14fe+15ge+16he	3ef+14gf+15ff+16hf	3eg+14fg+15gg+16gh	3eh+14fh+15gh+16hh
17aa+18ba+117ca+20da+	17ab+18bb+117cb+20db	17ac+18bc+117cc+20dc	17ad+18bd+117cd+20dd	17ae+18be+117ce+20de	17af+18bf+117cf+20df	17ag+18bg+117cg+20dg	17ah+18bh+117ch+20dh
21ea+22fa+23ga+24ha	+21eb+22fb+23gb+24hb	+21ec+22fc+23gc+24hc	+21ed+22fd+23gd+24hd	+21ee+22fe+23ge+24he	+21ef+22gf+23ff+24hf	+21eg+22fg+23gg+24gh	+21eh+22fh+23gh+24hh
25aa+26ba+27ca+28da+2	25ab+26bb+27cb+28db+	25ac+26bc+27cc+28dc+	25ad+26bd+27cd+28dd+	25ae+26be+27ce+28de+	25af+26bf+27cf+28df+	25ag+26bg+27cg+28dg+	25ah+26bh+27ch+28dh+
9ea+30fa+31ga+32ha	29eb+30fb+31gb+32hb	29ec+30fc+31gc+32hc`	29ed+30fd+31gd+32hd	29ee+30fe+31ge+32he	29ef+30gf+31ff+32hf	29eg+30fg+31gg+32gh	29eh+30fh+31gh+32hh
33aa+34ba+35ca+36da+3	33ab+34bb+35cb+36db+	33ac+34bc+35cc+36dc+	33ad+34bd+35cd+36dd+	33ae+34be+35ce+36de+	33af+34bf+35cf+36df+	33ag+34bg+35cg+36dg+	33ah+34bh+35ch+36dh+
7ea+38fa+39ga+40ha	37eb+38fb+39gb+40hb	37ec+38fc+39gc+40hc	37ed+38fd+39gd+40hd	37ee+38fe+39ge+40he	37ef+38gf+39ff+40hf	37eg+38fg+39gg+40gh	37eh+38fh+39gh+40hh
41aa+42ba+43ca+44da+4	41ab+42bb+43cb+44db+	41ac+42bc+43cc+44dc+	41ad+42bd+43cd+44dd+	41ae+42be+43ce+44de+	41af+42bf+43cf+44df+	41ag+42bg+43cg+44dg+	41ah+42bh+43ch+44dh+
5ea+46fa+47ga+48ha	45eb+46fb+47gb+48hb	45ec+46fc+47gc+48hc	45ed+46fd+47gd+48hd	45ee+46fe+47ge+48he	45ef+46gf+47ff+48hf	45eg+46fg+47gg+48gh	45eh+46fh+47gh+48hh
49aa+50ba+51ca+52da+5	49ab+50bb+51cb+52db+	49ac+50bc+51cc+52dc+	49ad+50bd+51cd+52dd+	49ae+50be+51ce+52de+	49af+50bf+51cf+52df+	49ag+50bg+51cg+52dg+	49ah+50bh+51ch+52dh+
3ea+54fa+55ga+56ha	53eb+54fb+55gb+56hb	53ec+54fc+55gc+56hc	53ed+54fd+55gd+56hd	53ee+54fe+55ge+56he	53ef+54gf+55ff+56hf	53eg+54fg+55gg+56gh	53eh+54fh+55gh+56hh
57aa+58ba+59ca+60da+6	57ab+58bb+59cb+60db+	57ac+58bc+59cc+60dc+	57ad+58bd+59cd+60dd+	57ae+58be+59ce+60de+	57af+58bf+59cf+60df+	57ag+58bg+59cg+60dg+	57ah+58bh+59ch+60dh+
1ea+62fa+63ga+64ha	61eb+62fb+63gb+64hb	61ec+62fc+63gc+64hc	61ed+62fd+63gd+64hd	61ee+62fe+63ge+64he	61ef+62gf+63ff+64hf	61eg+62fg+63gg+64gh	61eh+62fh+63gh+64hh
							8+8

Figure 1-2 Example of an 8x8 matrix multiplication

The code used to perform this matrix multiplication is shown in Example 1-1. The code multiplies each element of  $i^{th}$  row of A with each element of  $j^{th}$  column of B. The result of each multiplication is accumulated to generate a result in the  $i^{th}$  row and  $j^{th}$  column of C.

Example 1-1 Basic matrix multiplication

```
#include <stdio.h>
#include <stdlib.h>
void printF (const char *name, float *M, int m, int n) {
        printf ("\n**** Matrix %s****\n",name);
        for (int i=0; i< m; i++) {</pre>
                printf("| ");
                for (int j=0; j< n; j++) printf("%-25.4f", *(M++));</pre>
                printf("
                         \n");
        }
        }
int main (int argc, char **argv ) {
    if (argc < 4) {
        printf("Usage: %s <M> <N> <K> \n", argv[0]);
        return -1;
    }
    const int M = atoi(argv[1]);
    const int N = atoi(argv[2]);
    const int K = atoi(argv[3]);
    printf("Running: %s M=%s N=%s K=%s \n", argv[0], argv[1], argv[2], argv[3]);
    float A[M][K];
    float B[K][N];
    float C[M][N];
    for (int i=0; i<M; i++) for (int j=0; j<N; j++) C[i][j] = 0;
    int x = 1;
    for (int i=0; i<M; i++) for (int j=0; j<K; j++) A[i][j] = float(x++) * 7 / 15;
    for (int i=0; i<K; i++) for (int j=0; j<N; j++) B[i][j] = float(x++) * 3 / 17;
    for (int i=0; i<M; i++) {</pre>
        for (int j=0; j<N; j++) {</pre>
           for (int k=0; k<K; k++)</pre>
                C[i][j] += A[i][k] * B[k][j];
        }
    }
    printF("C", (float *)C, M, N);
    return 0;
}
```

## 1.2 Vector outer product operation

In Figure 1-3, the code cannot be vectorized and it is not optimal. To optimize the matrix multiplication operation the vector outer product is performed. Matrix A is transposed and then the computation is performed in a blocked manner, as follows:

- The first 8x4 block of transposed matrix A (A<sup>T</sup>) is iterated over the two 8x4 blocks of matrix B, computing outer products of the corresponding rows from blocks of A<sup>T</sup> and B, to generate two 4x4 results.
- The second 8x4 block of matrix A transposed is iterated over the same two blocks of matrix B to generate the next two 4x4 results. This operation is explained in Figure 1-3.

		Α	Т		В	~	
		1       9       17       25       3         2       10       18       26       3         3       11       19       27       3         4       12       20       28       3         5       13       21       29       3         6       14       22       30       3         7       15       23       31       3         8       16       24       32       4	3 41 49 57 4 42 50 58 5 43 51 59 6 44 52 60 7 45 53 61 8 46 54 62 9 47 55 63 0 48 56 64 8 x	aa ab a ba bb l ca cb d da db d ea eb d ga gb d ha hb l	ac ad ae af bc bd be bf cc cd ce cf dc dd de df ec ed ee ef fc fd fe ff gc gd ge gf hc hd he hf	ag ah bg bh cg ch dg dh eg eh fg fh gg gh hg hh	
1aa+2ba+3ca+4da+5ea+6			C = A	*B			
	lab+2bb+3cb+4db+5eb+	lac+2bc+3cc+4dc+5ec+	1ad+2bd+3cd+4dd+5ed+	lae+2be+3ce+4de+5ee+	laf+2bf+3cf+4df+5ef+	lag+2bg+3cg+4dg+5eg+	lah+2bh+3ch+4dh+5eh+
fa+7ga+8ha	1ab+2bb+3cb+4db+5eb+ 6fb+7gb+8hb	lac+2bc+3cc+4dc+5ec+ 6fc+7gc+8hc	1ad+2bd+3cd+4dd+5ed+ 6fd+7gd+8hd	lae+2be+3ce+4de+5ee+ 6fe+7ge+8he	laf+2bf+3cf+4df+5ef+ 6gf+7ff+8hf	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh
fa+7ga+8ha 9aa+10ba+11ca+12da+13 ea+14fa+15ga+16ha	1ab+2bb+3cb+4db+5eb+ 6fb+7gb+8hb 9ab+10bb+11cb+12db+1 3eb+14fb+15gb+16hb	lac+2bc+3cc+4dc+5ec+ 6fc+7gc+8hc 9ac+10bc+11cc+12dc+1 3ec+14fc+15gc+16hc	lad+2bd+3cd+4dd+5ed+ 6fd+7gd+8hd 9ad+10bd+11cd+12dd+1 3ed+14fd+15gd+16hd	lae+2be+3ce+4de+5ee+ 6fe+7ge+8he 9ae+10be+11ce+12de+1 3ee+14fe+15ge+16he	laf+2bf+3cf+4df+5ef+ 6gf+7ff+8hf 9af+10bf+11cf+12df+1 3ef+14gf+15ff+16hf	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh 9ag+10bg+11cg+12dg+1 3eg+14fg+15gg+16gh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh 9ah+10bh+11ch+12dh+1 3eh+14fh+15gh+16hh
fa+7ga+8ha 9aa+10ba+11ca+12da+13 ea+14fa+15ga+16ha 17aa+18ba+117ca+20da+ 21ea+22fa+23ga+24ha	1ab+2bb+3cb+4db+5eb+ 6fb+7gb+8hb 9ab+10bb+11cb+12db+1 3eb+14fb+15gb+16hb 17ab+18bb+117cb+20db +21eb+22fb+23gb+24hb	lac+2bc+3cc+4dc+5ec+ 6fc+7gc+8hc 9ac+10bc+11cc+12dc+1 3ec+14fc+15gc+16hc 17ac+18bc+117cc+20dc +21ec+22fc+23gc+24hc	lad+2bd+3cd+4dd+5ed+ 6fd+7gd+8hd 9ad+10bd+11cd+12dd+1 3ed+14fd+15gd+16hd 17ad+18bd+117cd+20dd +21ed+22fd+23gd+24hd	lae+2be+3ce+4de+5ee+ 6fe+7ge+8he 9ae+10be+11ce+12de+1 3ee+14fe+15ge+16he 17ae+18be+117ce+20de +21ee+22fe+23ge+24he	<pre>laf+2bf+3cf+4df+5ef+ 6gf+7ff+8hf 9af+10bf+11cf+12df+1 3ef+14gf+15ff+16hf 17af+18bf+117cf+20df +21ef+22gf+23ff+24hf</pre>	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh 9ag+10bg+11cg+12dg+1 3eg+14fg+15gg+16gh 17ag+18bg+117cg+20dg +21eg+22fg+23gg+24gh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh 9ah+10bh+11ch+12dh+1 3eh+14fh+15gh+16hh 17ah+18bh+117ch+20dh +21eh+22fh+23gh+24hh
fa+7ga+8ha 9aa+10ba+11ca+12da+13 ea+14fa+15ga+16ha 17aa+18ba+117ca+20da+ 21ea+22fa+23ga+24ha 25aa+26ba+27ca+28da+2 9ea+30fa+31ga+32ha	1ab+2bb+3cb+4db+5bb+ 6fb+7gb+8hb 9ab+10b+11cb+12db+1 3eb+14fb+15gb+16hb 17ab+18bb+117cb+20db +21eb+22fb+23gb+24hb 25ab+26bb+27cb+28db+ 29eb+30fb+31gb+32hb	lac+2bc+3cc+4dc+5ec+ 6fc+7gc+8bc 9ac+10bc+11cc+12dc+1 3ac+14fc+15gc+16bc 17ac+18bc+117cc+20dc +21ac+22fc+23gc+24bc 25ac+26bc+27cc+28dc+ 29ac+30fc+31gc+32bc'	1ad+2bd+3cd+4dd+5ed+ 6fd+7gd+8hd 9ad+10bd+11cd+12dd+1 3ad+14fd+15gd+16hd 17ad+18bd+117cd+20dd +21ed+22fd+23gd+24hd 25ad+26bd+27cd+28dd+ 29ed+30fd+31gd+32hd	lae+2be+3ce+4de+5ee+ 6fe+7ge+8be 9ae+10be+11ce+12de+1 3ee+14fe+15ge+18be +17ae+18be+117ce+20de +21ce+22fe+23ge+24be 25ae+26be+27ce+28de+ 29ee+30fe+31ge+32be	laf+2bf+3cf+4df+5ef+ 6gf+7ff+8bf 9af+10bf+11cf+12df+1 3ef+14gf+15ff+16hf 17af+18bf+117cf+20df +21ef+22gf+23ff+2df 25af+26bf+27cf+28df+ 29ef+30gf+31ff+32hf	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh 9ag+10bg+11cg+12dg+1 3eg+14fg+15gg+16gh 17ag+18bg+117cg+20dg +21cg+22fg+23gg+2dg 25ag+26bg+27cg+28dg+ 29eg+30fg+31gg+32gh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh 9ah+10bh+11ch+12dh+1 3eh+14fh+15gh+16hh 17ah+18bh+117ch+20dh +21eh+22fh+23gh+24hh 25ah+26bh+27ch+28dh+ 29eh+30fh+31gh+32hh
fa+7ga+8ha 9aa+10ba+11ca+12da+13 ar13ga+15ga+16ha 17aa+18ba+117ca+20da+ 21ea+22fa+23ga+24ha 25aa+26ba+27ca+28da+2 9ea+30fa+31ga+3ba 33aa+34ba+35ca+36da+3 7aa+36fa+39ga+40ha	1ab+2bb+3cb+4db+5ab+ 6fb+7dp+8bb 9ab+10bb+11cb+12db+1 3eb+14fb+15gb+16hb 17ab+18bb+117cb+20db +22tb+22db+23gb+24hb 25ab+26bb+27cb+28db+ 39ab+36b+33cb+35cb+36db+ 37eb+38fb+39gb+40hb	lac*2bc+3cc+4dc+Ssc+ 6fc+7qc+8bc 9ac+10bc+11cc+12dc+1 3ec+14fc+15qc+18hc 17ac+18bc+117cc+20dc +21ac+22fc+23qc+24hc 25ac+26bc+27cc+28dc+ 29ac+30fc+31gc+23bc <sup>2</sup> 33ac+34bc+35cc+36dc+ 33ac+34bc+35cc+36dc+	1ad+2bd+3cd+4dd+5ed+ 6fd+7gd+8hd 9ed+10bd+11cd+2dd+1 3ed+14fd+15gd+16hd 17ad+18bd+117cd+20dd +21ed+22fd+23gd+2hd 25ad+26bd+27cd+28dd+ 29ed+30fd+31gd+32hd 33ad+34bd+35cd+36dd+ 33rd+38fd+39gd+40hd	1ae+2be+3ce+4de+5ee+ 6fe+7ge+8he 9ae+10be+11ce+12de+1 3ae+14fe+15ge+16he 17ae+18be+117ce+20de +21ee+22fe+23ge+2dhe 25ae+26be+27ce+28de+ 29ee+30fe+31ge+32be 33ae+34be+35ce+36de+ 37ee+38fe+39ge+40be	laf+2Df+3of+4df+5ef+ 6gf+7ff+8hf 9af+10bf+11cf+12df+1 3af+13gf+13ff+10hf 17af+18bf+117cf+20df +21ef+22gf+23ff+23hf 25af+26bf+27cf+28df+ 29ef+30ff+31ff+32hf 33af+34bf+35cf+36df+ 37ef+38gf+39ff+40hf	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh 9ag+10bg+11cg+12dg+1 jag+14fg+15gg+16gh 17ag+18bg+117cg+20dg +21eg+22fg+23gg+2dgh 25ag+26bg+27cg+28dg+ 29ag+30fg+31gg+2gh 31ag+14bg+35cg+36dg+ 37cg+38fg+39gg+40gh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh 9ah+10bh+11ch+12dh+1 3eh+14fh+15gh+16hh 17ah+18bh+117ch+20dh +21eh+22fh+23gh+2dhh 25ah+26bh+27ch+28dh+ 29eh+30fh+31gh+32hh 37ah+3bfh+35ch+36dh+ 37eh+38fh+39gh+40hh
fa+7g+8ha 9a+10b+11c+12d+13 al+16b+11c+12d+13 17a+18b+11c+20d+ 21c+22fa+23b+24b+ 9a+10c+13g+2bha 33a+3bb+35c+15da+3 33a+3bb+35c+15da+3 41a+42b+43c+44da+4	1 ab 200 300 4 d0 5 600 4 6 fbr 7 0 5 6 8 b 9 ab 1 0 bb 1 1 cb 1 2 db 1 3 ab 1 4 fb 1 3 gb 1 4 fb 1 5 gb 1 6 bb 1 7 ab 1 8 bb 1 1 7 cb + 2 db 1 2 1 cb + 2 fb + 2 3 gb 2 db 4 2 7 cb + 2 db 4 2 9 ab 3 0 fb 3 1 gb + 3 2 bb 3 3 ab 3 4 bb + 3 5 cb 3 6 gb 4 0 bb 3 1 ab + 3 4 bb + 3 5 cb + 4 db 4 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 5 eb + 6 fb + 4 3 cb + 4 db 4 \\ 5 eb + 6 fb + 4 3 cb + 4 db 4 \\ 5 eb + 6 fb + 4 3 cb + 4 db 4 \\ 5 eb + 6 fb + 4 3 cb + 4 db 4 \\ 5 eb + 6 fb + 4 3 cb + 4 db 4 \\ 5 eb + 6 fb + 4 3 cb + 4 db 4 \\ 5 eb +	lact2bc+3cc+4dc+5sc+ 6fc+7gc+8bc 9ac+10bc+11cc+12dc+1 3ac+14fc+15gc+16bc 17ac+18bc+117cc+20dc +21ac+22fc+23gc+24bc 25ac+26bc+27cc+28dc+ 29ac+30fc+31gc+32bc' 33ac+3bc+35gc+43gc+40bc 41ac+42bc+43gc+40bc	1 ad+22bd+3cd+4dd+5ed+ 6fd+7gd+8hd 9ad+10bd+11cd+12dd+1 3ad+14fd+15gd+16hd 17ad+18bd+117cd+22dd+ 421ed+22fd+23gd+2dhd 25ad+26bd+27cd+28dd+ 25ed+30fd+31gd+32bd 33ad+3bd+35cd+35cd+ 37ed+38fd+39gd+40bd 41bd+22bd+32cd+43cd+44bd	1 20+220+320+420+580+ 6 20+730+830e 3 20+13 20+1120+1220+1 1 730+1380+1170+73006 2 2 2 2 0 + 2 2 2 2 0 + 2 2 2 0 + 2 2 2 0 + 2 2 0 + 2 2 0 + 2 2 0 + 2 2 0 + 2 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 + 2 0 +	l <u>df-2Df+3Df+df+5ef+</u> <u>6gf+7ff+8bf</u> <u>3df+10bf+11cf+12df+1</u> <u>3df+14gf+15ff+16bf</u> <u>17df+18bf+117cf+20df</u> <u>+21ef+23ff+23ff+2dhf</u> <u>25df+25ff+27cf+28df+</u> <u>33df+3bf+35cf+3ff+43bf</u> <u>37df+38gf+39ff+40bf</u> <u>41bf+42bf+43cf+44bf</u> <u>45ef+46gf+47ff+48bf</u>	lag+2bg+3cg+4dg+5eg+ 6fg+7gg+8gh           9ag+10bg+11cg+12dg+1           3ag+14fg+15gg+16gh           17ag+18bg+117cg+20dg           +21eg+22fg+23gg+2dgh           29ag+30bg+27cg+28dgh           33ag+34bg+32cg+26dgh           33ag+34bg+32cg+36dgh           37ag+38fg+39gg+40gh           41ag+42bg+43cg+44dgh           45eg+445g+43cg+44dgh	lah+2bh+3ch+4dh+5eh+ 6fh+7gh+8hh 3eh+10bh+11ch+12ch+1 3eh+14fh+15gh+16hh 17ah+18bh+117ch+20dh +21eh+22fh+23gh+2dh+ 29eh+30fh+31gh+32hh 33ah+34bh+35ch+35ch+ 37eh+38fh+39gh+40hh 4lah+42bh+43ch+4dh
fa+7g+8ha 9a+10b+11Ca+12da+13 17a+18b+117ca+20da+ 21a+22fa+23ga+24ha 9a+30fa+31ga+32ha 7aa+38fa+35ga+40ha 41aa+42ba+43ca+44da+ 9aa+50fa+55ga+56ha	1 ab 2200 300 4 40 5 500 4 50 4 70 40 8 50 4 3 ab 14 (bb + 11 cb + 12 db + 1 3 ab 14 (bb + 11 cb + 12 db + 1 3 ab 14 (bb + 12 bb + 12 bb + 1 2 1 ab + 22 fb + 2 3 gb + 2 4 bb 4 2 3 ab 2 4 fb + 27 cb + 2 3 gb + 2 4 bb 4 2 3 ab + 3 (bb + 27 cb + 2 3 gb + 2 4 bb 4 3 3 ab + 3 (bb + 32 cb + 3 4 db + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 5 3 2 bb + 5 3 bb + 5 4 2 bb + 3 2 bb + 5 2 bb + 5 3 bb + 5 4 2 bb + 3 2 bb + 5 2 bb + 5 3 2 bb + 5 3 bb + 5 4 2 bb + 3 2 bb + 5 2 bb + 5 2 bb + 5 3 bb + 5 4 2 bb + 3 2 bb + 5 2 bb	lac*2bc+3cc+4dc+Ssc+ 6fc+7qc+4kbc           9ac+10bc+11cc+12dc+1           3ac+14fc+13qc+14bc           17ac+18bc+117cc+20dc           421ac+22fc+23qc+24bc           25ac+26bc+27cc+28dc+           25ac+26bc+33qc+23bc*           33ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+34bc+35cc+36dc+           37ac+35bc+51cc+52dc+	1ad+2bd+3cd+4dd+5ed+           6fd+7gd+8hd           9ad+10bd+11cd+12dd+1           3ad+14fd+15gd+16hd           17ad+18bd+117cd+20dd           +21ed+22fd+23gd+2dhd           25ad+26bd+27cd+28dd+           29d+30fd+31gd+2bd           3ad+34bd+35cd+36dd+           3rd+34cd+32d+43cd+4dd+           45scd+46fd+47gd+48bd           41ad+2bd+32cd+36dd+           3bd+35cd+36dd+           3bd+35cd+36dd+           3bd+35cd+36dd+           3bd+35cd+36dd+           3bd+45bd+51cd+52dd+           3bd+45bd+51cd+52dd+           3bd+54fd+55gd+56bd	1 24220+326+426+56e+ 526+79e+8be 30e+1326+1126+1226+1 30e+1426+15ge+16be 17ae+182be+117ce+2026 25ae+26be+27ce+28de+ 25ae+26be+27ce+28de+ 30ae+34be+15ce+36de+ 37ae+382e+39ge+40be 41ae+42be+43ce+43de+ 45ae+562e+45ge+56be	1 af+2bf+3of+4df+Sef+ 6gf+7ff+8bf 3 af+10f+12df+1 3 af+14gf+15ff+16hf 4 21ef+22gf+33ff+2ahf 2 2sf+26f+27of+28df+ 2 2sf+30ff+35ff+3hf 3 3af+34bf+35cf+36df+ 3 7ef+38gf+39ff+40hf 4 3af+42bf+43cf+44df+ 4 45ef+46f+47ff+8bf 5 3ef+54gf+55ff+56hf	1ag+2bg+3cg+4dg+5eg+ 65g+7gg+8gh 9ag+10bg+11cg+12dg+1 1ag+145g+15gg+16gh 17ag+18bg+117cg+20dg +21eg+22fg+23gg+2dgh 25ag+26bg+27cg+28dg+ 29eg+30fg+31gg+22gh 33ag+34bg+35cg+36dg+ 37eg+38fg+39gf+40gh 41ag+42bg+42cg+44dg+ 45eg+46fg+47gg+48gh 53eg+54fg+55gg+56gh	lah+2bh+3ch+4dh+5eh+ 6fhr7gh+8hh 9ah+10bh+11ch+12dh+1 3eh+14fh+15gh+16hh 17ah+18bh+117ch+20dh +21eh+22fh+23gh+28dh+ 29eh+30fh+31gh+32hh 33ah+34bh+35ch+36dh+ 37eh+38fh+39gh+40hh 41ah+42bh+43ch+4ddh+ 45eh+46fh+31gh+88hh

Figure 1-3 Matrix multiplication through outer products

Two key benefits of this optimization are:

- Opens up the possibility of parallelizing the entire computation.
- ► Reduces the number of loads performed and improves data locality because of reuse.

Figure 1-4 explains, in detail, how a partial matrix product is generated using the vector outer product operation. Each of the elements of transposed matrix A (1, 9, 17, 25) is multiplied with each of the elements of matrix B (aa, ab, ac, ad) to generate 16 outputs, which become the partial result of the resultant 4x4 output submatrix. The same operation is performed for each of the rows in both transposed matrix A and matrix B and the subsequent results are accumulated to get the final result.



Figure 1-4 Generating a partial matrix product using a vector outer product

If this operation is repeated over all the blocks of both transposed matrix A and matrix B, then the final 8x8 results are generated, as shown in Figure 1-5. For a more detailed discussion on high-performance matrix multiplication, see Anatomy of High-Performance Matrix Multiply.

	$\begin{bmatrix} 1 & 9 & 17 & 25\\ [2 & 10 & 18 & 26\\ [3 & 11 & 19 & 27\\ [4 & 12 & 20 & 28\\ [5 & 13 & 21 & 29\\ [6 & 14 & 22 & 30\\ [7 & 15 & 23 & 31\\ [8 & 16 & 24 & 32 \end{bmatrix}$	] [ aa ab ac ad ] ] [ ba bb bc bd ] ] [ ca cb cc cd ] ] [ ca db dc dd ] ] [ aa db ec ed ] ] [ fa fb fc fd ] ] [ ga gb gc gd ] ] [ ha hb hc hd ]	$\mathbf{C} = \mathbf{A}$	A*B	$      \begin{bmatrix} 1 & 9 & 17 \\ 12 & 10 & 18 \\ 3 & 11 & 19 \\ 14 & 12 & 20 \\ 15 & 13 & 21 \\ 16 & 14 & 22 \\ 17 & 15 & 23 \\ 18 & 16 & 24 \\      \end{bmatrix} $	25] [ ae af ag ah 26] [ be bf bg bh 27] [ ce cf cg ch 28] [ de df dg dh 29] [ ee df eg eh 30] [ fe ff fg fh 31] [ ge gf gg gh 32] [ he hf hg hh	
laa+2ba+3ca+4da+5ea+6	lab+2bb+3cb+4db+5eb+	lac+2bc+3cc+4dc+5ec+	lad+2bd+3cd+4dd+5ed+	lae+2be+3ce+4de+5ee+	laf+2bf+3cf+4df+5ef+	lag+2bg+3cg+4dg+5eg+	lah+2bh+3ch+4dh+5eh+
fa+7ga+8ha	6fb+7gb+8hb	6fc+7gc+8hc	6fd+7gd+8hd	6fe+7ge+8he	6gf+7ff+8hf	6fg+7gg+8gh	6fh+7gh+8hh
9aa+10ba+11ca+12da+13	9ab+10bb+11cb+12db+1	9ac+10bc+11cc+12dc+1	9ad+10bd+11cd+12dd+1	9ae+10be+11ce+12de+1	9af+10bf+11cf+12df+1	9ag+10bg+11cg+12dg+1	9ah+10bh+11ch+12dh+1
ea+14fa+15ga+16ha	3eb+14fb+15gb+16hb	3ec+14fc+15gc+16hc	3ed+14fd+15gd+16hd	3ee+14fe+15ge+16he	3ef+14gf+15ff+16hf	3eg+14fg+15gg+16gh	3eh+14fh+15gh+16hh
17aa+18ba+117ca+20da+	17ab+18bb+117cb+20db	17ac+18bc+117cc+20dc	17ad+18bd+117cd+20dd	17ae+18be+117ce+20de	17af+18bf+117cf+20df	17ag+18bg+117cg+20dg	17ah+18bh+117ch+20dh
21ea+22fa+23ga+24ha	+21eb+22fb+23gb+24hb	+21ec+22fc+23gc+24hc	+21ed+22fd+23gd+24hd	+21ee+22fe+23ge+24he	+21ef+22gf+23ff+24hf	+21eg+22fg+23gg+24gh	+21eh+22fh+23gh+24hh
25aa+26ba+27ca+28da+2	25ab+26bb+27cb+28db+	25ac+26bc+27cc+28dc+	25ad+26bd+27cd+28dd+	25ae+26be+27ce+28de+	25af+26bf+27cf+28df+	25ag+26bg+27cg+28dg+	25ah+26bh+27ch+28dh+
9ea+30fa+31ga+32ha	29eb+30fb+31gb+32hb	29ec+30fc+31gc+32hc	29ed+30fd+31gd+32hd	29ee+30fe+31ge+32he	29ef+30gf+31ff+32hf	29eg+30fg+31gg+32gh	29eh+30fh+31gh+32hh
33aa+34ba+35ca+36da+3	33ab+34bb+35cb+36db+	33ac+34bc+35cc+36dc+	33ad+34bd+35cd+36dd+	33ae+34be+35ce+36de+	33af+34bf+35cf+36df+	33ag+34bg+35cg+36dg+	33ah+34bh+35ch+36dh+
7ea+38fa+39ga+40ha	37eb+38fb+39gb+40hb	37ec+38fc+39gc+40hc	37ed+38fd+39gd+40hd	37ee+38fe+39ge+40he	37ef+38gf+39ff+40hf	37eg+38fg+39gg+40gh	37eh+38fh+39gh+40hh
41aa+42ba+43ca+44da+4	41ab+42bb+43cb+44db+	41ac+42bc+43cc+44dc+	41ad+42bd+43cd+44dd+	41ae+42be+43ce+44de+	41af+42bf+43cf+44df+	41ag+42bg+43cg+44dg+	41ah+42bh+43ch+44dh+
5ea+46fa+47ga+48ha	45eb+46fb+47gb+48hb	45ec+46fc+47gc+48hc	45ed+46fd+47gd+48hd	45ee+46fe+47ge+48he	45ef+46gf+47ff+48hf	45eg+46fg+47gg+48gh	45eh+46fh+47gh+48hh
49aa+50ba+51ca+52da+5	49ab+50bb+51cb+52db+	49ac+50bc+51cc+52dc+	49ad+50bd+51cd+52dd+	49ae+50be+51ce+52de+	49af+50bf+51cf+52df+	49ag+50bg+51cg+52dg+	49ah+50bh+51ch+52dh+
3ea+54fa+55ga+56ha	53eb+54fb+55gb+56hb	53ec+54fc+55gc+56hc	53ed+54fd+55gd+56hd	53ee+54fe+55ge+56he	53ef+54gf+55ff+56hf	53eg+54fg+55gg+56gh	53eh+54fh+55gh+56hh
57aa+58ba+59ca+60da+6	57ab+58bb+59cb+60db+	57ac+58bc+59cc+60dc+	57ad+58bd+59cd+60dd+	57ae+58be+59ce+60de+	57af+58bf+59cf+60df+	57ag+58bg+59cg+60dg+	57ah+58bh+59ch+60dh+
1ea+62fa+63ga+64ha	61eb+62fb+63gb+64hb	61ec+62fc+63gc+64hc	61ed+62fd+63gd+64hd	61ee+62fe+63ge+64he	61ef+62gf+63ff+64hf	61eg+62fg+63gg+64gh	61eh+62fh+63gh+64hh
	$\begin{bmatrix} 33 & 41 & 49 & 57 \\ 34 & 42 & 50 & 56 \\ 155 & 43 & 51 & 55 \\ 36 & 44 & 52 & 66 \\ 137 & 45 & 53 & 61 \\ 138 & 46 & 54 & 65 \\ 139 & 47 & 55 & 63 \\ 140 & 48 & 56 & 64 \end{bmatrix}$	7]       [ aa ab ac ad ]         8]       [ ba bb bc bd ]         9]       [ ca cb cc cd ]         0]       [ da db dc dd ]         1]       [ ea db ec ed ]         2)       [ fa fb fc fd ]         3)       [ a db hc ed ]			$\begin{bmatrix} 33 & 41 & 49 \\ [34 & 42 & 50 \\ [35 & 43 & 51 \\ [36 & 44 & 52 \\ [37 & 45 & 53 \\ [38 & 46 & 54 \\ [39 & 47 & 55 \\ [40 & 48 & 56 \end{bmatrix}$	57) [ ae af ag af 58] [ be bf bg bf 59) [ ce cf cg cf 60] [ de df dg df 61] [ ee df eg ef 62] [ fe ff fg ff 63] [ ge gf gg gf 64] [ he hf hg hf	

Figure 1-5 Full 8x8 matrix-multiply from four 4x4 blocks

# **1.3 Introduction to Vector Scalar Extension**

Vector Scalar Extension (VSX) is the vector scalar extension capability introduced in the POWER ISA V2.06 and first implemented in the IBM POWER7® processor. In Figure 1-6, the VSX capability extends and unifies the 32 floating-point registers (FPRs) of Power ISA to 128 bits and combines with the existing 32 128-bit Vector Multimedia Extension (VMX) registers to create a single register file.



Figure 1-6 VSX registers

The VSX capability allows the VSX instruction to utilize 64 x 128-bit registers to perform its compute operations. VSX ISA supports numerous operations in both floating point and fixed point.

The key instructions that are useful in demonstrating the outer-product operations are:

- Vector loads and stores
- Splat instructions to replicate one element of the source vector register to all fields of the target vector register
- Vector floating point multiply-add instruction

### 1.4 Simple VSX code example for a vector outer product

Example 1-2 shows code that initializes matrices A, B, and C and gets a transform of matrix A.

Example 1-2 VSX Code example for a vector outer product

```
#include <stdio.h>
#include <stdlib.h>
#define KM 4
#define KN 4
extern "C" void sgemm_kernel_4x4(float*,float*,float*,int,int,int,int);
void sgemm(float *A, float *B, float *C, int M, int N, int K) {
    for (int i=0; i<M; i+=KM) {</pre>
    for (int j=0; j<N; j+=KN) {</pre>
        sgemm kernel 4x4(A+i, B+j, C+j, K, M, N, N);
       }
    C += N*KM;
    }
}
void printF (const char *name, float *M, int m, int n) {
        printf ("\n**** Matrix %s****\n",name);
        for (int i=0; i< m; i++) {
                printf("| ");
                for (int j=0; j< n; j++) printf("%-25.4f", *(M++));</pre>
                printf("
                          \n");
        }
        }
int main (int
                 argc, char **argv ) {
    if (argc < 4) {
        printf("Usage: %s <M> <N> <K> \n", argv[0]);
        return -1;
   }
    const int M = atoi(argv[1]);
    const int N = atoi(argv[2]);
    const int K = atoi(argv[3]);
    printf("Running: %s M=%s N=%s K=%s \n", argv[0], argv[1], argv[2], argv[3]);
    float A[M][K];
    float AT[K][M];
    float B[K][N];
    float C[M][N];
    for (int i=0; i<M; i++) for (int j=0; j<N; j++) C[i][j] = 0;
    int x = 1;
    for (int i=0; i<M; i++) for (int j=0; j<K; j++) A[i][j] = float(x++) * 7 / 15;
    for (int i=0; i<K; i++) for (int j=0; j<N; j++) B[i][j] = float(x++) * 3 / 17;
    for (int i=0; i<M; i++) for (int j=0; j<K; j++) AT[j][i] = A[i][j];</pre>
    sgemm((float*)AT, (float*)B, (float*)C, M, N, K);
    printF("C", (float *)C, M, N);
    return 0;
```

The **Sgemm** routine computes the result matrix C in blocks of 4x4 (Power ISA VSX register size is 128 bits and contains 4x32-bit floating-point (fp32) elements.) The **sgemm\_kerne1\_4x4** routine computes one 4x4 (KM x KN) block of the resultant matrix C.

The innermost loop of the kernel shown in Example 1-3 loads four elements from matrix A and four elements of matrix B and performs multiply-accumulate operations. This operation is performed 'K' times to compute one full 4x4 block of the resultant matrix C.

Example 1-3 Multiply-accumulate operation

```
".text"
.section
        .global sgemm kernel 4x4
        .type
                sgemm kernel 4x4, @function
sgemm kernel 4x4:
/* adjust lda, ldb, ldc for vector size 4 */
        slwi
                7, 7, 2
        slwi
                8, 8, 2
        slwi
                9, 9, 2
/* Reset VSX registers */
    xx1xor 0, 0, 0
   xxlxor 1, 1, 1
   xxlxor 2, 2, 2
    xx1xor 3, 3, 3
/* LOOP for K to 0 */
    K LOOP:
/* Load 4 elements of A, B */
              32, 0(3)
       1xv
              33, 0(4)
       1xv
/* Copy each A[i] 4 times */
       xxspltw
                  34, 32, 3
       xxspltw
                  35, 32, 2
                  36, 32, 1
       xxspltw
       xxspltw
                  37, 32, 0
/* Multiply-Add-Accumulate */
        xvmaddasp 0, 34, 33
        xvmaddasp 1, 35, 33
        xvmaddasp 2, 36, 33
        xvmaddasp 3, 37, 33
/* Update Loop count & A,B */
        add
                3, 3, 7
        add
                4, 4, 8
        addic. 6, 6, -1
           K LOOP
   bgt
/* Offsets of 4x4 C Matrix */
        slwi
                3, 9, 1
                4, 5, 9
        add
                6, 5, 3
        add
        add
                7, 4, 3
/* Store the 4x4 c Matrix */
            0, 0(5)
    stxv
            1, 0(4)
    stxv
            2, 0(6)
    stxv
    stxv
            3, 0(7)
blr
```

The following command demonstrates how to build this example with the main C file linked with a defined assembly function using the latest MMA-supported GNU Compiler Collection (GCC):

>> g++ -mcpu=power10 -02 sgemm\_4x4.cc sgemm\_vsx\_kernel.s -o sgemm\_vsx

# 2

# Matrix-Multiply Assist Architecture

The Matrix-Multiply Assist (MMA) architecture is introduced in Power ISA v3.1. Several new concepts are described in this chapter, such as:

- An introduction of accumulator registers
- New compute instructions for the matrix multiplication operation
- Support for lower precision arithmetic beyond single- and double-precision

# 2.1 Data types

MMA architecture supports both floating-point and integer data types. This support was introduced because of the growing requirements of future AI-inferencing models.

The following data types are floating-point:

- FP32 (IEEE single-precision)<sup>1</sup>
- FP64 (IEEE double-precision)<sup>2</sup>
- ► FP16 (IEEE half-precision)<sup>3</sup>
- bfloat16<sup>4</sup>

The following data types are integer:

- INT16 (16-bit integer)
- INT8 (8-bit integer)
- ► INT4 (4-bit integer)

## 2.2 Data layout in accumulators and VSRs

One of the key concepts of the MMA architecture is the accumulator register. There are eight such registers in Power ISA v3.1 and each accumulator is 512 bits. Currently, a clear association exists between accumulators and VSRs. Each VSR is 128 bits and four such VSRs are combined and shadowed to form one accumulator register. The first 32 VSRs are mapped to eight accumulator registers as shown in Figure 2-1 on page 12.

When programming with MMA instructions, one of the first design decisions is how to partition the space of VSRs and accumulators. It is important to keep the two separated. That is, if accumulator ACCx is being used (where x is from 0 - 7) from the accumulators. See 2.3, "Instructions" on page 13 for information on how to use instructions to transfer data.

VSR[0]  VSR[3]	ACC[0]
VSR[4]  VSR[7]	ACC[1]
VSR[8]  VSR[11]	ACC[2]
VSR[12]  VSR[15]	ACC[3]
VSR[16]  VSR[19]	ACC[4]
VSR[20]  VSR[23]	ACC[5]
VSR[24]  VSR[27]	ACC[6]
USR[20]	ACC[7]

Figure 2-1 Accumulator architecture in POWER ISA v3.1

<sup>&</sup>lt;sup>1</sup> 754-2019 - IEEE Standard for Floating-Point Arithmetic, found at:

https://ieeexplore.ieee.org/document/8766229

<sup>&</sup>lt;sup>2</sup> Ibid.

<sup>&</sup>lt;sup>3</sup> Ibid.

<sup>&</sup>lt;sup>4</sup> A transprecision floating-point platform for ultra-low power computing, found at: https://ieeexplore.ieee.org/abstract/document/8342167

**Important:** Do not mix instructions that use an accumulator (for example, ACC0) and the corresponding VSRs (for example, VSR[0:3]). You can mix instructions that use an accumulator (for example, ACC0) and VSRs that do not overlap it (for example, VSR[4:7]). VSR[32:63] should never overlap with an accumulator. This guideline is essential for both performance and guaranteed compatibility with future implementations.

## 2.3 Instructions

MMA architecture instructions are split into the following two categories:

- Accumulator operation instructions
- Outer product instructions

#### 2.3.1 Accumulator operation instructions

The first category of instructions is those that deal with moving values between the VSRs and their associated accumulator registers. The following three instructions are used to operate on the accumulator registers:

- xxmfacc: Moves the contents from the accumulator to the associated VSR.
- **xxmtacc**: Moves the contents from the associated VSR to the accumulator.
- xxsetaccz: Zeros-out the contents of the accumulator.

When the move is done, both VSRs and accumulator registers are tied and the VSRs' content becomes undefined. If an instruction tries to write content to VSRs, the accumulator content becomes undefined. When the MMA operations are done, the **xxmfacc** instruction is used to copy the content of the accumulator back to the VSRs and the VSRs become valid. For more information, see Power ISA Version 3.1.

Examples are:

xxmfacc AS

AS can be any value from 0 - 7, each referring to one accumulator register. The **xxmfacc AS** instruction moves the contents of accumulator AS to the corresponding 4 VSRs.

► xxmtacc AT

AT can be any value from 0 - 7, each referring to one accumulator register. The **xxmtacc** AT instruction moves the contents of 4 corresponding VSRs to accumulator AT.

xxsetaccz AT

AT can be any value from 0-7, each referring to one accumulator register. The **xxsetaccz AT** instruction sets the contents of accumulator **AT** to zero.

#### 2.3.2 Outer product instructions

The second category of instructions is those that are used to perform the actual arithmetic. Both integer and floating point arithmetic are supported in the MMA architecture at different precision levels as described in 2.1, "Data types" on page 12.

#### Instructions for 32-bit floating-point arithmetic

Two 32-bit FP arithmetic instructions are used to discuss the functionality of MMA. The two instructions that are used to perform a single precision matrix multiplication operation are: **xvf32ger** and **xvf32gerpp**.

The difference between the ger instruction and the gerpp instruction is as follows:

- ► The gerpp instruction *accumulates* the results in the accumulator register. This instruction requires the accumulator to already have a defined content.
- ► The ger instruction *overwrites* the results in the accumulator register. This instruction defines the content of an accumulator, similar to the xxmtacc and xxsetaccz instructions.

#### xvf32gerpp AT,XA,XB, where:

- AT refers to any of the eight accumulator registers (ACC0-ACC7).
- ► XA and XB refer to VSRs.

For the **xvf32gerpp AT,XA,XB** instruction, assume **AT=1**, **XA= 32**, and **XB=33**. The VSR 32 has four 32-bit single precision values and VSR 33 has four 32-bit single precision values. Each value in VSR 32 is multiplied with each value in VSR 33, generating a 4×4 array of 32-bit results (a total of 512 bits of output). The output is accumulated with the content of ACC1, as shown in Figure 2-2.



Figure 2-2 MMA xvf32gerpp instruction operation

#### Instructions for 8-bit arithmetic

MMA supports 8-bit integer operations. The 128-bit VSR register is split into 16 8-bit values. vi8ger4 and xvi8ger4pp are the two instructions that perform outer product operation on the 8-bit values.

xvi8ger4pp AT,XA,XB, where:

- ► AT refers to any of the eight accumulator registers (ACC0-ACC7).
- ► XA and XB refer to VSR registers.

Assume that AT=1, XA=32, and XB=33. VSR 32 has sixteen 8-bit integer values and VSR 33 has sixteen 8-bit integer values. The function of the 8-bit outer product instruction is a bit different than the 32-bit arithmetic. Each four 8-bit value in a word of XA is multiplied with each corresponding four 8-bit value in a word of XB, and the four partial products are added together to produce a 32-bit result.

The output generated by the xvi8ger4pp instruction is shown in Figure 2-3.



Figure 2-3 MMA xvi8ger4pp instruction operation

Though there are 16 8-bit values in each input VSR, the output generated still consists of 4x4 32-bit numbers and is 512 bits. The first four 8-bit elements are multiplied individually and the result of all four multiply operations is summed to generate one 32-bit value. Since there are 16 32-bit values produced, the result is still 512 bits. To use this instruction to accomplish an outer product operation, the input matrix needs to be reordered. The value formatting and how the outer product operation is performed is explained in detail in 3.3, "Mixed and lower precision matrix multiplication with MMA" on page 24.

### 2.3.3 Advanced feature: Lane masking

Lane masking is one of the advanced features available with the MMA architecture. The purpose of this feature is to perform an operation of a lower-sized input or to skip certain elements. For example, a 6x6 matrix multiplication as shown in Figure 2-4.

A 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36	B aa ab ac ad ae af ba bb bc bd be bf ca cb cc cd ce cf da bb cd de df ea eb ec ed ee ef fa fb fc fd fe ff	1 2 3 4 5 6	A <sup>T</sup> 7 13 19 25 31 8 14 20 26 32 9 15 21 27 33 10 16 22 28 34 11 17 23 29 35 12 18 24 30 36 A <sup>T</sup> aa ba ba ca ca fa fa	B ab ac ad ae af bb bc bd be bf bc cd ce cf bd bc dd de df eb ec ed ee ef bb fc fd fe ff	
-	$ \begin{bmatrix} 1 & 7 & 13 & 19 \\ [2 & 8 & 14 & 20 ] \\ [3 & 9 & 15 & 21 ] \\ [4 & 10 & 16 & 22 ] \\ [5 & 11 & 17 & 23 ] \\ [6 & 12 & 18 & 24 ] \\ \end{bmatrix} $	[ aa ab ac ad ] [ ba bb bc bd ] [ ca cb cc cd ] [ da db dc dd ] [ ea db ec ed ] [ fa fb fc fd ]		$ \begin{bmatrix} 1 & 7 & 13 & 19 \\ [2 & 8 & 14 & 20] \\ [3 & 9 & 15 & 21] \\ [4 & 10 & 16 & 22] \\ [5 & 11 & 17 & 23] \\ [6 & 12 & 18 & 24] \end{bmatrix}                                  $	ac af ] be bf ] ce cf ] de df ] ee ef ] fe ff ]
laa+2ba+3ca+4da+5ea+6f	1ab+2bb+3cb+4db+5eb+6f	lac+2bc+3cc+4cd+5ce+6f	1ad+2bd+3cd+4dd+5ed+6f	1ae+2be+3ce+4de+5ee+6f	laf+2bf+3cf+4df+5ef+6f
a	b	c	d	e	f
7aa+8ba+9ca+10da+11ea+	7ab+8bb+9cb+10db+11eb+	7ac+8bc+9cc+10cd+11ce+	7ad+8bd+9cd+10dd+11ed+	7ae+8be+9ce+10de+11ee+	7af+8bf+9cf+10df+11ef+
12fa	12fb	12fc	12fd	12fe	12ff
13aa+14ba+15ca+16da+17	13ab+14bb+15cb+16db+17	13ac+14bc+15cc+16cd+17	13ad+14bd+15cd+16dd+17	13ae+14be+15ce+16de+17	13af+14bf+15cf+16df+17
ea+18fa	eb+18fb	ce+18fc	ed+18fd	ee+18fe	ef+18ff
19aa+20ba+21ca+22da+23	19ab+20bb+21cb+22db+23	19ac+20bc+21cc+22cd+23	19ad+20bd+21cd+22dd+23	19ae+20be+21ce+22de+23	19af+20bf+21cf+22df+23
ea+24fa	eb+24fb	ce+24fc	ed+24fd	ee+24fe	ef+24ff
25aa+26ba+27ca+28da+29	25ab+26bb+27cb+28db+29	25ac+26bc+27cc+28cd+29	25ad+26bd+27cd+28dd+29	25ae+26be+27ce+28de+29	25af+26bf+27cf+28df+29
ea+30fa	eb+30fb	ce+30fc	ed+30fd	ee+30fe	ef+30ff
31aa+32ba+33ca+34da+35	31ab+32bb+33cb+34db+35	31ac+32bc+33cc+34cd+35	31ad+32bd+33cd+34dd+35	31ae+32be+33ce+34de+35	31af+32bf+33cf+34df+35
ea+36fa	eb+36fb	ce+36fc	ed+36fd	ee+36fe	ef+36ff
	[25 31] [ [26 32] [ [27 33] [ [28 34] [ [29 35] [ [30 36] [	aa ab ac ad ] ba bb bc bd ] ca cb cc cd ] da db dc dd ] fa fb fc fd		[25 31] [26 32] [27 33] [28 34] [29 35] [30 36]	[ ae af ] 6x6 [ be bf ] [ ce cf ] [ de df ] [ ee ef ] [ fe ff ]

Figure 2-4 Example of a 6x6 matrix multiplication

Figure 2-5 on page 16 shows the details of the following single precision lane-masking instruction:

Prefix::											
1	3	9		/	/				XMSK	YMS	SK
0	6	8	12	14	15	16			24	28	31
Suffix::											
59	A	r //	A	١			В		ХО	AX	3X /
0	6	9	11			16		21		29	30 31

pmxvf32gerpp AT,XA,XB,XMSK,YMSK

*Figure 2-5* Instruction word details of a prefix instruction (pmxvf32gerpp)

The instruction encoding is 64 bits. The *prefix* is the first 32 bits and the *suffix* is the next 32 bits. The prefix architecture is a new capability introduced in the Power ISA v3.1 to extend the capability of the previous 32-bit fixed instruction-size architecture. This architecture is helpful when you represent bigger instructions with more parameters. Power ISA v3.1 uses the prefix architecture in several categories of instruction.

In Figure 2-5, the MMA lane-masking instruction has a total of five input arguments in the prefix:

► The first three arguments (AT, XA, and XB) are the same as the regular 32-bit single-precision MMA instruction.

 The last two arguments (XMSK and YMSK) are the mask values for the input VSR XA and XB.

The mask values for this instruction are of size 4 bits each. Each bit masks represents one 32-bit value in the source register. The multiplication operation is performed only if both mask bits of the respective input element are set to 1. Otherwise, the respective results will be **0**.

For example, the output of the following command is shown in Figure 2-6:

pmxvf32gerpp 1,32,33,0xE,0xF

In Figure 2-6, the last 32-bit element of input register VSR[32] is skipped from the computation and the respective output to be accumulated in the ACC register 1 is 0.



Figure 2-6 MMA pmxvf32gerpp instruction operation

The gray output shown in Figure 2-6 is not computed and the corresponding four 32-bit elements in the accumulator register remain unchanged.

# Programming with Matrix-Multiply Assist

The Matrix-Multiply Assist (MMA) implementation of various kernels at different levels of precision is described in this chapter. The implementations that are shown use a single accumulator.

## 3.1 Single-precision GEMM using MMA

The innermost kernel of **sgemm\_kernel\_4x4** shown in Example 3-1 loads four elements of A, loads four elements of B, and performs an outer product MMA operation to produce one 4x4 partial result of C in one accumulator register.

Example 3-1 SGEMM kernel using MMA instructions

```
".text"
.section
        .global sgemm_kernel_4x4
                sgemm_kernel_4x4, @function
        .type
sgemm kernel 4x4:
/* adjust lda, ldb, ldc for vector size 4 */
        slwi
                7,7,2
        slwi
                8, 8, 2
        slwi
                9, 9, 2
 /* Reset accumulator */
    xxsetaccz 0
/* LOOP for K to 0 */
    K_L00P:
 /* Load 4 elements of A, B */
       1xv
              32, 0(3)
              33, 0(4)
       1xv
 /* Multiply-Add-Accumulate */
       xvf32gerpp 0, 32, 33
 /* Update Loop count & A,B */
        add
                3, 3, 7
        add
                4, 4, 8
        addic. 6, 6, -1
    bgt
           K LOOP
 /* Unprime the accumulator 0 */
    xxmfacc 0
 /* Offsets of 4x4 C Matrix */
                3, 9, 1
        slwi
        add
                4, 5, 9
        add
                6, 5, 3
                7, 4, 3
        add
 /* Store the 4x4 C Matrix */
            0, 0(5)
    stxv
            1, 0(4)
    stxv
            2, 0(6)
    stxv
    stxv
            3, 0(7)
blr
```

## 3.2 Double-precision GEMM using MMA and one accumulator

Each accumulator can hold eight (4×2) double-precision values and each VSR register (128-bit long) can contain two double-precision values (64-bit each). To produce an outer product with eight 64-bit values, we need three source VSRs so that four elements from A matrix can be multiplied with two elements from B matrix to generate a 4x2 result C submatrix. MMA enables this by taking a paired VSX register as a first operand and a single VSX register as a second operand.

For a double-precision **ger** instruction, the first operand is always an even VSX register and is considered as being paired with the next register.

Example 3-2 shows a simple example of a double-precision **gemm**. The code snippet shows the initialization of the input matrices and how an external dgemm kernel is referenced.

Example 3-2 dgemm example using 4x2 MMA kernel

```
#include <stdio.h>
#include <stdlib.h>
#define KM 4
#define KN 2
extern "C" void dgemm kernel 4x2 (double *, double *, double *, int, int, int,
int);
void dgemm(double *A, double *B, double *C, int M, int N, int K) {
    for (int i=0; i<M; i+=KM) {</pre>
    for (int j=0; j<N; j+=KN) {</pre>
       dgemm kernel 4x2(A+i, B+j, C+j, K, M, N, N);
       }
   C += N*KM;
    }
}
void printD (const char *name, double *M, int m, int n) {
        printf ("\n**** Matrix %s****\n",name);
        for (int i=0; i< m; i++) {
                printf("| ");
                for (int j=0; j< n; j++) printf("%-25.4f", *(M++));</pre>
                printf(" |\n");
        }
       }
int main (int argc, char **argv ) {
    if (argc < 4) {
       printf("Usage: %s < M > < N > < K > \n", argv[0]);
       return -1;
    }
    const int M = atoi(argv[1]);
    const int N = atoi(argv[2]);
    const int K = atoi(argv[3]);
```

```
printf("Running: %s M=%s N=%s K=%s \n", argv[0], argv[1], argv[2], argv[3]);
double A[M][K];
double AT[K][M];
double B[K][N];
double C[M][N];
for (int i=0; i<M; i++) for (int j=0; j<N; j++) C[i][j] = 0;
int x = 1;
for (int i=0; i<M; i++) for (int j=0; j<K; j++) A[i][j] = double(x++) * 7 /
15;
for (int i=0; i<K; i++) for (int j=0; j<K; j++) B[i][j] = double(x++) * 3 /
17;
for (int i=0; i<K; i++) for (int j=0; j<K; j++) AT[j][i] = A[i][j];
dgemm((double*)AT, (double*)B, (double*)C, M, N, K);
printD("C", (double *)C, M, N);
return=0;
```

The code in Example 3-3 defines the full routine of the **dgemm** 4x2 kernel in assembly using MMA instructions. This example can be compiled using the following command:

```
>> g++ -mcpu=power10 -02 dgemm_kernel_4x2.s dgemm_4x2.cc -o dgemm_mma
```

Example 3-3 A simple 4x2 dgemm kernel using MMA instructions

```
".text"
.section
        .global dgemm kernel 4x2
        .type
              dgemm_kernel_4x2, @function
dgemm kernel 4x2:
/* adjust lda, ldb, ldc for vector size 8 */
        slwi
                7, 7, 3
        slwi
                8, 8, 3
        slwi
                9, 9, 3
/* Reset acc0 */
    xxsetaccz
                 0
/* LOOP for K to 0 */
    K LOOP:
/* Load 4 elements of A, B */
       lxvp
               32, 0(3)
       lxv
               34, 0(4)
/* Multiply-Add-Accumulate */
        xvf64gerpp 0, 32, 34
/* Update Loop count & A,B */
        add
                3, 3, 7
        add
                4, 4, 8
        addic. 6, 6, -1
           K LOOP
    bgt
/* Unprime the accumulator 0 */
    xxmfacc 0
```

```
/* Offsets of 4x2 C Matrix */
       slwi
             3,9,1
       add
               4, 5, 9
       add
               6, 5, 3
       add
               7,4,3
/* Store the 4x2 C Matrix */
           0, 0(5)
   stxv
   stxv
           1, 0(4)
           2, 0(6)
   stxv
           3, 0(7)
   stxv
blr
```

# 3.3 Mixed and lower precision matrix multiplication with MMA

With mixed precision matrix multiplication, the source matrices A and B are of reduced/low precision (either 8- or 16-bit, with the same type for A and B) and the resulting matrix C is of 32-bit integer elements. MMA supports these mixed precision matrix multiplications with a broad set of instructions, as explained in the ISA.

#### 3.3.1 Source matrix reordering with Int8 as example

Consider an example of generating a 32-bit result matrix with two 8-bit lower-precision source matrices. Since VSX registers are 128 bits long, you can load 16 8-bit values from the two source matrices at each step. As the result matrix C type is 32-bit integer, accumulators still contain 16 elements (4x4 sub-matrix). The MMA operates slightly differently with lower-precision operations. Each set of four 8-bit values in each source VSX register is packed and considered as single unit. As a result, one source VSX register of 16 8-bit values is now assumed to hold four sets of four 8-bit values. Each set in the first VSX register is multiplied and accumulated with each set in second VSX register.

Within a set, the four 8-bit values are multiplied and accumulated in a one-to-one mapping, as described in Figure 3-1 on page 25 and Example 3-4 on page 25.



*Figure 3-1* A single xvi8ger instruction performs a 4x4 matrix-multiply

Example 3-4 shows how the four 8-bit values are multiplied and accumulated in a one-to-one mapping.

Example 3-4 Code showing matrix reorganization for an int-8 gemm with MMA kernel

```
#include <stdio.h>
#include <stdlib.h>
#include <inttypes.h>
```

#define KM 4
#define KN 4

#### #define Q 4

```
extern "C" void i8gemm mma 4x4 (int8 t *, uint8 t *, int32 t *, int, int, int);
void i8gemm(int8 t *A, uint8 t *B, int32 t *C, int M, int N, int K) {
    int8 t *At = (int8 t *) malloc(M*K); //transform A[M][K] --> At[K/Q][M*Q]
    uint8_t *Bt = (uint8_t *) malloc(N*K); //transform B[K][N] --> At[K/Q][N*Q]
    for (int i=0, x=0; i<K; i+=Q) {
        for(int j=0; j<M; j++) {</pre>
             for(int l=0; l<Q; l++) At[x++] = *(A+(j*K)+1);</pre>
        A = A+Q;
    }
    for (int i=0, x=0; i<K; i+=Q) {</pre>
        for(int j=0; j<N; j++) {</pre>
            for(int l=0; l<Q; l++) Bt[x++] = *(B+(l*N)+j);</pre>
        }
        B+=Q*N;
    }
    for (int i=0; i<M; i+=KM) {</pre>
    for (int j=0; j<N; j+=KN) {</pre>
        i8gemm_mma_4x4(At+(i*Q), Bt+(j*Q), C+j, K/Q, M, N, N);
        }
    C += N*KM;
    }
}
void printI (const char *name, int32_t *M, int m, int n) {
        printf ("\n**** Matrix %s****\n",name);
        for (int i=0; i< m; i++) \{
                printf("| ");
                for (int j=0; j< n; j++) printf("%-25d", *(M++));</pre>
                printf(" \\n");
        }
        printf("************************\n");
}
                 argc, char **argv ) {
int main (int
    if (argc < 4) {
        printf("Usage: %s < M > < N > < K > \n", argv[0]);
        return -1;
    }
    const int M = atoi(argv[1]);
    const int N = atoi(argv[2]);
    const int K = atoi(argv[3]);
    printf("Running: %s M=%s N=%s K=%s \n", argv[0], argv[1], argv[2], argv[3]);
    int8 t A[M][K];
    uint8_t B[K][N];
    int32_t C[M][N];
    for (int i=0; i<M; i++) for (int j=0; j<N; j++) C[i][j] = 0;
    int x = 1;
```

```
for (int i=0; i<M; i++) for (int j=0; j<K; j++) A[i][j] = (x++)%128;
for (int i=0; i<K; i++) for (int j=0; j<N; j++) B[i][j] = (x++)%256;
i8gemm((int8_t *)A, (uint8_t *)B, (int32_t *)C, M, N, K);
printI("C", (int32_t *)C, M, N);
return 0;
```

The code in Example 3-5 shows a simple 4x4 int-8 gemm kernel using MMA instructions.

Example 3-5 A 4x4 int-8 gemm kernel with xvi8ger MMA instruction

```
".text"
.section
        .global i8gemm_mma 4x4
        .type i8gemm mma 4x4, @function
i8gemm mma 4x4:
/* adjust lda, ldb, ldc for vector size 4 */
        slwi
                7, 7, 2
                8, 8, 2
        slwi
        slwi
                9, 9, 2
 /* Reset acc0 & prime */
    xxsetaccz 0
/* LOOP for K to 0 */
    K LOOP:
 /* Load 4 elements of A, B */
       1xv
              32, 0(3)
              33, 0(4)
       1xv
 /* Multiply-Add-Accumulate */
       xvi8ger4pp 0, 32, 33
 /* Update Loop count & A,B */
               3, 3, 7
       add
               4, 4, 8
       add
       addic. 6, 6, -1
           K_LOOP
    bgt
 /* Unprime the accumulator 0 */
    xxmfacc 0
 /* Offsets of 4x4 C Matrix */
        slwi
                3, 9, 1
                4, 5, 9
        add
        add
                6, 5, 3
        add
                7, 4, 3
 /* Store the 4x4 C Matrix */
            0, 0(5)
    stxv
    stxv
            1, 0(4)
            2, 0(6)
    stxv
            3, 0(7)
    stxv
blr
```

# 4

# Advanced programming concepts

This chapter describes how to optimize the matrix multiplication examples that are shown in this publication. It also explains some of the improved techniques for effectively using the compute units with load reuse and cache blocking.

# 4.1 Multiple accumulators SGEMM for load value reuse

Consider the example of the **sgemm** instruction, shown in Example 3-4 on page 25, which uses one accumulator to generate a 4x4 result, where two load operations are required for each **gerpp** instruction. This type of model restricts the MMA unit utilization, as it is limited by the number of load ports. This limitation can be overcome by using multiple accumulators to reuse the same loaded data multiple times, as shown in Figure 4-1 and Example 4-1. When you use all eight accumulators and reuse loads, the result is eight **gerpp** instructions for every six vector register loads. (The code uses the register pair **load lxvp** instruction, which is another feature of Power ISA v3.1.) This kernel generates an 8x16 submatrix of C (eight 4x4 accumulators).



Figure 4-1 Example of an 8x16 GEMM FP32 kernel execution

Example 4-1 GEMM FP32 kernel execution code

```
#define KM 8
#define KN 16
extern "C" void sgemm_mma_8x16 (float*, float*, float*, int, int, int, int);
void sgemm(float *A, float *B, float *C, int M, int N, int K) {
   for (int i=0; i<M; i+=KM) {
    for (int j=0; j<N; j+=KN) {
      sgemm_mma_8x16(A+i, B+j, C+j, K, M, N, N);
      }
   C += N*KM;
   }
}</pre>
```

Example 4-2 shows a code snippet for a **sgemm** kernel that reuses the loads and all eight accumulators for improved compute unit utilization. The code computes a resultant 8x16 submatrix.

Example 4-2 An 8-accumulator version of sgemm kernel

```
.section
               ".text"
        .global sgemm mma 8x16
        .type
               sgemm_mma_8x16, @function
sgemm mma 8x16:
/* adjust lda, ldb, ldc for vector size 4 */
       slwi
               7, 7, 2
               8, 8, 2
       slwi
               9, 9, 2
       slwi
 /* Reset acc0-7 & prime */
    xxsetaccz
                0
   xxsetaccz
                1
   xxsetaccz
                2
   xxsetaccz
                3
   xxsetaccz
                4
   xxsetaccz
                5
   xxsetaccz
                6
   xxsetaccz
               7
/* LOOP for K to 0 */
   K L00P:
/* Load 4 elements of A, B */
      lxvp
             32, 0(3)
      lxvp
             34, 0(4)
            36, 32(4)
      lxvp
/* Multiply-Add-Accumulate */
      xvf32gerpp 0, 33, 35
      xvf32gerpp 1, 33, 34
      xvf32gerpp 2, 33, 37
      xvf32gerpp 3, 33, 36
      xvf32gerpp 4, 32, 35
      xvf32gerpp 5, 32, 34
      xvf32gerpp 6, 32, 37
      xvf32gerpp 7, 32, 36
/* Update Loop count & A,B */
      add
              3, 3, 7
      add
              4, 4, 8
      addic. 6, 6, -1
```

```
bgt
           K LOOP
/* Unprime the accO-7 */
    xxmfacc 0
    xxmfacc 1
    xxmfacc 2
    xxmfacc 3
    xxmfacc 4
    xxmfacc 5
    xxmfacc 6
    xxmfacc 7
 /* Offsets of 4x4 C Matrix */
    slwi 3, 9, 1
    add
         4, 5, 9
    add
          6, 5, 3
    add
         7, 4, 3
 /* Store the 4x16 c Matrix */
            3, 0(5)
    stxv
    stxv
            2, 0(4)
    stxv
            1, 0(6)
    stxv
            0, 0(7)
    stxv
            7,
                16(5)
    stxv
            6,
                16(4)
    stxv
            5, 16(6)
    stxv
            4, 16(7)
            11, 32(5)
    stxv
            10, 32(4)
    stxv
    stxv
            9, 32(6)
    stxv
            8, 32(7)
            15, 48(5)
    stxv
            14, 48(4)
    stxv
            13, 48(6)
    stxv
            12, 48(7)
    stxv
 /* Update index of C */
    add
         5,7,9
    add
          4, 5, 9
    add
          6, 5, 3
    add
         7,4,3
 /* Store the 4x16 c Matrix */
    stxv
            19, 0(5)
    stxv
            18, 0(4)
    stxv
            17, 0(6)
            16, 0(7)
    stxv
            23, 16(5)
    stxv
            22, 16(4)
    stxv
    stxv
            21, 16(6)
    stxv
            20, 16(7)
    stxv
            27, 32(5)
            26, 32(4)
    stxv
    stxv
            25, 32(6)
    stxv
            24, 32(7)
    stxv
            31, 48(5)
    stxv
            30, 48(4)
    stxv
            29, 48(6)
    stxv
            28, 48(7)
blr
```

### 4.2 Multiple accumulators DGEMM for load value reuse

From the DGEMM example and as specified in MMA design, a single 512-bit accumulator result of a double-precision instruction can hold a resultant 4x2 submatrix. To perform one double-precision MMA instruction, two VSR registers (two loads, four elements) from matrix A and one VSR register (one load, two elements) from matrix B. By using all eight accumulators, you can generate a resultant 8x8 result. Including the reusing of loads, a total of eight VSR loads (four each from matrices A and B) need to be performed. A simple example of an dgemm\_kerne1\_8x8 is shown in Example 4-3.

Example 4-3 8x8 FP64 GEMM kernel code

```
#define DKM 8
#define DKN 8
extern "C" void dgemm_kernel_8x8 (double *, double *, double *, int, int, int);
void dgemm(double *A, double *B, double *C, int M, int N, int K) {
   for (int i=0; i<M; i+=DKM) {
    for (int j=0; j<N; j+=DKN) {
      dgemm_kernel_8x8(A+i, B+j, C+j, K, M, N, N);
      }
   C += N*DKM;
   }
}</pre>
```

The code snippet in Example 4-4 shows an improved **dgemm** kernel using eight accumulators. This code computes an 8x8 block of result matrix.

Example 4-4 An 8x8 dgemm kernel using 8 accumulators

```
".text"
        .section
       .global dgemm kernel 8x8
       .type dgemm kernel 8x8, @function
dgemm kernel 8x8:
/* adjust lda, ldb, ldc for vector size 8 */
       slwi
               7,7,3
       slwi
               8, 8, 3
       slwi
               9, 9, 3
/* Reset acc0-7 & prime */
   xxsetaccz
                0
   xxsetaccz
                1
   xxsetaccz 2
             3
   xxsetaccz
   xxsetaccz
               4
   xxsetaccz
                5
   xxsetaccz
                6
   xxsetaccz
                7
/* LOOP for K to 0 */
   K LOOP:
/* Load 4 elements of A, B */
      1xvp 32, 0(3)
      lxvp
             34, 32(3)
```

```
lxvp
              36, 0(4)
       lxvp
              38, 32(4)
/* Multiply-Add-Accumulate */
       xvf64gerpp 0, 32, 37
       xvf64gerpp 1, 32, 36
       xvf64gerpp 2, 32, 39
       xvf64gerpp 3, 32, 38
       xvf64gerpp 4, 34, 37
       xvf64gerpp 5, 34, 36
       xvf64gerpp 6, 34, 39
       xvf64gerpp 7, 34, 38
/* Update Loop count & A,B */
       add
               3, 3, 7
       add
               4, 4, 8
       addic. 6, 6, -1
  bgt
          K LOOP
/* Unprime the acc0-7 */
    xxmfacc 0
    xxmfacc 1
    xxmfacc 2
    xxmfacc 3
    xxmfacc 4
    xxmfacc 5
    xxmfacc 6
    xxmfacc 7
/* Offsets of 4x4 C Matrix */
    slwi
         3, 9, 1
    add
          4, 5, 9
          6, 5, 3
    add
    add
          7, 4, 3
/* Store the 4x16 c Matrix */
            3, 0(5)
    stxv
    stxv
            2, 0(4)
    stxv
            1, 0(6)
    stxv
            0, 0(7)
            7, 16(5)
    stxv
    stxv
            6, 16(4)
    stxv
            5, 16(6)
    stxv
            4, 16(7)
    stxv
            11, 32(5)
            10, 32(4)
    stxv
    stxv
            9, 32(6)
            8, 32(7)
    stxv
    stxv
            15, 48(5)
    stxv
            14, 48(4)
    stxv
            13, 48(6)
    stxv
            12, 48(7)
/* Update index of C */
    add
          5, 7, 9
    add
          4, 5, 9
          6, 5, 3
    add
    add
          7, 4, 3
/* Store the 4x16 c Matrix */
            19, 0(5)
    stxv
    stxv
            18, 0(4)
```

stxv	17,	0(6)
stxv	16,	0(7)
stxv	23,	16(5)
stxv	22,	16(4)
stxv	21,	16(6)
stxv	20,	16(7)
stxv	27,	32(5)
stxv	26,	32(4)
stxv	25,	32(6)
stxv	24,	32(7)
stxv	31,	48(5)
stxv	30,	48(4)
stxv	29,	48(6)
stxv	28,	48(7)
b1	r	

**Programming note:** Instead of the resultant 8x8 submatrix shown in Example 4-4 on page 33, you can use an alternate approach to generate a resultant 4x16 submatrix for dgemm using eight accumulators. In this approach, the kernel requires ten vector loads, instead of eight vector loads. Though the kernel stresses the load units, there might be a slight advantage over the cache blocking method, which is discussed in 4.3, "SGEMM performance with advanced cache-blocking" on page 35.

## 4.3 SGEMM performance with advanced cache-blocking

Consider the multiplication of two float (single-precision) 256 x 256 matrices using an **sgemm\_kerne1\_8x16**, which computes the resultant C matrix in 8x16 blocks. In this example, M=256, K=256, and N=256. To compute the full resultant 256 x 256 matrix C, the program loops over A in blocks of size 8x256 and over B in blocks of size 16x256, as shown in Example 4-5.

Example 4-5 Full computation of matrix C in blocks of size 8x16

```
#define KM 8
#define KN 16
void sgemm(float *A, float *B, float *C, int M, int N, int K) {
    for (int i=0; i<M; i+=KM) {
        for (int j=0; j<N; j+=KN) {
            sgemm_kernel_8x16(A+i, B+j, C+j, K, M, N, N);
        }
        C += N*KM;
    }
}</pre>
```

To carry out the computation of one 8x16 block of matrix C in a processing core (call to **sgemm\_kernel\_8x16** in the innermost kernel of Example 4-5), you need one 8x256 block of matrix A (8x256x4 bytes = 8 KiB) and one 16x256 block of matrix B (16x256x4 bytes = 16 KiB) to be loaded into the L1 cache. You can then compute one 8x16 block of resultant matrix C. In the second iteration of the innermost loop, the 8-KiB block of matrix A is retained and the next block of 16-KiB data of matrix B is freshly loaded onto L1 to compute the next 8x16 block of matrix C.

To improve the performance of the computation, retain the bigger block of data (from matrix B) in L1 cache and move through the smaller blocks (from matrix A). This is achieved by simply interchanging the loops, as shown in Example 4-6.

Example 4-6 Loops interchanged to keep the bigger block in L1 cache

```
#define KM 8
#define KN 16
void sgemm(float *A, float *B, float *C, int M, int N, int K) {
    for (int j=0; j<N; j+=KN) {
        for (int i=0; i<M; i+=KM) {
            sgemm_kernel_8x16(A+i, B+j, C+(N*KM), K, M, N, N);
        }
        C += KN;
    }
}</pre>
```

Assume an architecture with an L1 cache of 32 KiB. In the 256x256 matrix-multiplication example shown in Example 4-6, 16 KiB of matrix B and 8 KiB of matrix A are consumed by each execution of **sgemm\_kernel\_8x16**. The total of 24 KiB fits well within the L1 cache. Now consider an example of K=1024. The size of matrix B block becomes 16x1024x4 bytes = 64 KiB and the size of matrix A block grows to 32 KiB. The total exceeds the L1 cache size by a factor of three. Therefore, to keep the computation well within the L1 cache, there needs to be a third loop to iterate over 'K' in chunks of size KS. If KS=256 in this case, you need to repeat the above set of iterations K/KS = 1024/256 = 4 times.

# 5

# Matrix-Multiply Assist programming with compiler built-ins

Support for Matrix-Multiply Assist (MMA) instructions and built-ins is enabled in the latest version of GCC and LLVM compilers, which are publicly available. You can use GCC 10 and later and LLVM 12.0 to build programs with MMA. Though the compiler does not yet generate MMA instructions from the direct existing source code, support is available to recognize inline MMA assembly instructions. The compilers can also generate binaries for **genm** programs with MMA compiler built-ins.

For MMA programming, in addition to the already supported vector data types, the new \_\_vector\_quad data type is introduced to represent an accumulator. This data type represents a set of four vector registers forming an accumulator.

The following set of built-ins is used to move values to and from vector registers to accumulators:

void \_\_builtin\_mma\_xxmtacc (\_\_vector\_quad \*);

```
void __builtin_mma_xxmfacc (__vector_quad *);
```

The following built-in can be used to set the accumulator to zero:

```
void __builtin_mma_xxsetaccz (__vector_quad *);
```

The following built-ins can be used to collate and dismantle the accumulator register (\_\_vector\_quad) to four independent registers to be further used in the program:

- ▶ void \_\_builtin\_mma\_assemble\_acc (\_\_vector\_quad \*, vec\_t, vec\_t, vec\_t, vec\_t);
- void \_\_builtin\_mma\_disassemble\_acc (void \*, \_\_vector\_quad \*);

The builtin\_mma\_xv\* can be used to perform the matrix-multiply-and-accumulate operations. The following code is an example built-in for 32-bit floating-point ger operation:

```
void __builtin_mma_xvf32gerpp (__vector_quad *, vec_t, vec_t);
```

### 5.1 Simple MMA SGEMM example using built-ins

Use the **gcc** compiler with the following command to compile the code example shown in Example 5-1. For a list of MMA compiler built-ins, see Appendix A, "List of Matrix-Multiply Assist compiler built-ins" on page 41.

> gcc -o sgemm -02 -mcpu=power10 -mtune=power10 sgemm\_intrinsics.cc

Example 5-1 Sample GEMM code using MMA compiler built-ins

```
#define KM 4
#define KN 4
typedef vector unsigned char
                                 vec t;
typedef __vector_quad acc_t;
void sgemm kernel 4x4 (float *a, float *b, float *c, int K, int lda, int ldb, int
ldc) {
        int i;
        vec_t vec_A, vec_B, vec_C[4];
        acc t acc 0;
        builtin mma xxsetaccz(&acc 0);
        for (i=0; i<K; i++) {
                vec A = *((vec t *)(a+(i*lda)));
                vec B = *((vec t *)(b+(i*ldb)));
                __builtin_mma_xvf32gerpp(&acc_0, vec_A, vec_B);
        }
         __builtin_mma_disassemble_acc(vec_C, &acc_0);
        *((vec_t *)(c)) = vec_C[0];
        *((vec t *)(c+ldc)) = vec C[1];
        *((vec t *)(c+(2*ldc))) = vec C[2];
        *((vec_t *)(c+(3*1dc))) = vec_C[3];
}
void sgemm(float *A, float *B, float *C, int M, int N, int K) {
    for (int i=0; i<M; i+=KM) {</pre>
    for (int j=0; j<N; j+=KN) {</pre>
        sgemm kernel 4x4(A+i, B+j, C+j, K, M, N, N);
        }
    C += N*KM;
    }
```

**Note:** The following notes provide important programming information for successful use of the built-ins. Follow these suggestions to avoid having syntax errors reported from the front end:

- When passing arrays to built-ins that expect a void \* pointer, there needs to be an explicit cast.
- When declaring vectors that are passed to built-ins, use Altivec vector syntax (such as vector unsigned char and vector double), rather than another generic vector syntax (such as \_\_attribute\_\_((vector\_size(16)) or similar).

**Performance note:** Additional care should be taken to avoid Pipeline flushes while programming MMA.

The following Pipeline flushes degrade gemm performance and should be avoided:

- A VSR Conflict Flush occurs when you access a VSR associated with a primed accumulator.
- An Accumulator Conflict Flush occurs when you perform an MMA ger instruction without priming the accumulator.

# Α

# List of Matrix-Multiply Assist compiler built-ins

**Programming note:** Type vec\_t is defined to be a normal vector unsigned char type. The uint2, uint4, and uint8 parameters are 2-bit, 4-bit, and 8-bit unsigned integer constants, respectively. The compiler verifies that they are constants and that their values are within range.

```
void __builtin_mma_xvi4ger8 (__vector_quad *, vec_t, vec_t);
void __builtin_mma_xvi8ger4 (__vector_quad *, vec_t, vec_t);
void __builtin_mma_xvi16ger2 (__vector_quad *, vec_t, vec_t);
void builtin mma xvi16ger2s ( vector quad *, vec t, vec t);
void __builtin_mma_xvf16ger2 (__vector_quad *, vec t, vec t);
void builtin mma xvbf16ger2 ( vector quad *, vec t, vec t);
void __builtin_mma_xvf32ger (__vector_quad *, vec_t, vec_t);
      _builtin_mma_xvi4ger8pp (__vector_quad *, vec_t, vec_t);
void
      _builtin_mma_xvi8ger4pp (__vector_quad *, vec_t, vec_t);
void
      _builtin_mma_xvi8ger4spp(__vector_quad *, vec_t, vec_t);
void
void _
      _builtin_mma_xvi16ger2pp (__vector_quad *, vec_t, vec_t);
void __builtin_mma_xvi16ger2spp (__vector_quad *, vec_t, vec_t);
void __builtin_mma_xvf16ger2pp (__vector_quad *, vec_t, vec_t);
void __builtin_mma_xvf16ger2pn (__vector_quad *, vec_t, vec_t);
void builtin mma xvf16ger2np ( vector quad *, vec t, vec t);
void builtin mma xvf16ger2nn ( vector quad *, vec t, vec t);
void builtin mma xvbf16ger2pp ( vector quad *, vec t, vec t);
void __builtin_mma_xvbf16ger2pn (__vector_quad *, vec_t, vec_t);
void
      _builtin_mma_xvbf16ger2np (__vector_quad *, vec_t, vec_t);
void
      _builtin_mma_xvbf16ger2nn (__vector_quad *, vec_t, vec_t);
      _builtin_mma_xvf32gerpp (__vector_quad *, vec_t, vec_t);
_builtin_mma_xvf32gerpn (__vector_quad *, vec_t, vec_t);
void
void
void __builtin_mma_xvf32gernp (__vector_quad *, vec_t, vec_t);
void builtin mma_xvf32gernn (__vector_quad *, vec_t, vec_t);
void __builtin_mma_pmxvi4ger8 (__vector_quad *, vec_t, vec_t, uint4, uint4, uint8);
```

void \_\_builtin\_mma\_pmxvi4ger8pp (\_\_vector\_quad \*, vec\_t, vec\_t, unit, unit, unit, unit8);

void \_\_builtin\_mma\_pmxvi8ger4 (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint4); void builtin mma pmxvi8ger4pp ( vector quad \*, vec t, vec t, uint4, uint4, uint4); void builtin mma pmxvi8ger4spp( vector quad \*, vec t, vec t, uint4, uint4, uint4); void \_\_builtin\_mma\_pmxvi16ger2 (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); \_builtin\_mma\_pmxvil6ger2s (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void void \_builtin\_mma\_pmxvf16ger2 (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void \_\_builtin\_mma\_pmxvbf16ger2 (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void \_\_builtin\_mma\_pmxvi16ger2pp (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void \_\_builtin\_mma\_pmxvi16ger2spp (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint2); void \_\_builtin\_mma\_pmxvf16ger2pp (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void \_\_builtin\_mma\_pmxvf16ger2pn (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void builtin mma pmxvf16ger2np ( vector quad \*, vec t, vec t, uint4, uint4, uint2); void builtin mma pmxvf16ger2nn ( vector quad \*, vec t, vec t, uint4, uint4, uint2); void \_\_builtin\_mma\_pmxvbfl6ger2pp (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void \_builtin\_mma\_pmxvbf16ger2pn (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); \_builtin\_mma\_pmxvbf16ger2np (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void void \_\_builtin\_mma\_pmxvbf16ger2nn (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4, uint2); void \_builtin\_mma\_pmxvf32ger (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4); void \_\_builtin\_mma\_pmxvf32gerpp (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4); void \_\_builtin\_mma\_pmxvf32gerpn (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4); void \_\_builtin\_mma\_pmxvf32gernp (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4); void \_\_builtin\_mma\_pmxvf32gernn (\_\_vector\_quad \*, vec\_t, vec\_t, uint4, uint4); void builtin mma\_xvf64ger (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t); void \_\_builtin\_mma\_xvf64gerpp (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t); void \_\_builtin\_mma\_xvf64gerpn (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t); \_builtin\_mma\_xvf64gernp (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t); void void \_\_builtin\_mma\_xvf64gernn (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t); void \_builtin\_mma\_pmxvf64ger (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t, uint4, uint2); void \_\_builtin\_mma\_pmxvf64gerpp (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t, uint4, uint2); void \_\_builtin\_mma\_pmxvf64gerpn (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t, uint4, uint2); void \_\_builtin\_mma\_pmxvf64gernp (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t, uint4, uint2); void \_\_builtin\_mma\_pmxvf64gernn (\_\_vector\_quad \*, \_\_vector\_pair, vec\_t, uint4, uint2); void builtin mma xxmtacc ( vector quad \*); void builtin mma xxmfacc ( vector quad \*); void \_\_builtin\_mma\_xxsetaccz (\_\_vector\_quad \*); void \_\_builtin\_mma\_assemble\_acc (\_\_vector\_quad \*, vec\_t, vec\_t, vec\_t, vec\_t); void \_\_builtin\_mma\_disassemble\_acc (void \*, \_\_vector\_quad \*); void \_builtin\_mma\_assemble\_pair (\_\_vector\_pair \*, vec\_t, vec\_t); void \_\_builtin\_mma\_disassemble\_pair (void \*, \_\_vector\_pair \*); vec\_t \_\_builtin\_xvcvspbf16 (vec\_t); vec t builtin xvcvbf16sp (vec t);

# Β

# List of Matrix-Multiply Assist instructions in Power ISA v3.1

Table 5-1lists the available Matrix-Multiply Assist (MMA) instructions defined in Power ISA v3.1. For details and syntax, see Power ISA Version 3.1.

MMA instruction type	Traditional instructions 32-bit encoding	Prefix instructions 64-bit encoding
Data movement	xxmfacc xxmtacc xxsetaccz	
64-bit floating-point inputs (IEEE double-precision)	xvf64ger2 xvf64ger2nn xvf64ger2np xvf64ger2pn xvf64ger2pp	pmxvf64ger2 pmxvf64ger2nn pmxvf64ger2np pmxvf64ger2pn pmxvf64ger2pp
32-bit floating-point inputs (IEEE single-precision)	xvf32ger2 xvf32ger2nn xvf32ger2np xvf32ger2pn xvf32ger2pp	pmxvf32ger2 pmxvf32ger2nn pmxvf32ger2np pmxvf32ger2pn pmxvf32ger2pp
16-bit floating-point inputs (IEEE half-precision)	xvf16ger2 xvf16ger2nn xvf16ger2np xvf16ger2pn xvf16ger2pp	pmxvf16ger2 pmxvf16ger2nn pmxvf16ger2np pmxvf16ger2pn pmxvf16ger2pp
16-bit floating-point inputs (bfloat16 format)	xvbf16ger2 xvbf16ger2nn xvbf16ger2np xvbf16ger2pn xvbf16ger2pp	pmxvbf16ger2 pmxvbf16ger2nn pmxvbf16ger2np pmxvbf16ger2pn pmxvbf16ger2pp

Table 5-1 MMA instructions defined in Power ISA v3.1

MMA instruction type	Traditional instructions 32-bit encoding	Prefix instructions 64-bit encoding
16-bit integer inputs (modulo arithmetic)	xvi16ger2 xvi16ger2pp	pmxvi16ger2 pmxvi16ger2pp
16-bit integer inputs (saturating arithmetic)	xvi16ger2s xvi16ger2spp	pmxvi16ger2s pmxvi16ger2spp
8-bit integer inputs (modulo/saturating)	xvi8ger4 xvi8ger4pp xvi8ger4spp	pmxvi8ger4 pmxvi8ger4pp pmxvi8ger4spp
4-bit integer inputs	xvi4ger8 xvi4ger8pp	pmxvi4ger8 pmxvi4ger8pp

# **Related publications**

The publications listed in this section are considered particularly suitable for a more detailed discussion of the topics covered in this paper.

#### **Online resources**

These websites are also relevant as further information sources:

- Power ISA Version 3.1 https://wiki.raptorcs.com/w/images/f/f5/PowerISA public.v3.1.pdf
- Anatomy of High-Performance MatrixMultiplication https://www.cs.utexas.edu/users/flame/pubs/GotoTOMS final.pdf
- 754-2019 IEEE Standard for Floating-Point Arithmetic https://ieeexplore.ieee.org/document/8766229
- A transprecision floating-point platform for ultra-low power computing https://ieeexplore.ieee.org/abstract/document/8342167

### Help from IBM

IBM Support and downloads **ibm.com**/support IBM Global Services **ibm.com**/services



REDP-5612-00

ISBN 0738459453

Printed in U.S.A.



**Get connected** 

