

---

# Ramanujan Bipartite Graph Products for Efficient Block Sparse Neural Networks

---

**Dharma Teja Vooturi, Girish Varma, Kishore Kothapalli**  
Center for Security Theory and Algorithmic Research  
International Institute of Information Technology Hyderabad, India  
dharmateja.vooturi@research.iiit.ac.in

## Abstract

Sparse neural networks are shown to give accurate predictions competitive to denser versions, while also minimizing the number of arithmetic operations performed. However current hardware like GPU's can only exploit structured sparsity patterns for better efficiency. Hence the run time of a sparse neural network may not correspond to the arithmetic operations required.

In this work, we propose RBGP( Ramanujan Bipartite Graph Product) framework for generating structured multi level block sparse neural networks by using the theory of Graph products. We also propose to use products of Ramanujan graphs which gives the best connectivity for a given level of sparsity. This essentially ensures that the i.) the networks has the structured block sparsity for which runtime efficient algorithms exists ii.) the model gives high prediction accuracy, due to the better expressive power derived from the connectivity of the graph iii.) the graph data structure has a succinct representation that can be stored efficiently in memory. We use our framework to design a specific connectivity pattern called RBGP4 which makes efficient use of the memory hierarchy available on GPU. We benchmark our approach by experimenting on image classification task over CIFAR dataset using VGG19 and WideResnet-40-4 networks and achieve 5-9x and 2-5x runtime gains over unstructured and block sparsity patterns respectively, while achieving the same level of accuracy.

## 1 Introduction

Sparsity is an essential tool for generating compute and memory efficient neural networks. Despite this, the predominant choice of deep neural networks in production are dense instead of sparse. This is mainly because sparse neural networks tend to have poor runtime performance on the widely used dense AI hardware like GPU/TPU, that are primarily designed for accelerating dense neural networks. So in order to truly uncover the potential of sparsity in production, it is necessary to generate sparse neural networks, that are in harmony with the dense AI hardware.

Pruning [16, 11, 10, 9] is one of the widely used approach for generating sparse neural networks. In element pruning, individual parameters/elements are removed from a pre-trained dense neural network based on some criterion such as magnitude, and then the resultant sparse network is finetuned to recover accuracy. Significant number of parameters can be removed by using element pruning with minimal loss in model accuracy. But the main issue with element pruning is that the generated sparse neural networks have irregular compute and memory access patterns due to unstructured sparsity pattern, and thus cannot be efficiently mapped onto dense AI hardware. Structured pruning methods [18, 26, 12, 22, 23, 36, 4, 33] are proposed to improve the runtime performance of sparse neural networks. Unlike element pruning, where parameters are removed at an individual level, in structured pruning, parameters are first divided into structural units like filter, channel, block, multi-block etc and

Preprint. Under review.

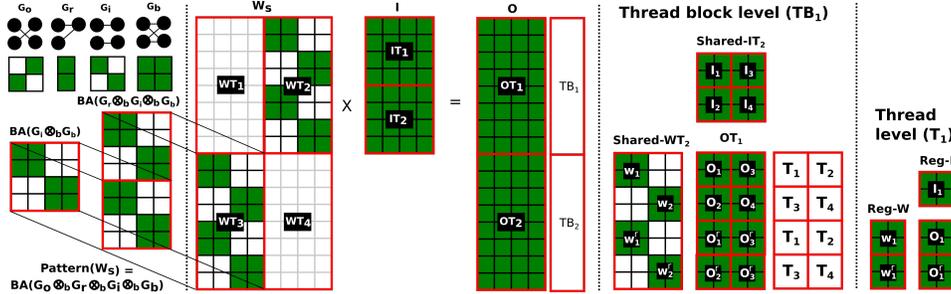


Figure 1: Tiled matrix multiplication of RBGP4 sparse matrix  $W_s$  with a dense matrix  $I$  ( $O = W_s \times I$ ) on GPU. A tile in  $O$  ( $OT$ ) is mapped to a thread block  $TB$ , and each thread in  $TB$  is mapped to a 2D strided grid of element blocks in  $OT$ , where the number of strides, and the size of the element block in row dimension are set to  $|G_r \cdot U|$  and  $|G_b \cdot U|$  respectively.  $OT$  is computed in steps, where in each step, tiles  $WT$  and  $IT$  are first loaded into shared memory from DRAM, and a thread in  $TB$  loads corresponding elements from shared memory to registers before performing the computation.

then are removed at a unit level based on the strength of the unit. Structured sparse neural networks have better run-time performance than unstructured sparse neural networks. But this improvement in run-time performance comes at the cost of accuracy due to the imposed structural constraints while removing parameters from a trained model. For example, Mao et al. [23] have shown that for a given amount of pruning, model accuracy decreases and run-time performance increases with increase in coarsity of structural unit from 0D to 3D in pruning 4D weight tensors in convolutional neural networks. This trade-off between run-time and accuracy limits the possibility of generating efficient structured sparse neural networks using structured pruning methods. Structured sparse neural networks can also be generated using structure aware training (STAT) methods [35, 29, 19, 14, 34, 15], where structure is part of the training process. Because the structure is coupled with the training process, STAT methods are better placed than structured pruning methods in generating efficient structured sparse neural networks.

Runtime of a sparse neural network on a given hardware is dependent on the efficiency with which SDMM (Multiplication of a Sparse Matrix with a Dense matrix) operation can be implemented. On a hardware like GPU with memory hierarchy (Registers > Shared memory > L2 cache > DRAM), SDMM operation will have good runtime efficiency if and only if it maximizes data accesses from faster memory through data reuse. And for a structured sparse neural network, the amount of reuse depends on the choice of the structured sparsity pattern. Additionally, the chosen pattern should be well connected to allow for good flow of information in the neural network. In this work, we address these requirements and generate structured sparse networks that are performant and connected. Following are our main contributions:

- Proposed RBGP (Ramanujan Bipartite Graph Product) framework for generating structured sparse neural networks that have multiple levels of block sparsity, good connectivity, and takes less memory for storage.
- Using RBGP framework, we proposed RBGP4 structured sparsity pattern for the GPU, a representative dense hardware, and achieve good runtime efficiency for the SDMM (Multiplication of a sparse matrix with a dense matrix) operation on GPU.
- We demonstrate the utility of RBGP4 sparsity pattern on image classification task over CIFAR dataset and achieve 5-9x and 2-5x runtime gains over unstructured and block sparsity patterns respectively, while achieving the same level of accuracy.

## 2 Related work

**Post training:** Generating sparse neural network from a trained dense model dates back to decades old work of Lecun et al. [16] and Hassibi & Stork [11] where they use second-derivative information to prune weights from a dense model. The idea of pruning was revived by Han et al. [10, 9] by simply pruning weights based on their magnitude. To improve runtime performance on dense AI hardware, structured pruning methods [18, 26, 12, 22, 23, 36, 4, 33] are proposed with various structured sparsity patterns like filter, channel, block and multi-block.

**During training:** Sparse neural networks are generated during the training process either by gradually removing the connections or rearranging existing set of connections [32, 28, 2, 25, 27, 17, 6]. Similarly, structured sparse networks are generated by removing elements at a structural unit level during training. Wen et al. [35] used group Lasso regularization to induce channel and filter sparsity in CNNs. Narang et al. [29] used gradual pruning along with group Lasso regularization to induce block sparsity pattern in RNNs. In [19, 14, 34], structure is induced by assigning a learnable parameter for each structural unit and removing them gradually through regularization and pruning.

**Before training(predefined):** Sparsity can be incorporated apriori to the training process by choosing a mask(choice of connections) in each layer of the sparse neural network and keeping it fixed through out the training. Prior works in predefined approach differ in the way the mask is chosen. Prabhu et al.[30] makes use of expander graphs, and generates a random mask with row uniformity pattern, where all the rows in the mask have equal number of non zeros. Sourya et al.[7] generates a random mask with both row and column uniformity. Frankle et al. [8] uses an unstructured mask generated by pruning a trained dense model. Kepner et al. [15] uses the idea of radix topology to generate a mask with cyclical diagonal pattern. Blocking pattern is the key requirement for achieving runtime performance on dense AI hardware, and none of the above works incorporate block sparsity pattern. In this work, we impose block sparsity pattern at multiple levels using RBGP framework, and achieve good runtime performance on GPU, a representative dense AI hardware.

### 3 Preliminaries

In this section, we setup various definitions and notations used throughout the paper. First we define various types of block sparsity patterns.

**Block Sparse (BS) matrix:** A BS matrix  $W_{bs}$  is a sparse matrix, where non zero elements are structured in the form of blocks of size  $(bh, bw)$ . Matrix  $W_{bs}$  has  $(W_{bs}.rows/bh \times W_{bs}.columns/bw)$  number of blocks, and a block in  $W_{bs}$  is either a zero block with all zeros or a non-zero block with some or all elements as non-zeros.

**Uniform Block Sparse (UBS) matrix:** A UBS matrix  $W_{ubs}$  is a block sparse matrix with block size  $(bh, bw)$ , where all the row/column blocks of size  $(bh, W_{ubs}.columns)/(W_{ubs}.rows, bw)$  have equal number of non-zero blocks of size  $(bh, bw)$ .

**Cloned Block Sparse (CBS) matrix:** A CBS matrix is a block sparse matrix with block size  $(bh, bw)$ , where all the non zero blocks of size  $(bh, bw)$  have the same non-zero pattern.

**Cloned Uniform Block Sparse (CUBS) matrix:** A CUBS matrix is a block sparse matrix with block size  $(bh, bw)$  that is both UBS and CBS matrix with block size  $(bh, bw)$ .

**Recursive CUBS (RCUBS) matrix:** An RCUBS matrix  $W_s$  is a sparse matrix with  $K$  levels of blocking  $B_1, \dots, B_K$  and following recursion:  $W_s$  is a CUBS matrix with block size  $B_1$ , and a non zero block of size  $B_i$  in  $W_s$  is again a CUBS matrix with block size  $B_{i+1}$ . Figure 3 shows an example of RCUBS matrix with three levels of blocking.

We consider the Bipartite graph  $G = (U, V, E)$  representation of matrices (with dimension  $|U| \times |V|$ ). In a biregular bipartite graph, all the vertices in  $U$  and  $V$  have same degree  $d_l$  and  $d_r$  respectively. The degree also characterizes the sparsity of such graphs. The eigenvalues of a graph  $G$  are the eigenvalues of its adjacency matrix and they characterize many graph properties including connectivity [5]. Bipartite graph with  $N$  vertices have Eigen values  $\pm\lambda_1, \dots, \pm\lambda_{N/2}$ , where  $\lambda_1 \geq \lambda_2 \dots \geq \lambda_{N/2}$ . The *spectral gap* between  $\lambda_1, \lambda_2$  is a measure of the connectivity properties of the graph [1]. *Ramanujan Graphs* are the graphs with the optimal connectivity (as measured by the spectral gap) for a given level of sparsity [21].

**Ramanujan bipartite graph:** A Ramanujan bipartite graph is a  $(d_l, d_r)$ -biregular bipartite graph, where the second largest eigenvalue  $\lambda_2$  is less than or equal to  $(\sqrt{d_l - 1} + \sqrt{d_r - 1})$ .

**Bipartite Graph Product ( $\otimes_b$ ):** Bipartite graph product( $G_p = G_1 \otimes_b G_2$ ) takes two bipartite graphs,  $G_1(U_1, V_1, E_1)$  and  $G_2(U_2, V_2, E_2)$  as the input and produces a bigger bipartite graph  $G_p(U_p, V_p, E_p)$ , where  $U_p = U_1 \times U_2$ ,  $V_p = V_1 \times V_2$ , and  $E_p$  is constructed using cross product of edges from  $G_1$  and  $G_2$  i.e,  $E_p = \{((u_1, u_2), (v_1, v_2)) | ((u_1, v_1) \in E_1 \& (u_2, v_2)) \in E_2\}$ . Bipartite graph product can also be viewed from a matrix viewpoint in the following way:

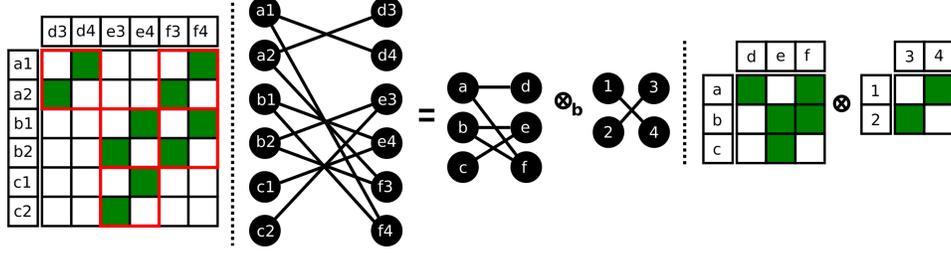


Figure 2: Bipartite graph product operation( $\otimes_b$ ) along with matrix view. Biadjacency matrix of the product graph has CBS(Cloned Block Sparse) pattern with block size (2,2).

A bipartite graph  $G(U, V, E)$  can be represented as a bi-adjacency matrix  $BA$  of size  $(|U|, |V|)$ , with  $BA_{uv} = 1$  if  $(u, v) \in E$ , and zero otherwise. For the bipartite graph product( $G_p = G_1 \otimes_b G_2$ ), bi-adjacency matrix of  $G_p$  is equal to the Tensor product( $\otimes$ ) of the bi-adjacency matrices of the input bipartite graphs  $G_1$  and  $G_2$  i.e,  $BA_p = BA_1 \otimes BA_2$ . Figure 2 shows an example of bipartite graph product both from the viewpoint of both graph and matrix.

#### 4 Ramanujan Bipartite Graph Product Framework

The connectivity between neurons in a layer  $L$  of a sparse neural network can be captured using a bipartite graph  $G$ , where left/right neurons in  $L$  corresponds to left/right vertices in  $G$ , and the connections between left and right neurons in  $L$  corresponds to undirected edges between left and right vertices in  $G$ . The core idea in RBGP (Ramanujan Bipartite Graph Product) framework is to express  $G$  as a bipartite graph product of Ramanujan bipartite graphs i.e ( $G = G_1 \otimes_b \dots \otimes_b G_K$ ), where  $K$  is the number of base graphs. In the rest of the section, we show how expressing connectivity of a layer using bipartite graph products leads to sparse neural networks that have structured sparsity, good connectivity, and memory efficiency.

**Structured sparsity.** In bipartite graph product ( $G_p = G_1 \otimes_b G_2$ ), the biadjacency matrix of  $G_p$  is equal to the Tensor product( $\otimes$ ) of the biadjacency matrices of  $G_1$  and  $G_2$  i.e,  $BA_p = BA_1 \otimes BA_2$ . And in Tensor product,  $BA_p$  is constructed by replacing each non zero element in  $BA_1$  with  $BA_2$  matrix, and each zero element in  $BA_1$  with zero matrix of size  $BA_2$ . As  $BA_2$  is repeated,  $BA_p$  will have CBS (Cloned Block Sparse) sparsity pattern with block size equal to the size of  $BA_2$  or  $(|G_2.U|, |G_2.V|)$ . Figure 2 shows an example of bipartite graph product, where the biadjacency matrix of the product graph has CBS pattern with block size (2, 2). Additionally, when  $G_1$  is a biregular bipartite graph,  $BA_p$  will have CUBS (Cloned Uniform Block Sparse) sparsity pattern as  $BA_1$  will have equal number of elements in all rows, and all columns. In RBGP framework, the bipartite graph  $G$  of a layer  $L$  in the neural network is constructed by performing a series of  $(K - 1)$  bipartite graph products on  $K$  base biregular bipartite graphs ( $G = G_1 \otimes_b \dots \otimes_b G_K$ ) that are Ramanujan. Bipartite graph  $G$  can be rewritten as  $G = G_1 \otimes_b CG_2$ , where  $CG_2 = (G_2 \otimes_b \dots \otimes_b G_K)$ . As  $G_1$  is a biregular bipartite graph,  $BA$  (biadjacency matrix of  $G$ ) will have CUBS sparsity pattern with block size  $(\pi_{i=2}^{i=K} |G_i.U|, \pi_{i=2}^{i=K} |G_i.V|)$ . Going deeper, as  $CG_i = (G_i \otimes_b CG_{(i+1)})$ , and also as all the base graphs are biregular,  $BA$  will have RCUBS (Recursive Cloned Uniform Block Sparse) sparsity pattern with  $(K - 1)$  blocking levels  $B_1 \dots B_{(K-1)}$ , where  $B_j = (\pi_{i=j+1}^{i=K} |G_i.U|, \pi_{i=j+1}^{i=K} |G_i.V|)$ . Figure 3 shows an example bipartite graph generated using RBGP framework that uses four base graphs and has three block sizes (16, 16), (8, 8), and (2, 2).

**Memory efficiency.** A sparse neural network can be efficiently stored by only storing the information related to the connections that are present in the sparse layers. For a sparse layer  $L$  and its associated bipartite graph  $G$ ,  $|E(G)|$  memory is required for storing the parameters corresponding to connections, and another  $|E(G)|$  memory is required for storing connectivity information in the form of adjacency list of  $G$ . Thus a total of  $2 \times |E(G)|$  memory is required for storing the information of a layer in a sparse neural network. But in a RBGP sparse neural network, the memory requirement can be reduced by reducing the memory required for storing connectivity information. In RBGP sparse neural network, as  $G$  is constructed using  $K$  base bipartite graphs ( $G = G_1 \otimes_b \dots \otimes_b G_K$ ), the connectivity information of  $G$  can be reduced from  $E(G)(\prod_{i=1}^{i=K} |E(G_i)|)$  to  $\sum_{i=1}^{i=K} |E(G_i)|$ , by only storing the connectivity information of the individual base graphs. For example, the bipartite

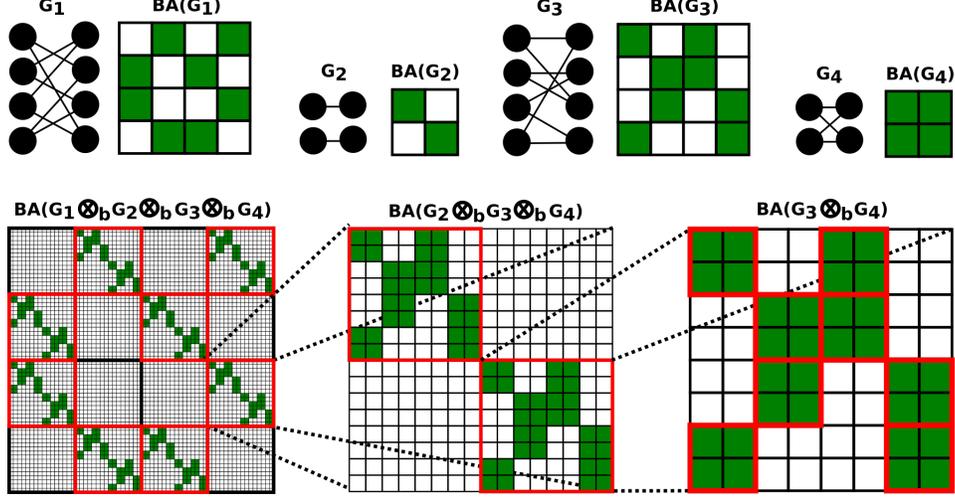


Figure 3: Biadjacency matrix  $BA$  of a bipartite graph generated using RBGP framework.  $BA$  has RCUBS(Recursive Cloned Uniform Block Sparse) sparsity pattern with three blocking levels  $(16, 16)$ ,  $(8, 8)$  and  $(2, 2)$

graph  $G$  generated using RBGP framework in Figure 3 has 512 edges  $(8 \times 2 \times 8 \times 4)$ , but it only requires storing 22 edges  $(8 + 2 + 8 + 4)$  from the base graphs to construct the connectivity information of  $G$ , thus leading to a 23x reduction in memory requirement for storing the connectivity information when compared to a random bipartite graph with same number of edges as  $G$ .

**Good connectivity.** Connectivity in a sparse neural network is key for ensuring good flow of information. It is well known [1] that connectivity of the graph is characterized by the *spectral gap* between the largest and second largest eigenvalue (in absolute terms) of the adjacency matrix. In this section, we show that the spectral gap for the block sparse graph we construct using graph products, are optimal for any level of sparsity, for large graphs.

For a  $d$ -regular bipartite graph the largest eigenvalue in absolute value is  $d$  and  $-d$ . The next largest eigenvalue is considered as the second largest eigenvalue  $\lambda_2$ . The spectral gap is  $d - \lambda_2$  and larger this quantity, the better connected the graph. Suppose the bipartite graph has  $n$  vertices on both sides, the degree  $d$  is  $\alpha n$  where  $\alpha$  is the fractional sparsity. For a given value of  $d$ , the best possible spectral gap of  $d - 2\sqrt{d} - 1$  is achieved by Ramanujan Graphs. We construct block sparse graphs using graph products of smaller Ramanujan Graphs and show below that this construction has similar spectral gap as  $n \rightarrow \infty$ . For simplicity we consider the case where the bipartite graph  $G$  is the graph product of  $G_1, G_2$  which are bipartite graphs with  $n$  vertices on each sides and degree  $d = \alpha n$ . Note that  $G$  has degree  $d^2$  and sparsity  $1 - (1 - \alpha)^2$ .

**Theorem 1.** Let  $G = G_1 \otimes_b G_2$  where  $G_i$  are bipartite graphs with  $n$  vertices on each sides and degree  $d = \alpha n$ . Then for any fixed level of sparsity  $\alpha$ ,

$$\frac{IdealSpectralGap_{d^2}}{SpectralGap(G)} \rightarrow 1 \quad \text{as } n \rightarrow \infty \quad (1)$$

where  $IdealSpectralGap_{d^2} = d^2 - 2\sqrt{d^2} - 1$  is the best possible spectral gap for  $d^2$ -regular graphs and  $SpectralGap(G)$  is the spectral gap of the block sparse graph  $G$  that we construct.

*Proof.* The biadjacency matrix of  $G$  is the tensor product of biadjacency matrices of  $G_1, G_2$ . Hence the eigenvalues of the biadjacency matrix is the product of eigenvalues of biadjacency matrices of  $G_1, G_2$ . Since  $G_1, G_2$  are Ramanujan Graphs, their second largest eigenvalue is  $2\sqrt{d} - 1$ . Hence second largest eigenvalue of  $G$  is  $\lambda_2(G) = d \times 2\sqrt{d} - 1$ . The ideal value of second largest eigenvalue for graphs of degree  $d^2$  is  $2\sqrt{d^2} - 1$ . Hence Equation 1, becomes

$$\frac{d^2 - 2\sqrt{d^2} - 1}{d^2 - 2d\sqrt{d} - 1} = \frac{1 - 2\sqrt{1/d^2} - 1/d^4}{1 - 2\sqrt{1/d} - 1/d^2}.$$

Hence for any fixed level of sparsity  $\alpha$ ,  $n \rightarrow \infty$  (large matrices),  $d \rightarrow \infty$ , the LHS of Equation 1  $\rightarrow 1$ .  $\square$

## 5 RBGP framework for GPU

A GPU is fundamentally a many core architecture with thousands of cores, and have multiple memory subsystems(DRAM, L2 cache, L1 cache/shared memory, and registers) with data access times decreasing in that order.The reason for having many memory subsystems is to feed data into cores at a higher rate by avoiding data accesses to slower memory say DRAM, when data is already available on faster memory say L2 cache. On GPU, a computational task can have good runtime efficiency, if it can avoid idling of cores by maximizing memory accesses from faster memories through data reuse. Sparse neural networks with unstructured sparsity pattern offers limited data reuse due to irregular memory access patterns, and thus has poor runtime performance on GPU. The only way for sparse neural networks to achieve good runtime performance on GPU is by embracing structured sparsity patterns. In this section, using our proposed RBGP framework, we design RBGP4 structured sparsity pattern to effectively use memory subsystems on GPU by facilitating data reuse, and achieve good runtime performance for RBGP4 sparse neural networks.

**RBGP4 sparsity pattern.** In RBGP framework, bipartite graph  $G(G = G_1 \otimes_b \dots \otimes_b G_K)$  corresponding to a layer in the sparse neural network is configured by the number of base graphs( $K$ ), and for each base graph  $G_i$ , it's type(sparse or complete). RBGP4 sparsity pattern corresponds to a specific configuration, where  $G$  is constructed using four base Ramanujan bipartite graphs ( $G = G_o \otimes_p G_r \otimes_p G_i \otimes_p G_b$ ), with graphs  $G_o$  and  $G_i$  being sparse, and  $G_r$  and  $G_b$  being complete bipartite graphs. Figure 1 shows an example of RBGP4 sparsity pattern, where  $G_o$  and  $G_i$  are 50% sparse, and  $G_r$  and  $G_b$  are (2,1) and (2,2) complete bipartite graphs respectively.

**GPU Implementation.** Compute in each layer of an RBGP4 sparse neural network is composed of RBGP4MM(Multiplication of a sparse matrix  $W_s$  with RBGP4 sparsity pattern, and a dense matrix  $I$ ) operation ( $O = W_s \times I$ ), where  $W_s$ ,  $I$ , and  $O$ , corresponds to sparse weight matrix, batched input activations, and batched output activations respectively. We use tiling approach for efficiently processing RBGP4MM operation. In tiling approach, matrices are divided into tiles, and  $OT$ (a tile in  $O$ ) is computed in steps, where each step is comprised of matrix multiplication of  $WT_s$ (a sparse tile in  $W$ ) with  $IT$ (a dense tile in  $I$ ) i.e,  $OT+ = WT_s \times IT$ . For RBGP4MM, we set tile size in  $W_s$  is set to be  $(|G_t.U|, |G_t.V|)$ , where  $G_t = (G_r \otimes_b G_i \otimes_b G_b)$ . On GPU, we associate computation of  $OT$  to a thread block, and with in a thread block, each thread maps to a strided 2D grid of element blocks in  $OT$ , with  $|G_r.U|$  number of strides and  $|G_b.U|$  element block size in row dimension. We exploit the data reuse offered by RBGP4 sparsity pattern and make efficient use of memory hierarchy on GPU, by first loading tiles  $WT_s$  and  $IT$  into shared memory in each step of  $OT$ , and each thread loads it's share of data into registers from shared memory before performing the computation. Figure 1 shows an example of using tiling approach for RBGP4MM operation on GPU. A more detailed GPU algorithm can be found in Appendix.

**Why RBGP4 ?** RBGP4 sparsity pattern ( $G = G_o \otimes_p G_r \otimes_p G_i \otimes_p G_b$ ) is designed to achieve runtime efficiency for SDMM operation ( $O = W_s \times I$ ) on GPU. Towards that, all the four base graphs  $G_o, G_r, G_i$ , and  $G_b$  in RBGP4 sparsity pattern have a specific role to play.

The role of  $G_o$  is to reduce the number of steps required to process  $OT$ (a tile in  $O$ ) by inducing sparsity at the tile level in  $W_s$ . Performing bipartite product to the left of  $G_t$  with  $G_o$  i.e, ( $G = G_o \otimes_b G_t$ ) results in block sparsity pattern in  $W_s$  with block size  $(|G_t.U|, |G_t.V|)$ . As we set tile size in  $W_s$  to be the block size, sparsity is induced at the tile/block level in  $W_s$ , which inturn reduces the number of steps for processing  $CT$  by skipping computation corresponding to zero tiles in  $W_s$ . For example in Figure 1, we can see that the number of steps required to compute  $OT$  is reduced from two to one, as  $W_s$  has only two non zero tiles out of four tiles due to 50% sparsity in  $G_o$ .

The role of graphs  $G_r$  and  $G_b$  in RBGP4 sparsity pattern is to maximize data reuse from registers in GPU threads by inducing row repetition in  $WT_s$ (a tile in  $W_s$ ). In row repetition, rows are divided into groups of equal size, where all the rows in a group have non zeros at the same locations. Having row repetition pattern in  $WT_s$  implies that all the rows in a group will have same memory access patterns into  $IT$ , and thus allows for reuse of data from  $WT_s$  and  $IT$ . Performing bipartite graph product to the left and right of  $G_i$  with complete graphs  $G_r$  and  $G_b$  respectively i.e, ( $G_t = G_r \otimes G_i \otimes G_b$ ) results in row repetition in  $WT_s$  with  $|G_i.U|$  groups, and  $|G_r.U| \times |G_b.U|$  rows in each group. For example in Figure 1, we can see that as  $G_r$  and  $G_b$  are complete bipartite graphs with (2, 1) and (2, 2) sizes, the sparsity pattern of  $WT_s$ , has row repetition pattern with 4 rows. In computation associated

Sparsity	Pattern	VGG19				WideResnet-40-4			
		CF10	CF100	Mem	Time	CF10	CF100	Mem	Time
00.00	Dense	93.14	70.64	77.39	22	95.01	77.20	34.10	40
50.00	Unstructured	92.67	70.31	77.39	165	95.42	77.92	34.10	241
	Block	92.45	70.75	41.12	94	95.49	77.52	18.12	165
	RBGP4	92.58	70.48	38.76	20	95.34	78.27	17.13	32
75.00	Unstructured	91.99	69.32	38.71	86	95.10	76.89	17.05	135
	Block	91.93	68.72	20.57	48	94.92	76.50	9.07	85
	RBGP4	91.99	68.34	19.40	13	94.72	76.80	8.57	20
87.50	Unstructured	90.88	65.41	19.37	79	94.48	75.21	8.53	102
	Block	90.62	65.37	10.30	25	94.56	74.55	4.54	45
	RBGP4	90.48	65.39	9.72	8	94.38	75.25	4.30	16
93.75	Unstructured	90.01	62.33	9.70	50	93.57	73.09	4.27	69
	Block	89.40	62.90	5.16	14	93.55	71.86	2.27	26
	RBGP4	89.32	62.79	4.88	6	93.53	72.44	2.16	14

Table 1: Image classification on CIFAR10 (CF10) and CIFAR100 (CF100) datasets using VGG19 and WideResnet-40-4 networks. Models are trained using predefined approach with unstructured, block, and RBGP4 sparsity patterns. For block pattern, we set block size to be (4, 4). Memory (Mem) is given in MB, and time is given in milliseconds for one forward pass in training.

with thread  $T_1$  in  $O$ , rows (1, 2, 5, 6) have same non zero pattern in  $WT_s$ , and this allows us to load two  $2 \times 2$  blocks from  $WT_s$  and one  $2 \times 2$  block from  $IT$  into register blocks  $RegW$  and  $RegI$  respectively and reuse each elements from  $RegW$  and  $RegI$  for 2 and 4 times respectively.

The role of  $G_i$  in RBGP4 sparsity pattern is to allow  $W_s$  to have any level of sparsity even when the tile size in  $W_s$  is big. When the tile size in  $W_s$  is relatively large when compared to the size of  $W_s$ , it is not possible to obtain desired level of sparsity if a non zero tile in  $W_s$  is dense. For example, if a tile in  $W_s$  is of size (64, 64), and  $W_s$  is of size (128, 64), only by allowing tiles in  $W_s$  to be sparse, can sparsity greater than 50% can be obtained. Bipartite graph  $G_t$  corresponds to sparsity pattern of  $WT_s$ , and in RBGP4 sparsity pattern  $G_t = (G_r \otimes_b G_i \otimes_b G_b)$ . As  $G_r$  and  $G_b$  are dense/complete,  $G_i$  has to be sparse to achieve a desired level of sparsity in  $W_s$ .

## 6 Results

We study the effect of RBGP4 sparsity pattern on model accuracy for the task of image classification and compare with unstructured and block structured sparsity patterns. Further more, we study the effect of changing configuration of base graphs in RBGP4 sparsity pattern on runtime. We perform all our experiments on V100 GPU, where we benchmark unstructured and block sparsity patterns using cuSparse library, and dense pattern using cuBLAS library from NVIDIA.

**Image classification benchmark.** In this benchmark, we perform the image classification task on CIFAR dataset using VGG19[31] as adapted by Liu et al. [20], and WideResnet-40-4[37] networks. To train the models, we use predefined approach, where the mask(choice of connections) is chosen a priori to the training process. As a sparse neural network has less number of parameters, we first train the dense model and guide the sparse neural network using knowledge distillation [13]. For all our experiments, we incorporate equal amount of sparsity in all layers, except for the first layer connected to input and the final classifier layer. For the optimizer, we use SGD optimizer with momentum of 0.9 and weight decay of  $1e-4$ . VGG19/WideResnet-40-4 model is trained for 160/200 epochs with batch size of 256/128. Initial learning rate is set to 0.1. For VGG19, learning rate is multiplied by 0.1 at epochs 60, 120, and 160. And for WideResnet-40-4, learning rate is multiplied by 0.2 at epochs

60,120, and 160. From Table 1, we can see that RBGP4 is as accurate as unstructured and block sparsity patterns, but takes 2x less memory and is 5-9x faster when compared to unstructured, and is 2-5x faster when compared to block sparsity pattern.

**RBGP4 runtime characteristics.** RBGP4 sparse matrix  $W_s$  of a given size and sparsity can be obtained in multiple ways by varying the sizes of base graphs  $G_o, G_p, G_i, G_b$ , and sparsities of  $G_o$  and  $G_i$ . For example, setting sparsities of  $(G_o, G_i)$  to either  $(0, 75\%)$  or  $(50\%, 50\%)$  leads to 75% sparsity in  $W_s$ , and setting sizes of base graphs to either  $((8, 4), (1, 1), (8, 4), (1, 1))$  or  $((8, 4), (2, 2), (4, 2), (1, 1))$  leads to  $W_s$  of size  $(64, 16)$ . In this section, we study the effect of RBGP4 configuration on runtime of SDMM operation ( $O = W_s \times I$ ). For all our experiments, we set sizes of matrices  $O, W_s$ , and  $I$  to be  $4096 \times 4096$ .

*Sparsity distribution* : In RBGP4 sparsity pattern, sparsity is solely due to presence of sparse graphs  $G_o$  and  $G_i$ , as  $G_r$  and  $G_b$  are dense or complete graphs. We run experiments with 75%, 87.5%, and 93.75% sparsity amounts distributed between  $G_o$  and  $G_i$ , while keeping sizes of  $G_o, G_r, G_i, G_b$  fixed to  $(32, 128), (4, 1), (32, 32), (1, 1)$ . From Table 2, we can see that for a given sparsity, as sparsity of  $G_o$  increases, the runtime decreases. This is because sparsity in  $G_o$  incorporates sparsity at the tile level, and this reduces runtime due to skipping of computation and memory loads associated with zero tiles. For dense case (0% sparsity), we use cuBLAS library from NVIDIA.

*Row repetition* : In row repetition, matrix  $W_s$  can be divided into row groups of equal size, where all the rows in a row group have non zeros exactly at the same locations. Having row repetitions allows us to effectively reuse data from  $I$  as rows have same non zero pattern.  $G_r$  and  $G_b$  in RBGP4 introduces  $|G_r.U| \times |G_b.U|$  amount of row repetition in  $W_s$ . We run experiments with 1, 2, and 4 repetition amounts, while keeping size of  $G_i (G_r \otimes G_i \otimes G_b)$  fixed at  $(128, 32)$ , and sparsity of  $G_o$  at 50%. From Table 3, we can see that increasing the size of  $G_r$  or  $G_b$  or both leads to improved runtime performance as repetition amount increases.

Sp(G)%	Sp( $G_o$ ) %	Sp( $G_i$ ) %	Time(ms)
0	0	0	11.2 (1x)
75.00	0.00	75.00	5.64 (2x)
	50.00	50.00	4.44 (2.5x)
87.50	0.00	87.50	4.31 (2.6x)
	50.00	75.00	2.74 (4.1x)
	75.00	50.00	2.29 (4.9x)
93.75	0.00	93.75	3.76 (3x)
	50.00	87.50	1.93 (5.8x)
	75.00	75.00	1.44 (7.8x)
	87.50	50.00	1.22 (9.2x)

Table 2: Effect of varying sparsities of sparse graphs  $G_o$  and  $G_i$  in RBGP4 sparsity pattern on runtime.

Sizes		Time(ms) for Sp(G)%		
$G_r$	$G_b$	75.00	87.50	93.75
(1,1)	(1,1)	7.07	3.91	2.45
(2,1)	(1,1)	4.89	3.02	1.97
(4,1)	(1,1)	4.47	2.75	1.92
(1,1)	(2,1)	4.85	3.01	2.03
(1,1)	(4,1)	4.47	2.84	2.02
(2,1)	(2,1)	4.41	2.75	1.98

Table 3: Effect of varying sizes of complete graphs  $G_r$  and  $G_b$  in RBGP4 sparsity pattern on runtime.

## 7 Conclusion

We used ideas from extremal graph theory and combinatorics to make sparse neural networks runtime efficient. Ramanujan graphs which gives the optimal connectivity for a given level of sparsity are used to model connections in a neural network layer. Furthermore, we obtain structured block sparsity by using products of Ramanujan graphs. We prove that the product graph also has the optimal connectivity for large matrices. For the specific case of GPUs, we describe how the block sparsity can be efficiently implemented in hardware, by exploiting the memory hierarchy through data reuse. Benchmarks of this implementation is shown to give significant runtime improvements. Similar ideas could be used for generating structured sparsity patterns that results in runtime efficient implementations in other hardware as well. For the future work, generating combinatorial structured sparsity patterns like RBGP4 during the training process could lead to more accurate models as structure is induced in a gradual manner.

## References

- [1] Alon, N.: Eigenvalues and expanders. *Combinatorica* **6**(2), 83–96 (1986)
- [2] Bellec, G., Kappel, D., Maass, W., Legenstein, R.: Deep rewiring: Training very sparse deep networks. arXiv preprint arXiv:1711.05136 (2017)
- [3] Bilu, Y., Linial, N.: Lifts, discrepancy and nearly optimal spectral gap. *Combinatorica* **26**(5), 495–519 (Oct 2006). <https://doi.org/10.1007/s00493-006-0029-7>, <https://doi.org/10.1007/s00493-006-0029-7>
- [4] Cao, S., Zhang, C., Yao, Z., Xiao, W., Nie, L., Zhan, D., Liu, Y., Wu, M., Zhang, L.: Efficient and effective sparse lstm on fpga with bank-balanced sparsity. In: *Proceedings of the 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. pp. 63–72 (2019)
- [5] Chung, F.R.K.: *Spectral Graph Theory*. American Mathematical Society (1997)
- [6] Dettmers, T., Zettlemoyer, L.: Sparse networks from scratch: Faster training without losing performance. *CoRR* **abs/1907.04840** (2019), <http://arxiv.org/abs/1907.04840>
- [7] Dey, S., Huang, K.W., Beerel, P.A., Chugg, K.M.: Pre-defined sparse neural networks with hardware acceleration. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems* **9**(2), 332–345 (2019)
- [8] Frankle, J., Carbin, M.: The lottery ticket hypothesis: Finding sparse, trainable neural networks. arXiv preprint arXiv:1803.03635 (2018)
- [9] Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In: Bengio, Y., LeCun, Y. (eds.) *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings* (2016), <http://arxiv.org/abs/1510.00149>
- [10] Han, S., Pool, J., Tran, J., Dally, W.: Learning both weights and connections for efficient neural network. In: *Advances in neural information processing systems*. pp. 1135–1143 (2015)
- [11] Hassibi, B., Stork, D.G., Wolff, G.: Optimal brain surgeon: Extensions and performance comparisons. In: *Advances in neural information processing systems*. pp. 263–270 (1994)
- [12] He, Y., Zhang, X., Sun, J.: Channel pruning for accelerating very deep neural networks. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1389–1397 (2017)
- [13] Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. In: *NIPS Deep Learning and Representation Learning Workshop* (2015), <http://arxiv.org/abs/1503.02531>
- [14] Huang, Z., Wang, N.: Data-driven sparse structure selection for deep neural networks. In: *The European Conference on Computer Vision (ECCV)* (September 2018)
- [15] Kepner, J., Robinett, R.: Radix-net: Structured sparse matrices for deep neural networks. In: *2019 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. pp. 268–274. IEEE (2019)
- [16] LeCun, Y., Denker, J.S., Solla, S.A.: Optimal brain damage. In: *Advances in neural information processing systems*. pp. 598–605 (1990)
- [17] Lee, N., Ajanthan, T., Torr, P.H.S.: SNIP: single-shot network pruning based on connection sensitivity. *CoRR* **abs/1810.02340** (2018), <http://arxiv.org/abs/1810.02340>
- [18] Li, H., Kadav, A., Durdanovic, I., Samet, H., Graf, H.P.: Pruning filters for efficient convnets. arXiv preprint arXiv:1608.08710 (2016)
- [19] Liu, Z., Li, J., Shen, Z., Huang, G., Yan, S., Zhang, C.: Learning efficient convolutional networks through network slimming. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 2736–2744 (2017)
- [20] Liu, Z., Sun, M., Zhou, T., Huang, G., Darrell, T.: Rethinking the value of network pruning (2018)
- [21] Lubotzky, A., Phillips, R., Sarnak, P.: Ramanujan graphs. *Combinatorica* **8**(3), 261–277 (1988)
- [22] Luo, J.H., Wu, J., Lin, W.: Thinet: A filter level pruning method for deep neural network compression. In: *Proceedings of the IEEE international conference on computer vision*. pp. 5058–5066 (2017)

- [23] Mao, H., Han, S., Pool, J., Li, W., Liu, X., Wang, Y., Dally, W.J.: Exploring the granularity of sparsity in convolutional neural networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (July 2017)
- [24] Marcus, A.W., Spielman, D.A., Srivastava, N.: Interlacing families i: Bipartite ramanujan graphs of all degrees. *Annals of Mathematics* **182**(1), 307–325 (2015), <http://www.jstor.org/stable/24523004>
- [25] Mocanu, D.C., Mocanu, E., Stone, P., Nguyen, P.H., Gibescu, M., Liotta, A.: Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications* **9**(1), 1–12 (2018)
- [26] Molchanov, P., Tyree, S., Karras, T., Aila, T., Kautz, J.: Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440* (2016)
- [27] Mostafa, H., Wang, X.: Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. In: International Conference on Machine Learning. pp. 4646–4655 (2019)
- [28] Narang, S., Elsen, E., Damos, G., Sengupta, S.: Exploring sparsity in recurrent neural networks. *arXiv preprint arXiv:1704.05119* (2017)
- [29] Narang, S., Undersander, E., Damos, G.: Block-sparse recurrent neural networks. *arXiv preprint arXiv:1711.02782* (2017)
- [30] Prabhu, A., Varma, G., Nambodiri, A.: Deep expander networks: Efficient deep networks from graph theory. In: The European Conference on Computer Vision (ECCV) (September 2018)
- [31] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (2015)
- [32] Srinivas, S., Subramanya, A., Venkatesh Babu, R.: Training sparse neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. pp. 138–145 (2017)
- [33] Vooturi, D.T., Kothapalli, K.: Efficient sparse neural networks using regularized multi block sparsity pattern on a gpu. In: High Performance Computing and Data Analytics (HiPC) (December 2019)
- [34] Vooturi, D.T., Varma, G., Kothapalli, K.: Dynamic block sparse reparameterization of convolutional neural networks. In: The IEEE International Conference on Computer Vision (ICCV) Workshops (Oct 2019)
- [35] Wen, W., Wu, C., Wang, Y., Chen, Y., Li, H.: Learning structured sparsity in deep neural networks. In: Lee, D.D., Sugiyama, M., Luxburg, U.V., Guyon, I., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 29, pp. 2074–2082. Curran Associates, Inc. (2016), <http://papers.nips.cc/paper/6504-learning-structured-sparsity-in-deep-neural-networks.pdf>
- [36] Yu, R., Li, A., Chen, C.F., Lai, J.H., Morariu, V.I., Han, X., Gao, M., Lin, C.Y., Davis, L.S.: Nisp: Pruning networks using neuron importance score propagation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
- [37] Zagoruyko, S., Komodakis, N.: Wide residual networks. In: Richard C. Wilson, E.R.H., Smith, W.A.P. (eds.) *Proceedings of the British Machine Vision Conference (BMVC)*. pp. 87.1–87.12. BMVA Press (September 2016). <https://doi.org/10.5244/C.30.87>, <https://dx.doi.org/10.5244/C.30.87>

## 8 Appendix

### 8.1 Ramanujan Bipartite Graph Generation

A construction for Ramanujan Bipartite graph(RBG) was given by Bilu et al. [3]. The proof that this construction obtains the optimal eigenvalue gap was given by Marcus et al. [24]. We use algorithms(graph lifts) derived from these construction to generate Ramanujan Bipartite Graphs for a given sparsity.

**2-lift operation:** A 2-lift is an operation applied on a graph  $G$  to produce a bigger graph  $G_L$  that is twice as big as  $G$  in both vertices and edges. In the 2-lift operation, a clone graph  $G^c$  is first created and the vertex set of  $G_L$  is set to be the union of vertex sets of  $G$  and  $G^c$  i.e,  $V(G_L) = V(G) \cup V(G^c)$ . The edge set of  $G_L$  i.e,  $E(G_L)$  is then constructed in the following way: For an edge  $(u, v) \in G$ , and it's corresponding clone edge  $(u^c, v^c) \in G^c$ , either the identity edge pair  $\{(u, v), (u^c, v^c)\}$  or the crossover edge pair  $\{(u, v^c), (u^c, v)\}$  is chosen at random and added to  $E(G_L)$ . Figure 4 shows an example of 2-lift operation.

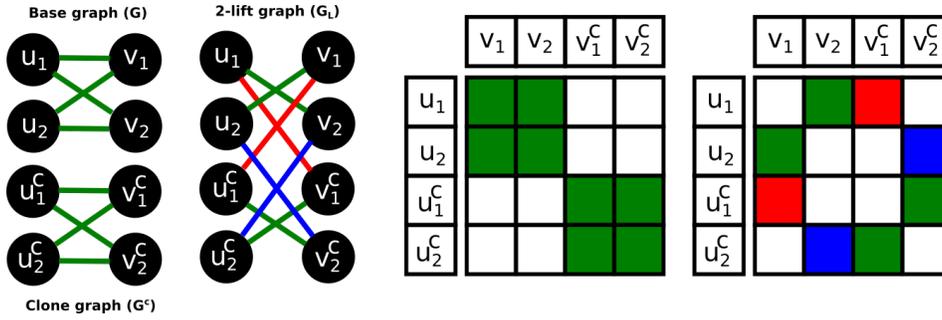


Figure 4: 2-lift operation on graph  $G$ . Clone graph  $G^c$  is first created and edges  $(u_1, v_1)$  and  $(u_2, v_2)$  are randomly chosen to cross over with the corresponding edges  $(u_1^c, v_1^c)$  and  $(u_2^c, v_2^c)$  respectively in the clone graph.

**Generating sparse biregular bipartite graph:** A 2-lift operation when applied on a biregular bipartite graph also results in a biregular bipartite graph that is twice as big with same left and right degrees. A biregular graph  $G(U, V, E)$  with sparsity  $(1.0 - |E(G)|/(|G.U| \times |G.V|))$   $sp$ , can be generated by repeatedly applying  $\log_2(1/(1 - sp))$  2-lift operations on a complete bipartite graph with  $(1 - sp) \times |G.U|$  left and  $(1 - sp) \times |G.V|$  right vertices.

**Generating RBG graph:** A Ramanujan bipartite graph is first a biregular bipartite graph with an additional constraint on second largest eigenvalue of the adjacency matrix of the graph. To generate an RBG graph, we sample sparse biregular bipartite graphs generated using 2-lift operations until the sampled graph is Ramanujan. We found that an RBG graph with sizes in the order of thousands can be generated in the order of minutes. For a layer in RBGP sparse neural network, the base Ramanujan graphs are generated only once before training and hence sampling approach is not a bottleneck.

### 8.2 Pseudo code for RBGP4MM operation on GPU

Computation in each layer of a sparse neural network is an SDMM(Multiplication of a sparse matrix with a dense matrix) operation ( $C = A_s \times B$ ). RBGP4MM is an SDMM operation where  $A_s$  has RBGP4 sparsity pattern. Algorithm 1 describes the pseudo code for RBGP4MM operation on a GPU. As RBGP4 sparsity pattern has equal number of non zero elements in each row, non zero elements in  $A_s$  can be stored using *data* array of size  $(A_s.rows, (1 - sp) \times A_s.columns)$ , and the index information of  $A_s$  is captured by storing adjacency lists of base bipartite graphs.

---

**Algorithm 1** GPU algorithm for RBGP4MM( $C = A_s \times B$ ) operation using tiling approach. Tile sizes for  $A_s, B$ , and  $C$  are chosen to be  $(TM, TK), (TK, TN)$ , and  $(TM, TN)$  respectively. On GPU, each tile in  $C$  is mapped to a thread block, and each thread in the thread block is mapped to a group of  $(RM \times BM \times RN \times BN)$  number of elements in a tile of  $C$ . Variables TM,TK,RM,RK,BM,BK are set based on RBGP4 configuration( $G = G_o \otimes_b G_r \otimes_b G_i \otimes_b G_b$ ) of  $A_s$ .

---

```

1: function LBFM(matrix, (bi, bj), (BH, BW))                                ▷ Load Block From Matrix
2:   block[BH][BW]
3:   for i in [0, BH) do
4:     for j in [0, BW) do
5:       block[i][j] = matrix[bi * BH + i][bj * BW + j]
6:     end for
7:   end for
8:   return block
9: end function

10:  $G_t = G_r \otimes_b G_i \otimes_b G_b$ 
11:  $TM, TK = |G_t.U|, |G_t.V|$                                 ▷ Number of left and right vertices of bipartite graph  $G_t$ 
12:  $RM, RK = |G_r.U|, |G_r.V|$ 
13:  $BM, BK = |G_b.U|, |G_b.V|$ 
14: gridBlockDim = (C.rows/TM, C.cols/TN)                    ▷ 2D grid block
15: threadBlockDim = (TM/ $(RM \times BM)$ , TN/ $(RN \times BN)$ )        ▷ 2D thread block
16: for (tbm, tbn) in [(0, 0) : gridBlockDim) do                ▷ Mapped to thread blocks
17:   for (thm, thn) in [(0, 0) : threadBlockDim) do          ▷ Mapped to threads
18:     Areg[RM][BM][BK]                                       ▷ Registers
19:     Breg[RN][BK][BN]                                       ▷ Registers
20:     Creg[RM][RN][BM][BN]                                   ▷ Registers
21:     for outk in [0,  $G_o.d_i$ ) do                                ▷  $G_o.d_i$  is left degree of biregular bipartite graph  $G_o$ 
22:       oind =  $G_o.adj\_list[tbm][outk]$ 
23:       Atile = LBFM(A_s.data, (tbm, outk), (TM,  $G_t.d_i$ )) ▷ DRAM to shared memory
24:       Btile = LBFM(B, (oind, tbn), (TK, TN)) ▷ DRAM to shared memory(shMem)
25:       __syncthreads()
26:       for rk, ink in [0, RK)  $\times$  [0,  $G_i.d_i$ ) do
27:         for rm in [0 : RM) do
28:           bm =  $rm \times |G_i.U| + thm$ 
29:           bk =  $rk \times G_i.d_i + ink$ 
30:           Areg[rm] = LBFM(Atile, (bm, bk), (BM, BK)) ▷ ShMem to registers
31:         end for
32:         for rn in [0, RN) do
33:           bk =  $rk \times |G_i.V| + G_i.adj\_list[thm][ink]$ 
34:           bn =  $rn \times TN / (RN \times BN) + thn$ 
35:           Breg[rn] = LBFM(Btile, (bk, bn), (BK, BN)) ▷ ShMem to registers
36:         end for
37:         for rm, rn in [0, RM)  $\times$  [0, RN) do
38:           Creg[rm][rn] += Areg[rm]  $\times$  Breg[rn]          ▷ Computation
39:         end for
40:       __syncthreads()
41:     end for
42:   end for
43:   for rm, rn in [0, RM)  $\times$  [0, RN) do
44:     for m, n in [0 : BM)  $\times$  [0 : BN) do
45:       row =  $tbm \times TM + rm \times (TM/RM) + thm \times BM + m$ 
46:       col =  $tbn \times TN + rn \times (TN/RN) + thn \times BN + n$ 
47:       C[row][col] += Creg[rm][rn][m][n]
48:     end for
49:   end for
50: end for
51: end for

```

---