

# Sparse DNN Results with the MATLAB interface to SuiteSparse:GraphBLAS

Tim Davis

Sept 2, 2019 (updated Feb 20, 2020)

The tables below report the results for the 12 sparse deep neural network problems. Problems 1-3 use 1024 neurons, 4-6 use 4,096 neurons, 7-9 use 16K neurons, and the last use 64K neurons. Each group of three uses 120, 480, and 1920 layers, respectively.

## 1 The MATLAB code

The MATLAB code (with GraphBLAS) is very simple, even simpler than the MATLAB reference implementation. All MATLAB variables are GraphBLAS GrB objects. In this case, they represent sparse matrices with the GrB\_FP32 floating point type, stored in CSR format (by row).

```
function Y = dnn_gb (W, bias, Y0)
Y = Y0 ;
for i=1:length(W)
    Y = GrB.select ('>0', GrB.mxm ('+.*', Y * W {i}, bias {i})) ;
    M = Y > 32 ;
    if (nnz (M) > 0)
        Y (M) = 32 ;
    end
end
end
```

For comparison, here is the MATLAB reference implementation at <http://graphchallenge.org>. It is about 60x to 70x slower than the two methods using GraphBLAS. Applying the bias is more complex than the GraphBLAS code:

```

function Y = dnn_matlab (W, bias, Y0)
Y = Y0 ;
for i=1:length(W)
    % Propagate through layer.
    Z = Y * W {i} ;
    % Apply bias to non-zero entries.
    Y = Z + (double(logical(Z)) .* bias {i}) ;
    % Threshold negative values.
    Y (Y < 0) = 0 ;
    % Threshold maximum values.
    Y (Y > 32) = 32 ;
end

```

The code to convert from MATLAB sparse matrices to GraphBLAS GrB objects is shown below. This time is not included since the problem could have been read in as GraphBLAS matrices to begin with. In any case, the conversion time is trivial. Problem 12 is converted from MATLAB to GraphBLAS in 30.4 seconds, or about 1% of the time to solve the DNN with 40 threads. Problem 1 is converted in 0.3 seconds.

```

function [W, bias, Y0] = dnn_mat2gb (W, bias, Y0)
n = size (Y0, 2) ;
Y0 = GrB (Y0, 'single') ;
for i=1:length(W)
    W {i} = GrB (W {i}, 'single') ;
    bias {i} = GrB.build (1:n, 1:n, bias {i}, n, n, '+', 'single') ;
end

```

## 2 The C code

The C version in LAGraph, minus error checking and with a few other trivial simplifications, is shown below. Is straight-forward but more complex than either the pure MATLAB version (`dnn_matlab.m`) or the MATLAB+GraphBLAS version (`dnn_gb.m`).

```
#include "LAGraph.h"
void ymax32 (float *z, const float *x)
{
    (*z) = fminf ((*x), (float) 32.0) ;
}
GrB_Info LAGraph_dnn // returns GrB_SUCCESS if successful
(
    GrB_Matrix *Yhandle, // Y, created on output
    GrB_Matrix *W, // W [0..nlayers-1], each nneurons-by-nneurons
    GrB_Matrix *Bias, // Bias [0..nlayers-1], diagonal nneurons-by-nneurons
    int nlayers, // # of layers
    GrB_Matrix Y0 // input features: nfeatures-by-nneurons
)
{
    GrB_Matrix Y = NULL, M = NULL ;
    GrB_Index nfeatures, nneurons ;
    GrB_Matrix_nrows (&nfeatures, Y0) ;
    GrB_Matrix_ncols (&nneurons, Y0) ;
    GrB_Matrix_new (&Y, type, nfeatures, nneurons) ;
    GrB_Matrix_new (&M, GrB_BOOL, nfeatures, nneurons) ;
    GrB_UnaryOp Ymax32 ;
    GrB_UnaryOp_new (&Ymax32, ymax32, GrB_FP32, GrB_FP32) ;
    for (int layer = 0 ; layer < nlayers ; layer++)
    {
        // Y = Y * W [layer], using the conventional PLUS_TIMES semiring
        GrB_mxm (Y, NULL, NULL, GxB__PLUS_TIMES_FP32,
            ((layer == 0) ? Y0 : Y), W [layer], NULL) ;
        // Y = Y * Bias [layer], using the PLUS_PLUS semiring.
        GrB_mxm (Y, NULL, NULL, GxB__PLUS_PLUS_FP32, Y, Bias [layer], NULL) ;
        // delete entries from Y: keep only those entries greater than zero
        GxB_select (Y, NULL, NULL, GxB_GT_ZERO, Y, NULL, NULL) ;
        // threshold maximum values: Y (Y > 32) = 32
        GrB_apply (Y, NULL, NULL, Ymax32, Y, NULL) ;
    }
    GrB_free (&M) ;
    (*Yhandle) = Y ;
    return (GrB_SUCCESS) ;
}
```

### 3 Run time results

Run time in seconds on an Intel Xeon E5-2698v4 @ 2.2GHz, with 20 hardware cores and 256GB of RAM, using the GCC 5.4.0 compiler, and Ubuntu 16.04. The fastest time in bold, for one and 40 threads. Lower is better. MATLAB R2018a was used for both v3.1.0 and v3.2.0.

#### 3.1 GraphBLAS v3.1.0, August 2019

Prob	one thread			40 threads		
	MATLAB	LAGraph	M/L	MATLAB	LAGraph	M/L
1	<b>24</b>	24.2	0.97	3	<b>2.4</b>	1.17
2	<b>68</b>	68.2	0.99	9	<b>4.5</b>	2.07
3	<b>242</b>	243.3	1.00	34	<b>16.4</b>	2.09
4	<b>98</b>	108.1	0.90	10	<b>9.3</b>	1.07
5	<b>293</b>	330.1	0.89	<b>31</b>	30.7	1.00
6	<b>1076</b>	1222.5	0.88	<b>117</b>	117.6	0.99
7	766	<b>741.4</b>	1.03	58	<b>51.0</b>	1.15
8	2684	<b>2552.1</b>	1.05	201	<b>175.0</b>	1.15
9	10381	<b>9783.1</b>	1.06	783	<b>690.1</b>	1.13
10	<b>3777</b>	4536.3	0.83	254	<b>245.3</b>	1.04
11	<b>13817</b>	16447.9	0.84	971	<b>926.4</b>	1.05
12	<b>54701</b>	65492.3	0.84	3829	<b>3743.3</b>	1.02

#### 3.2 GraphBLAS v3.2.0, Feb 2020

Prob	one thread			40 threads		
	MATLAB	LAGraph	M/L	MATLAB	LAGraph	M/L
1		25.1	.		1.4	.
2		85.6	.		4.8	.
3		328.7	.		18.3	.
4		102.0	.		6.2	.
5		356.6	.		21.7	.
6		1395.9	.		84.1	.
7		722.0	.		35.0	.
8		2653.2	.		131.5	.
9		10407.3	.		504.7	.
10		4001.3	.		229.4	.
11		15124.5	.		918.9	.
12		.	.		.	.

## 4 Rate results

The rate is equal to the number of edges in the DNN, times the number of features (60,000 for all cases), divided by the run time. Rate is reported in terms of billions of edges/sec. Best rate shown in bold; higher is better.

### 4.1 GraphBLAS v3.1.0, August 2019

Prob	one thread			40 threads		
	MATLAB	LAGraph	M/L	MATLAB	LAGraph	M/L
1	<b>10.0</b>	9.7	1.03	85.8	<b>100.4</b>	0.85
2	<b>14.0</b>	13.8	1.01	101.3	<b>209.2</b>	0.48
3	<b>15.6</b>	15.5	1.00	110.1	<b>230.2</b>	0.48
4	<b>9.7</b>	8.7	1.11	94.6	<b>101.1</b>	0.93
5	<b>12.9</b>	11.4	1.13	<b>123.3</b>	122.9	1.00
6	<b>14.0</b>	12.4	1.14	<b>129.1</b>	128.4	1.01
7	4.9	<b>5.1</b>	0.97	64.6	<b>74.0</b>	0.87
8	5.6	<b>5.9</b>	0.95	75.2	<b>86.3</b>	0.87
9	5.8	<b>6.2</b>	0.94	77.2	<b>87.5</b>	0.88
10	<b>4.0</b>	3.3	1.20	59.4	<b>61.5</b>	0.96
11	<b>4.4</b>	3.7	1.19	62.2	<b>65.2</b>	0.95
12	<b>4.4</b>	3.7	1.20	63.1	<b>64.5</b>	0.98

### 4.2 GraphBLAS v3.2.0, Feb 2020

Prob	one thread			40 threads		
	MATLAB	LAGraph	M/L	MATLAB	LAGraph	M/L
1	.	9.4	.	.	164.9	.
2	.	11.0	.	.	198.5	.
3	.	11.5	.	.	205.8	.
4	.	9.2	.	.	152.8	.
5	.	10.6	.	.	174.1	.
6	.	10.8	.	.	179.6	.
7	.	5.2	.	.	107.9	.
8	.	5.7	.	.	114.8	.
9	.	5.8	.	.	119.7	.
10	.	3.8	.	.	65.8	.
11	.	4.0	.	.	65.7	.
12	.	.	.	.	.	.

## 5 Comparison

### 5.1 v3.1.0, August 2019

When using 40 threads, the performance of the two methods (MATLAB+LAGraph, vs pure C in LAGraph) is almost identical, except for problems 2 and 3, where LAGraph is about twice as fast as the MATLAB `dnn_gb.m`. The two codes differ in how the max threshold of 32 is implemented. The MATLAB interface doesn't allow for user-defined operators, so a mask `M` is used for the `dnn_gb` function. This actually seems to be faster in many cases when using a single thread, as compared to the method used in `LAGraph_dnn`. The latter uses a user-defined operator, `ymax32` and `GrB_apply`. With 40 threads, the `GrB_apply` is faster.

It appears that very little is lost, if any, in the MATLAB interface. To be certain of this, the `LAGraph_dnn.c` function would need to be modified to use the masked assignment method used in `dnn_gb.m`. However, each function was written using the most natural approach available, and since the MATLAB interface does not allow for user-defined operators, it was most natural to write the `max32` threshold as masked assignment. In that sense, this is a fair comparison between MATLAB+GraphBLAS and LAGRAPH+GraphBLAS.

### 5.2 v3.2.0, Feb 2020

The new method in v3.2.0 is more parallelizable than the method in v3.1.0. However, for the smaller problems, v3.1.0 is faster when using a single thread. When using 40 threads, v3.2.0 (in LAGraph) is almost always significantly faster than the method in v3.1.0 (except for problems 2 and 3).

The sequential performance in v3.1.0 is likely higher because that version of GraphBLAS kept a set of global workspaces that it reused for each matrix multiplication (the Sauna). Those have been removed in v3.2.0. For the next release of GraphBLAS, I will explore how to improve the sequential performance of the Sauna-free method in v3.2.x, to see if I can match the performance of the Sauna-based method in v3.1.0.

Using a single thread on the largest problem takes about 18 hours, and thus results were still in progress at the time v3.2.0 was released. This document will be updated after the formal release of v3.2.0 when the runs complete. See the master branch at <https://github.com/DrTimothyAldenDavis/GraphBLAS> for the updated document.