



# GPU SEMIRING PRIMITIVES FOR SPARSE NEIGHBORHOOD METHODS

Corey J. Nolet, Divye Gala, Edward Raff, Joe Eaton, Brad Rees, John Zedlewski, Tim Oates

AUGUST 2022

## Publicly Available Libraries for Computing Sparse Distances are Inflexible & Inefficient

---

- Computing distances between vectors is core to many machine learning algorithms
  - Computing them efficiently for sparse datasets can be very hard
  - Little easier when it can be computed with standard matrix multiplication (i.e., Euclidean)
- There are many distance metrics people want to use and a large variety of sparsity patterns and interactions to contend with

# Publicly Available Libraries for Computing Sparse Distances are Inflexible & Inefficient

- We present a single unified framework for computing several important sparse pairwise distances
  - Fast and memory efficient across many different sparsity patterns.
  - Provides reusable building blocks for composing many different important metrics in ML.
  - Can be extended to different execution patterns by optimizing specific sparsity patterns.
  - Already available to you in RAPIDS!

<https://github.com/rapidsai/raft>

## RAPIDS RAFT: Reusable Accelerated Functions and Tools

RAFT contains fundamental widely-used algorithms and primitives for data science and machine learning. The algorithms are CUDA-accelerated and form building-blocks for rapidly composing analytics.

By taking a primitives-based approach to algorithm development, RAFT

- accelerates algorithm construction time
- reduces the maintenance burden by maximizing reuse across projects, and
- centralizes core reusable computations, allowing future optimizations to benefit all algorithms that use them.

While not exhaustive, the following general categories help summarize the accelerated functions in RAFT:

Category	Examples
Data Formats	sparse & dense, conversions, data generation
Dense Linear Algebra	matrix arithmetic, norms, factorization, least squares, svd & eigenvalue problems
Spatial	pairwise distances, nearest neighbors, neighborhood graph construction
Sparse Operations	linear algebra, eigenvalue problems, slicing, symmetrization, labeling
Basic Clustering	spectral clustering, hierarchical clustering, k-means
Solvers	combinatorial optimization, iterative solvers
Statistics	sampling, moments and summary statistics, metrics
Distributed Tools	multi-node multi-gpu infrastructure

RAFT provides a header-only C++ library and pre-compiled shared libraries that can 1) speed up compile times and 2) enable the APIs to be used without CUDA-enabled compilers.

RAFT also provides 2 Python libraries:

- `pylibraft` - low-level Python wrappers around RAFT algorithms and primitives.
- `pyraft` - reusable infrastructure for building analytics, including tools for building both single-GPU and multi-node multi-GPU algorithms.

# Semirings and Relation To Matrix Multiplication

---

- A **monoid** contains an associative binary relation, such as addition ( $\oplus$ ), and an identity element ( $id_{\oplus}$ )
- A **semiring**, denoted  $(S, R, \{\oplus, id_{\oplus}\}, \{\otimes, id_{\otimes}\})$ , is a tuple containing additive ( $\oplus$ ) and multiplicative ( $\otimes$ ) monoids where
  1.  $\oplus$  is commutative, distributive, and has an identity element 0
  2.  $\otimes$  distributes over  $\oplus$
- Given two sparse vectors  $a, b \in R^k$ , a semiring with  $(S, R, \{\oplus, 0\}, \{\otimes, 1\})$  and  $annihilator_{\otimes} = 0$  is a **standard matrix multiplication**.
- Sparse Matrix-Vector multiplication (**SPMV**) is fundamental low-level BLAS routine in sparse matrix multiplication. Our contribution is a CUDA-accelerated Sparse Matrix-Sparse Vector (**SPSV**) multiplication primitive.

# The Euclidean Semiring

---

- Let vector  $a = [1,0,1]$  and  $b = [0,1,0]$
- Take the formula for computing Euclidean distance

$$\sqrt{\sum_{i=1}^n |x_i - y_i|^2}$$

- We can use the distributive property to compute more efficiently in parallel:
  - $x^2 - 2\langle x, y \rangle + y^2$
  - Can compute with a simple dot product (plus-times semiring) and L2 norms of  $x$  and  $y$ .

$$\sum(|a - b|) = \quad (4)$$

$$\sum([|1 - 0|, |0 - 1|, |1 - 0|]) = \quad (5)$$

$$\sum([1, 1, 1]) = 3 \quad (6)$$

$$\sum(|a - b|) = \quad (7)$$

$$\sum([|1 - 0|, |0 - 1|, |1 - 0|]) = \quad (8)$$

$$\sum([0, 0, 0]) = 0 \quad (9)$$

# The Manhattan Semiring and Non-Annihilating Multiplicative Monoid (NAMM)

---

- Let vector  $a = [1,0,1]$  and  $b = [0,1,0]$
- We take the sum of the absolute value of their differences (eqs 4, 5, 6)
- Semiring libraries rely on the detail that the multiplicative annihilator is equal to the additive identity.
  - If we follow this detail in our example, we end up with the following result of Eqs. 7, 8, 9 (if any side is 0, the arithmetic evaluates to 0).
- What we need here instead is for the multiplicative identity to be **non-annihilating**, evaluating to the other side when either side is zero and evaluating to 0 only in the case where both sides have the same value. i.e., :

$$\begin{array}{l} |1 - 0| = 1 \\ |0 - 1| = 1 \\ |0 - 0| = 0 \\ |1 - 1| = 0 \end{array}$$

$$\sum(|a - b|) = \quad (4)$$

$$\sum([|1 - 0|, |0 - 1|, |1 - 0|]) = \quad (5)$$

$$\sum([1, 1, 1]) = 3 \quad (6)$$

$$\sum(|a - b|) = \quad (7)$$

$$\sum([|1 - 0|, |0 - 1|, |1 - 0|]) = \quad (8)$$

$$\sum([0, 0, 0]) = 0 \quad (9)$$

# Semirings of Several Important Distances

Distance	Formula	NAMM	Norm	Expansion
Correlation	$1 - \frac{\sum_{i=0}^k (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=0}^k x_i - \bar{x}^2} \sqrt{\sum_{i=0}^k y_i - \bar{y}^2}}$		$L_1, L_2$	$1 - \frac{k\langle x \cdot y \rangle - \ x\  \ y\ }{\sqrt{(k\ x\ _2 - \ x\ ^2)(k\ y\ _2 - \ y\ ^2)}}$
Cosine	$\frac{\sum_{i=0}^k x_i y_i}{\sqrt{\sum_{i=0}^k x_i^2} \sqrt{\sum_{i=0}^k y_i^2}}$		$L_2$	$1 - \frac{\langle x \cdot y \rangle}{\ x\ _2 \ y\ _2}$
Dice-Sorensen	$\frac{2 \sum_{i=0}^k x_i y_i }{(\sum_{i=0}^k x_i)^2 + (\sum_{i=0}^k y_i)^2}$		$L_0$	$\frac{2\langle x \cdot y \rangle}{ x ^2 +  y ^2}$
Dot Product	$\sum_{i=0}^k x_i y_i$			$\langle x \cdot y \rangle$
Euclidean	$\sqrt{\sum_{i=0}^k  x_i - y_i ^2}$		$L_2$	$\ x\ _2^2 - 2\langle x \cdot y \rangle + \ y\ _2^2$
Canberra	$\sum_{i=0}^k \frac{ x_i - y_i }{ x_i  +  y_i }$	$\{\frac{ x-y }{ x + y }, 0\}$		
Chebyshev	$\sum_{i=0}^k \max(x_i - y_i)$	$\{\max(x - y), 0\}$		
Hamming	$\frac{\sum_{i=0}^k x_i \neq y_i}{k}$	$\{x \neq y, 0\}$		
Hellinger	$\frac{1}{\sqrt{2}} \sqrt{\sum_{i=0}^k (\sqrt{x_i} - \sqrt{y_i})^2}$			$1 - \sqrt{\langle \sqrt{x} \cdot \sqrt{y} \rangle}$
Jaccard	$\frac{\sum_{i=0}^k x_i y_i}{(\sum_{i=0}^k x_i^2 + \sum_{i=0}^k y_i^2 - \sum_{i=0}^k x_i y_i)}$		$L_0$	$1 - \frac{\langle x \cdot y \rangle}{(\ x\  + \ y\  - \langle x \cdot y \rangle)}$
Jensen-Shannon	$\sqrt{\frac{\sum_{i=0}^k x_i \log \frac{x_i}{\mu_i} + y_i \log \frac{y_i}{\mu_i}}{2}}$	$\{x \log \frac{x}{\mu} + y \log \frac{y}{\mu}, 0\}$		
KL-Divergence	$\sum_{i=0}^k x_i \log(\frac{x_i}{y_i})$			$\langle x \cdot \log \frac{x}{y} \rangle$
Manhattan	$\sum_{i=0}^k  x_i - y_i $	$\{ x - y , 0\}$		
Minkowski	$(\sum_{i=0}^k  x_i - y_i ^p)^{1/p}$	$\{ x - y ^p, 0\}$		
Russel-Rao	$\frac{k - \sum_{i=0}^k x_i y_i}{k}$			$\frac{k - \langle x \cdot y \rangle}{k}$

# SPSV CUDA Kernel: Load-Balanced Hybrid CSR+COO

---

1. **Load-balancing** using a row index array in coordinate format (COO) for B, coalescing the loads from each vector from A
2. Lowered memory footprint by **removing the need to transpose B.**
3. **Two-pass execution**

---

**Algorithm 3** Load-balanced Hybrid CSR+COO SPMV.

---

**Input:**  $A_i, B, product\_op, reduce\_op$

**Result:**  $C_{ij} = d(A_i, B_j)$

read  $A_i$  into shared memory

cur\_row=rowidx[ind]

ind = idx of first elem to be processed by this thread

c = product\_op(A[ind], x[colidx[ind]])

**for**  $i \leftarrow 1$  **to**  $nz\_per\_chunk$ ; **by**  $warp\_size$  **do**

    next\_row = cur\_row +  $warp\_size$

**if**  $next\_row \neq cur\_row$  —  $is\_final\_iter?$  **then**

        v = segmented\_scan(cur\_row, c, product\_op)

**if**  $is\_segment\_leader?$  **then**

            atomic\_reduce(v, reduce\_op)

**end**

        c = 0

**end**

    cur\_row = next\_row

    ind +=  $warp\_size$

    c = product\_op(A[ind], x[colidx[ind]])

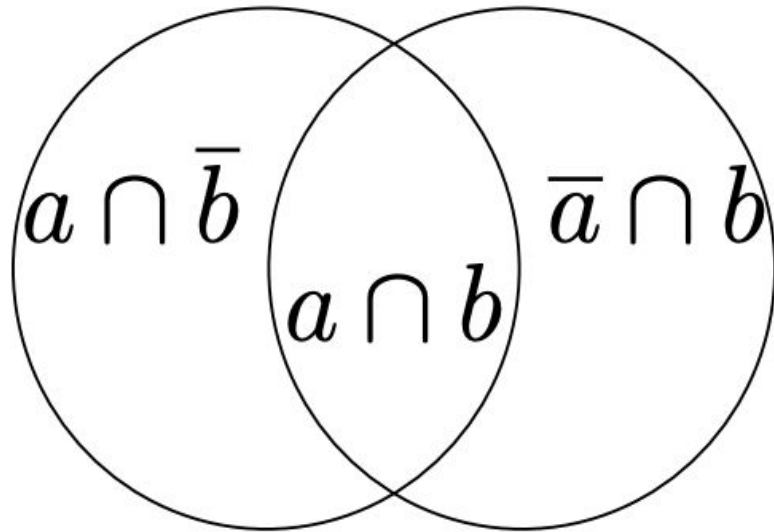
**end**

---



## Implement NAMM with Multiple Passes of an SPMV

---



$$a \cup b = \{a \cap b\} \cup \{\bar{a} \cap b\} \cup \{a \cap \bar{b}\}$$

- A single pass computes the intersection  $a \cap b$  between nonzero columns from each vector  $a$ , and  $b$  so long as  $\otimes$  is applied to all nonzero columns of  $b$
- A second pass can compute the remaining symmetric difference required for the full union between non-zero column
- $id_{\otimes}$  in  $B$  is skipped in the second pass.

## Performance- It's Fast But Also Memory Efficient

---

- Benchmarks were performed on a DGX1 containing dual 20-core Intel Xeon ES-2698 CPUs (80 total threads) at 2.20GHZ and a Volta V100 GPU running CUDA 11.0 for both the driver and toolkit.
- Each benchmark performs a k-nearest neighbors query to test our primitives end-to-end and allow scaling to datasets where the dense pairwise distance matrix may not otherwise fit in the memory of the GPU
- We used the brute-force *NearestNeighbors* estimator from RAPIDS cuML for the GPU benchmarks since it makes direct use of our primitive
- We used Scikit-learn's corresponding brute-force *NearestNeighbors* estimator as a CPU baseline and configured it to use all the available CPU cores

## Performance- Fast And Memory Efficient

Table 3: Benchmark Results for all datasets under consideration. All times are in seconds, best result in **bold**. The first italicized set of distances can all be computed as dot products, which are already highly optimized for sparse comparisons today. This easier case we are still competitive, and sometimes faster, than the dot-product based metrics. The Non-trivial set of distances that are not well supported by existing software are below, and our approach dominates amongst all these metrics.

Distance		MovieLens		scRNA		NY Times Bag of Words		SEC Edgar	
		Baseline	RAFT	Baseline	RAFT	Baseline	RAFT	Baseline	RAFT
Dot Product Based	<i>Correlation</i>	130.57	<b>111.20</b>	<b>207.00</b>	235.00	<b>257.36</b>	337.11	134.79	<b>87.99</b>
	<i>Cosine</i>	131.39	<b>110.01</b>	<b>206.00</b>	233.00	<b>257.73</b>	334.86	127.63	<b>87.96</b>
	<i>Dice</i>	130.52	<b>110.94</b>	<b>206.00</b>	233.00	<b>130.35</b>	335.49	134.36	<b>88.19</b>
	<i>Euclidean</i>	131.93	<b>111.38</b>	<b>206.00</b>	233.00	<b>258.38</b>	336.63	134.75	<b>87.77</b>
	<i>Hellinger</i>	129.79	<b>110.82</b>	<b>205.00</b>	232.00	<b>258.22</b>	334.80	134.11	<b>87.83</b>
	<i>Jaccard</i>	130.51	<b>110.67</b>	<b>206.00</b>	233.00	<b>258.24</b>	336.01	134.55	<b>87.73</b>
	<i>Russel-Rao</i>	130.35	<b>109.68</b>	<b>206.00</b>	232.00	<b>257.58</b>	332.93	134.31	<b>87.94</b>
Non-Trivial Metrics	Canberra	3014.34	<b>268.11</b>	4027.00	<b>598.00</b>	4164.98	<b>819.80</b>	505.71	<b>102.79</b>
	Chebyshev	1621.00	<b>336.05</b>	3907.00	<b>546.00</b>	2709.30	<b>1072.35</b>	253.00	<b>146.41</b>
	Hamming	1635.30	<b>229.59</b>	3902.00	<b>481.00</b>	2724.86	<b>728.05</b>	258.27	<b>97.65</b>
	Jensen-Shannon	7187.27	<b>415.12</b>	4257.00	<b>1052.00</b>	10869.32	<b>1331.37</b>	1248.83	<b>142.96</b>
	KL Divergence	5013.65	<b>170.06</b>	4117.00	<b>409.00</b>	7099.08	<b>525.32</b>	753.56	<b>87.72</b>
	Manhattan	1632.05	<b>227.98</b>	3904.00	<b>477.00</b>	2699.91	<b>715.78</b>	254.69	<b>98.05</b>
	Minkowski	1632.05	<b>367.17</b>	4051.00	<b>838.00</b>	5855.79	<b>1161.31</b>	646.71	<b>129.47</b>

## It Is Used In The RAPIDS cuML Library

---

- Enables several clustering and manifold learning algorithms to accept sparse inputs.
- Also being used in cuML's Sparse k-Nearest Neighbors estimator.
- Already available in current RAPIDS, no hard work required.

<https://github.com/rapidsai/cuml>

```
from cuml.neighbors import
    ↪ NearestNeighbors
nn = NearestNeighbors().fit(X)
dists, inds = nn.kneighbors(X)
from cuml.metrics import
    ↪ pairwise_distances
dists = pairwise_distances(X,
    ↪ metric='cosine')
```

Figure 2: Excluding data loading and logging, all the code needed to perform the same GPU accelerated sparse distance calculations done in this paper are contained within these two snippets. Top shows k-NN search, bottom all pairwise distance matrix construction. These are the APIs that most would use.

# RAFT Library Provides C++ API for Defining New Distance Semirings

---

Just define monoids!

```
#include
↪ <raft/sparse/distance/coo_spmv.cuh>
#include <raft/sparse/distance/operators.h>

using namespace raft::sparse::distance

distances_config_t<int, float> conf;

// Use conf to set input data arguments...

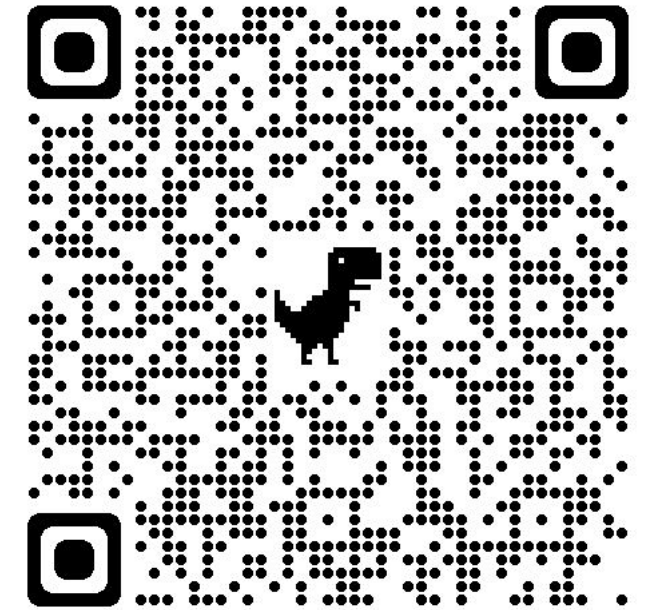
balanced_coo_pairwise_generalized_spmv(
    out_dists, conf, coo_rows_a,
    AbsDiff(), Sum(), AtomicSum());

balanced_coo_pairwise_generalized_spmv_rev(
    out_dists, conf, coo_rows_b,
    AbsDiff(), Sum(), AtomicSum());
```

## Conclusion / Questions?

---

- Semirings provide us a framework for unifying many important distances in ML applications.
- Our SPSV kernel is state of the art in performance, efficiency and flexibility



Check out the paper for details!

@cjnolet 

<https://github.com/rapidsai/raft>

<https://github.com/rapidsai/cuml>