

# gSSJoin: a GPU-based Set Similarity Join Algorithm

Sidney R. Junior, Rafael D. Quirino, Leonardo A. Ribeiro,  
Wellington S. Martins



INSTITUTO DE  
INFORMÁTICA  
UFG

# Agenda

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

- 1 Introduction
- 2 Background
- 3 The gSSJoin Algorithm
- 4 Experiments
- 5 Related Work
- 6 Conclusion and Future Work

# Introduction

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

Set similarity join returns all pairs of similar sets from a dataset. Sets are considered similar if the value returned by a set similarity function for them is not less than a given threshold.

# Applications

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

- Integration
- Cleaning
- Plagiarism identification
- Data Mining

# Current Algorithms

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

Most state-of-the-art algorithms uses a technique called *prefix filtering* to prune dissimilar sets by inspecting only a fraction of them.

# Limitations

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

Prefix filtering is only effective at high threshold values. When the threshold value is low, the fraction needed to be inspected is bigger and the remaining sets needed to be verified after the pruning step is much higher.

# Contributions

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

- A fine-grained parallel algorithm for both indexing the data and performing the set similarity joins.
- A highly threaded GPU implementation that takes advantage of intensive occupation, hierarchical memory, and coalesced memory access.
- A scalable multi-GPU implementation that exploits both data parallelism and task parallelism.
- Extensive experimental work with standard real-world textual datasets.

# Data Representation

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

We map the strings to sets of tokens called *q-grams*. *q-grams* are sub-strings of length  $q$  obtained by sliding a window over the characters of the input string.

3-grams of "gSSJoin": {"gSS", "SSJ", "SJo", "Joi", "oin"}



# Formal Definition

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

Given two set collections  $\mathcal{R}$  and  $\mathcal{S}$ , a set similarity function  $sim$ , and a threshold  $\tau$ , the *set similarity join* between  $\mathcal{R}$  and  $\mathcal{S}$  returns all scored set pairs  $\langle (r, s), \tau' \rangle$  s.t.  $(r, s) \in \mathcal{R} \times \mathcal{S}$  and  $sim(r, s) = \tau' \geq \tau$ .

# Similarity Function

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

A set similarity function  $sim(r, s)$  returns a value in  $[0, 1]$  to represent their similarity. A greater value indicates that  $r$  and  $s$  have higher similarity.

# Similarity Function

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

$$\text{Jaccard Similarity: } J(r, s) = \frac{|r \cap s|}{|r \cup s|}$$

# Overlap Constraint

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

A predicate involving the Jaccard similarity and a threshold  $\tau$  can be equivalently rewritten into a set overlap constraint:

$$J(r, s) \geq \tau \iff |r \cap s| \geq \frac{\tau}{1+\tau} (|r| + |s|).$$

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

For a threshold  $\tau$ , we can identify all candidate matches of a given set  $r$  using a prefix of length  $|r| - \lceil |r| \cdot \tau \rceil + 1$ .

---

## Algorithm 1: General set similarity join algorithm.

---

**Input:** A sorted set collection  $\mathcal{C}$ , a threshold  $\tau$

**Output:** A set  $S$  containing all pairs  $(r, s)$  s.t.  $\text{sim}(r, s) \geq \tau$

```
1  $I_1, I_2, \dots, I_{|\mathcal{U}|} \leftarrow \emptyset, S \leftarrow \emptyset$ 
2 foreach  $r \in \mathcal{C}$  do
3   foreach  $t \in \text{pref}_\beta(r)$  do
4     foreach  $s \in I_t$  do
5       if not  $\text{filter}(r, s)$ 
6          $S \leftarrow S \cup \text{refine}(r, s)$ 
7        $I_t \leftarrow I_t \cup \{r\}$ 
8 return  $S$ 
```

---

# Limitations of Current Algorithms

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

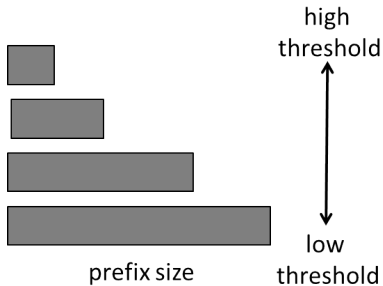
**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work



**threshold effect**

# Limitations of Current Algorithms

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

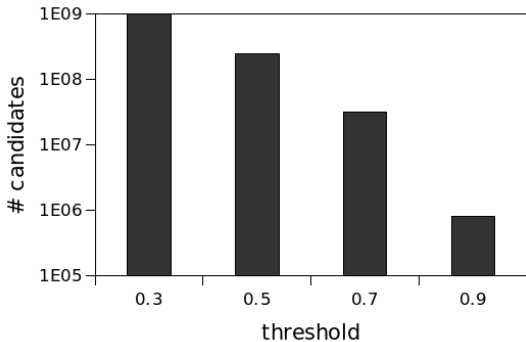
**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work





# Limitations of Current Algorithms

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

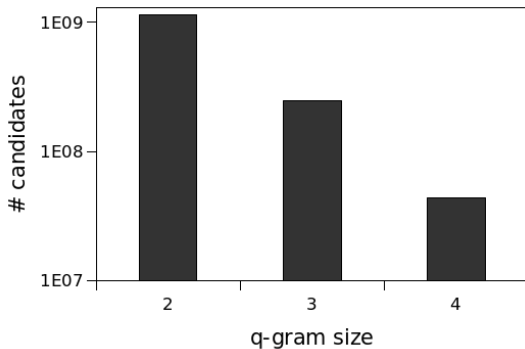
**Background**

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work



# The gSSJoin Algorithm

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

**The gSSJoin  
Algorithm**

Experiments

Related  
Work

Conclusion  
and Future  
Work

- Exploits massive parallelism instead of filtering
- Builds an inverted index on the pre-processing phase
- Performs a set similarity search for each set

# Pre-processing

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

---

## Algorithm 2: *DataIndexing*( $E$ )

---

**input** : token-set pairs in  $E[0..|\mathcal{E}|-1]$ .  
**output**: *count*, *index*, *cardinality*, *invertedIndex*.

- 1 array of integers *count*[ $0..|\mathcal{V}|-1$ ] // count array, initialized with zeros.
- 2 array of integers *index*[ $0..|\mathcal{V}|-1$ ].
- 3 array of integers *cardinality*[ $0..|\mathcal{S}|-1$ ].
- 4 *invertedIndex*[ $0..|\mathcal{E}|-1$ ] // the inverted index
- 5 Count the occurrences of each token in parallel on the input and accumulates in *count*.
- 6 Perform an exclusive parallel prefix sum on *count* and stores the result in *index*.
- 7 Access in parallel the pairs in  $E$ , with each processor performing the following tasks:
  - 8     **begin**
  - 9     |     Contribute on the cardinality computation of each set in *cardinality*[*s*].
  - 10    |     Store in *invertedIndex* the entries corresponding to pairs in  $E$ , according to *index*.
  - 11    |     **end**
- 12 **Return** the arrays: *count*, *index*, *cardinality* and *invertedIndex*.

---

# Set Similarity Search

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

---

## Algorithm 3: *SimilaritySearch(invertedIndex, s)*

---

```
input : invertedIndex, count, index, cardinality, threshold, inputset  $s[0..|V_s| - 1]$ .  
output: Jaccard similarity array jac_sim $[0..|S| - 1]$  initialized with zeros.  
1 array of integers counts $[0..|V_s| - 1]$  initialized with zeros  
2 array of integers indexs $[0..|V_s| - 1]$   
3 for each token  $t_i \in s$ , in parallel do  
4   | counts $[i] = \text{count}[t_i]$ ;  
5   | indexs $[i] = \text{index}[t_i]$ ;  
6 end  
7 Perform an inclusive parallel prefix sum on counts and stores the results in indexs  
8 foreach processor  $p_i \in \mathcal{P}$  do  
9   | for  $x \in [i \lceil \frac{|E_s|}{|\mathcal{P}|} \rceil, \min((i+1) \lceil \frac{|E_s|}{|\mathcal{P}|} \rceil - 1, |E_s| - 1)]$  do  
10  |   | // Map  $x$  to the correct position indInvPos of the invertedIndex  
11  |   |  $pos = \min(i : index_s[i] > x)$ ;  
12  |   | if  $pos = 0$  then  
13  |   |   |  $p = 0$ ; offset =  $x$ ;  
14  |   | else  
15  |   |   |  $p = index_s[pos - 1]$ ; offset =  $x - p$ ;  
16  |   | end  
17  |   | indInvPos = starts $[pos] + offset$   
18  |   | uses s $[pos]$  and invertedIndex $[indInvPos]$  in the partial computation of the intersection  
19  |   | between s and the set associated to invertedIndex $[indInvPos]$   
20  |   | end  
21  |   | for  $x \in [i \lceil \frac{|S|}{|\mathcal{P}|} \rceil, \min((i+1) \lceil \frac{|S|}{|\mathcal{P}|} \rceil - 1, |S| - 1)]$  do  
22  |   |   | // Jaccard similarity calculation for  $\mathcal{S}$  sets using  $\mathcal{P}$  processors  
23  |   |   | uses intersection and union (through cardinality) to compute the Jaccard similarity  
24  |   |   | flag sets with Jaccard similarity above the threshold  
25  |   | end  
26 end  
27 Perform an exclusive parallel prefix sum on the flagged sets to compact the selected sets  
28 Return the array: jac_sim with the selected jaccard similarities.
```

- The inverted index is created in all GPUs
- The search sets are divided between the GPUs

---

## Algorithm 4: *MultiGPU Search*( $E$ )

---

**input** : token-set pairs in  $E[0..|\mathcal{E}|-1]$ .

**output**: A list of the most similar, one for each set.

```
1 for each  $i \in g$ , in parallel do
2   |    $set.gpu.device(i)$ ;
3   |    $DataIndexing(E)$ ;
4 end
5 Allocate memory for the biggest set;
6 for each  $j \in g$ , in parallel do
7   |    $set.gpu.device(j)$ ;
8   |   for each set  $s \in (m/g)$  in parallel do
9     |    $SimilaritySearch(invertedIndex, s)$ ;
10  |   end
11 end
12 Return A list of the most similar sets, one for each set.
```

---

# Experimental Setup

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

**Experiments**

Related  
Work

Conclusion  
and Future  
Work

We implemented gSSJoin using the CUDA Toolkit version 7.5. The experiments were conducted on a machine running CentOS 7.2.1511 64-bits, with 24 Intel Xeon E5-2620, 16GB of ECC RAM, and four GeForce Zotac Nvidia GTX Titan Black, with 6GB of RAM and 2,880 CUDA cores each.

Table : Dataset statistics.

Database	Number of sets	Max	Mean	Standard deviation
DBLP	100k	205	70	23.9
Netflix	200k	54	30.54	8.2



# Performance Results

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

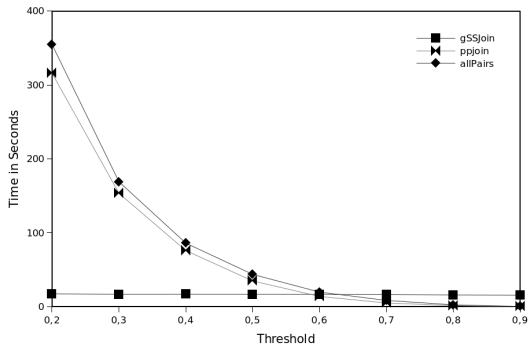


Figure : DBLP dataset using 3-grams.

# Performance Results

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

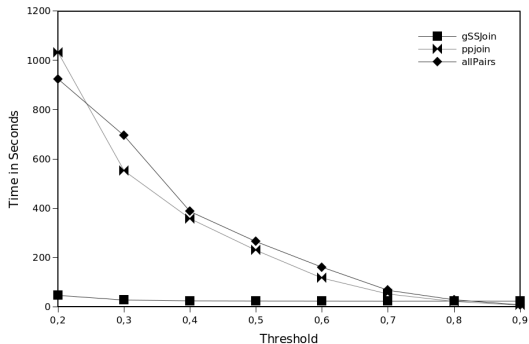


Figure : Netflix dataset using 3-grams.

# Performance Results

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

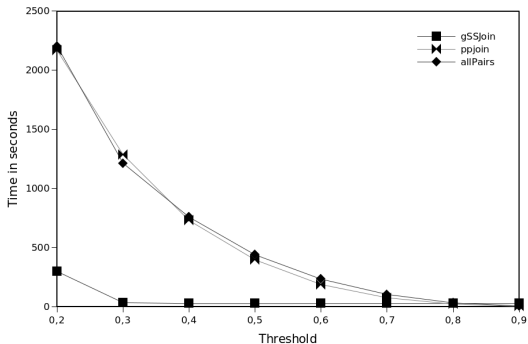


Figure : DBLP dataset using 2-grams.

# Performance Results

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

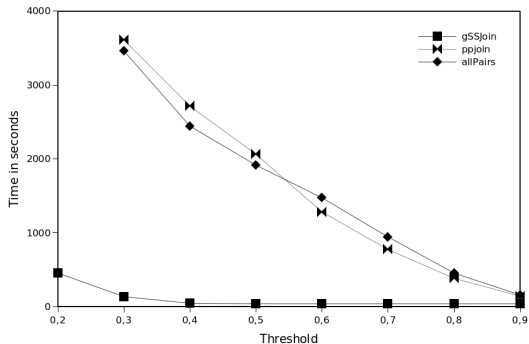


Figure : Netflix dataset using 2-grams.

# Performance Results

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

Table : Execution Time (secs) Using Multiple GPUs, Threshold = 0.5

Number of GPUs:	1	2	3	4
DBLP 2gram	27.56	14.37	10.4	7.2
DBLP 3gram	16.73	8.77	6	4.47
Netflix 2gram	41.10	21.63	14.37	11.2
Netflix 3gram	23.52	12.69	8.5	6.5

# Related Work

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

- The state-of-the-art algorithms are intrinsically sequential [Sarawagi and Kirpal, Chaudhuri et al. 2006, Bayardo et al. 2007, Xiao et al. 2011, Ribeiro and Härder 2011, Wang et al. 2012].
- A distributed memory MapReduce framework to perform similarity joins was proposed by [Vernica et al. 2010].
- Parallel solutions such as [Cruz et al. 2015] perform approximate set similarity joins using MinHash and Locality Sensitive Hashing.

# Conclusion

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

**Conclusion  
and Future  
Work**

The gSSJoin algorithm has a much better performance than the state-of-the-art CPU-based algorithms when the threshold is low and the database is more uniform.

In future work, we plan to incorporate candidate filtering techniques to obtain even greater speed-ups.

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

# Questions?

email: [srjsoftware@gmail.com](mailto:srjsoftware@gmail.com)



# References

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work



[Cruz et al. 2015]

Cruz, M. S. H., Kozawa, Y., Amagasa, T., and Kitagawa, H. (2015). GPU Acceleration of Set Similarity Joins. In DEXA, pages 384–398.



[Bayardo et al. 2007]

Bayardo, R. J., Ma, Y., and Srikant, R. (2007). Scaling up All Pairs Similarity Search. In WWW, pages 131–140.



[Chaudhuri et al. 2006]

Chaudhuri, S., Ganti, V., and Kaushik, R. (2006). A primitive operator for similarity joins in data cleaning. In ICDE, page 5.



[Ribeiro and Härder 2011]

Ribeiro, L. A. and Härder, T. (2011). Generalizing Prefix Filtering to Improve Set Similarity Joins. Information Systems, 36(1):62–78.



[Sarawagi and Kirpal 2004]

Sarawagi, S. and Kirpal, A. Efficient Set Joins on Similarity Predicates. In SIGMOD.

# References

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work



[[Vernica et al. 2010](#)]

Vernica, R., Carey, M. J., and Li, C. (2010). Efficient Parallel Set Similarity Joins using MapReduce. In SIGMOD, pages 495–506.



[[Xiao et al. 2011](#)]

Xiao, C., Wang, W., Lin, X., Yu, J. X., and Wang, G. (2011). Efficient Similarity Joins for Near-duplicate Detection. TODS, 36(3):15.



[[Wang et al. 2012](#)]

Wang, J., Li, G., and Feng, J. (2012). Can We Beat the Prefix Filtering?: An Adaptive Framework for Similarity Join and Search. In SIGMOD, pages 85–96.

# GPU Architecture

Mestrado em  
Ciência da  
Computação

Sidney R.  
Junior,  
Rafael D.  
Quirino,  
Leonardo A.  
Ribeiro,  
Wellington  
S. Martins

Introduction

Background

The gSSJoin  
Algorithm

Experiments

Related  
Work

Conclusion  
and Future  
Work

