Similarity Search Set Similarity Join

Nikolaus Augsten

nikolaus.augsten@plus.ac.at Department of Computer Science University of Salzburg



WS 2025/26

Version November 6, 2025

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Filters for the Set Similarity Join Motivation

Outline

- Filters for the Set Similarity Join
 - Motivation
 - Signature-based Filtering
 - Signatures for Overlap Similarity
 - Signatures for Hamming Distance
- 2 Implementations of Set Similarity Joins
 - Other Similarity Functions
 - Table of Set Similarity Join Algorithms and their Signatures
- Conclusion

Outline

- Filters for the Set Similarity Join
 - Motivation
 - Signature-based Filtering
 - Signatures for Overlap Similarity
 - Signatures for Hamming Distance
- 2 Implementations of Set Similarity Joins
 - Other Similarity Functions
 - Table of Set Similarity Join Algorithms and their Signatures
- 3 Conclusion

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

4 / 48

Application Scenario

• Scenario: A social network company stores user interests.

Filters for the Set Similarity Join Motivation

• Example: user table with interests:

R

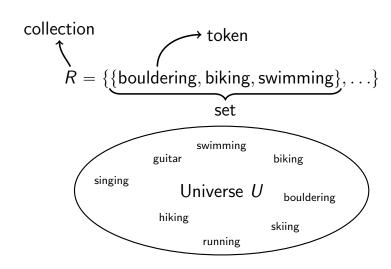
id	name	interests
S	Sebastian	{bouldering, biking, swimming}
n	Nathan	{bouldering, swimming, guitar, singing}
р	Philippides	{hiking, running}
m	Maria	{bouldering, hiking, running}
r	Rosa	{bouldering, skiing, hiking}

Task: Recommend new friends based on similar interests!

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 Augsten (Univ. Salzburg) 3 / 48

Similarity Search WS 2025/26





Filters for the Set Similarity Join Motivation

Filters for the Set Similarity Join Motivation

WS 2025/26

The Join Approach

• Solution: Compute the set similarity join

Definition (Set Similarity Join)

Given two collections of sets R and S, the set similarity join computes

$$R \approx S = \{(r,s) \in R \times S \mid Sim(r,s) \geqslant t\}$$

for a similarity function Sim or

$$R \approx S = \{(r, s) \in R \times S \mid \mathsf{Dis}(r, s) \leqslant t\}$$

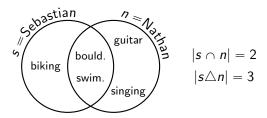
for a distance function Dis and threshold t.

- Naive Approach:
 - 1. Compute all pairs $R \times S$
 - 2. Test if $Sim(r, s) \ge t$ or $Dis(r, s) \le t$ on each tuple

Filters for the Set Similarity Join Motivation

Measuring Similarity of Sets

- Goal: measure the similarity of two sets r, s
- Similarity Function:
 - Sim(r, s) is high for similar sets, low for dissimilar sets
 - Example: Overlap $|r \cap s|$
- Distance Function:
 - Dis(r, s) is *low* for similar sets, *high* for dissimilar sets
 - Example: Hamming distance $|r \triangle s| = |(r \backslash s) \cup (s \backslash r)| = |r \cup s| |r \cap s|$
- Example:



WS 2025/26

Naive Join Example

• Example: self-join $R \bowtie R$, overlap similarity, threshold t = 2

Filters for the Set Similarity Join Motivation

K						
id	name	interests				
S	Sebastian	{bouldering, biking, swimming}				
n	Nathan	{bouldering, swimming, guitar, singing}				
р	Philippides	{hiking, running}				
m	Maria	{bouldering, hiking, running}				
r	Rosa	{bouldering, skiing, hiking}				

$$R \cong R = \{(s,n),(p,m),(m,r)\}$$

• 10 (non-reflexive, non-symmetric) comparisons!

Augsten (Univ. Salzburg) WS 2025/26 Similarity Search 7 / 48 Augsten (Univ. Salzburg) Similarity Search Filters for the Set Similarity Join Motivation

Demonstration

- Experiment: Naive approach
 - self-join with varying |R|
 - average set size 10
 - universe size 1000, uniformly distributed
 - overlap similarity with threshold t = 4

R	#comparisons	runtime [s]
1000	$5 \cdot 10^{5}$	0.022
10000	$5 \cdot 10^7$	2.288
100000	$5\cdot 10^9$	218.773

- A *single* similarity computation is fast (\approx 150 CPU cycles)
- But the search space grows fast: $\Theta(|R|^2)$

Augsten (Univ. Salzburg

Filters for the Set Similarity Join Signature-based Filtering

WS 2025/26

Augsten (Univ. Salzburg)

Conclusion

Outline

Motivation

Filters for the Set Similarity Join

Signature-based Filtering

Signatures for Overlap Similarity

Signatures for Hamming Distance

2 Implementations of Set Similarity Joins Other Similarity Functions

Filters for the Set Similarity Join Signature-based Filtering

• Table of Set Similarity Join Algorithms and their Signatures

Filters for the Set Similarity Join Signature-based Filtering

WS 2025/26

Reducing the Search Space using Filters

- Filtering: Reduce the search space by removing dissimilar pairs of sets
- Set similarity Join: Most filters are signature-based

Definition (Signature Scheme)

A signature scheme Sign is a function that maps a set of tokens to a set of signatures such that for any two sets of tokens, r and s:

$$Sim(r, s) \ge t \Rightarrow Sign(r) \cap Sign(s) \ne \emptyset$$

for a similarity function Sim and

$$\mathsf{Dis}(r,s) \leqslant t \Rightarrow \mathsf{Sign}(r) \cap \mathsf{Sign}(s) \neq \emptyset$$

for a distance function Dis.

• Intuition: Similar sets share at least one signature.

Signature-based Set Similarity Join

- Idea: Similar sets share signatures.
 - 1. Find all pairs sharing signatures (candidates)
 - 2. Test if $Sim(r, s) \ge t$ or $Dis(r, s) \le t$ on each tuple
- How do we find pairs sharing signatures?
 - 1. Compute all pairs $R \times S$
 - 2. Test if $Sign(r) \cap Sign(s) \neq \emptyset$ on each tuple
- Likely slower than naive approach!
- Index: Build a simple index to find sets for each signature

Augsten (Univ. Salzburg) WS 2025/26 11 / 48 WS 2025/26 12 / 48 Similarity Search Augsten (Univ. Salzburg) Similarity Search

Inverted-list Index

- Inverted-list Index: Stores mappings from content (e.g., signatures) to locations (e.g., sets)
 - 1. Compute signatures Sign(s) for set s
 - 2. Store a pointer to s in the list I_{sig} of each signature $sig \in Sign(s)$
- Example:

$$R = \{s_1, s_2, s_3, s_4\}$$

$$Sign(s_1) = \{\star, \spadesuit\}$$

$$Sign(s_2) = \{\circ, \clubsuit, \heartsuit\}$$

$$Sign(s_3) = \{\circ, \spadesuit, \diamondsuit\}$$

$$Sign(s_4) = \{\clubsuit, \spadesuit, \heartsuit, \diamondsuit\}$$

$$Sign(s_4) = \{\clubsuit, \spadesuit, \heartsuit, \diamondsuit\}$$

• A good signature scheme is both easy to compute and results in few false positives (= number of unnecessary verifications).

Augsten (Univ. Salzburg)

Similarity Search Filters for the Set Similarity Join Signatures for Overlap Similarity

WS 2025/26

Augsten (Univ. Salzburg)

Filters for the Set Similarity Join Signatures for Overlap Similarity

WS 2025/26

Outline

- Filters for the Set Similarity Join
 - Motivation
 - Signature-based Filtering
 - Signatures for Overlap Similarity
 - Signatures for Hamming Distance
- 2 Implementations of Set Similarity Joins
 - Other Similarity Functions
 - Table of Set Similarity Join Algorithms and their Signatures
- Conclusion

Filters for the Set Similarity Join Signature-based Filtering

Signature-based Framework

Algorithm 1: Signature-based Framework

Data: Collection *R*, threshold *t* **Result:** All similar pairs $M \subseteq R \times S$

 $I_{sig} \leftarrow I_{sig} \cup \{s\}$

forall $s \in S$ do

 $M \leftarrow \emptyset, C \leftarrow \emptyset$

forall $r \in R$ do

 $I \leftarrow \emptyset$ // inverted list index

forall signatures $sig \in Sign(s)$ **do**

forall signatures $sig \in Sign(r)$ **do**

 $C \leftarrow C \cup \{(r,s) \mid s \in I_{sig}\}$ **forall** candidate pairs $(r, s) \in C$ **do**

 $M \leftarrow M \cup (r, s)$ if $Sim(r, s) \ge t$

// indexing

// probing

(or $Dis(r, s) \leq t$)

Identity Signature

return M

• Simplest signature scheme (for overlap) is identity (Sign = Id):

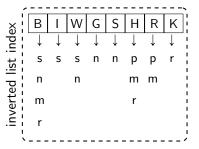
$$|r \cap s| \ge t \Rightarrow \operatorname{Id}(r) \cap \operatorname{Id}(s) \ne \emptyset$$

 $\Leftrightarrow |r \cap s| \ge t \Rightarrow |r \cap s| \ge 1$ assuming $t \ge 1$

- Every token is a signature
- Example:

Augsten (Univ. Salzburg)

id	interests
S	{ B ould., b l king, s W im.}
n	$\mid \{ \mathbf{B} $ ould., s \mathbf{W} im., \mathbf{G} uitar, \mathbf{S} ing. $\}$
p	{ H iking, R unning}
m	{Bould., Hiking, Running}
r	{Bould., sKiing, Hiking}



Filters for the Set Similarity Join Signatures for Overlap Similarity

Identity Signature

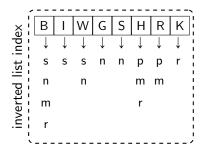
• Simplest signature scheme (for overlap) is identity (Sign = Id):

$$|r \cap s| \ge t \Rightarrow \operatorname{Id}(r) \cap \operatorname{Id}(s) \ne \emptyset$$

 $\Leftrightarrow |r \cap s| \ge t \Rightarrow |r \cap s| \ge 1$ assuming $t \ge 1$

- Every token is a signature
- Example:

id	interests
S	{B, I, W}
n	{B, W, G, S}
р	{H, R}
m	{B, H, R}
r	{B, K, H}



Augsten (Univ. Salzburg

WS 2025/26

Filters for the Set Similarity Join Signatures for Overlap Similarity

Demonstration

- Experiment: Identity signature¹
 - self-join with varying |R|
 - average set size 10
 - universe sizes |U| = 1000 and |U| = 10000, uniformly distributed
 - overlap similarity with threshold t = 4

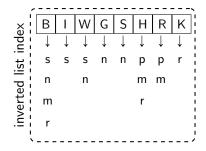
R	U	#comparisons	runtime [s]
1000	1000	$4.7 \cdot 10^4$	0.0
	10000	$4.8 \cdot 10^{3}$	0.0
10000	1000	$4.7 \cdot 10^{6}$	0.01
	10000	$5 \cdot 10^5$	0.005
100000	1000	$4.7 \cdot 10^{8}$	1.6
	10000	$4.9 \cdot 10^{7}$	0.19
1000000	1000	$4.7 \cdot 10^{10}$	241
	10000	$4.9 \cdot 10^9$	23

- Large improvement over naive approach
- Universe size heavily influences performance

Filters for the Set Similarity Join Signatures for Overlap Similarity

Identity Signature Example

id	interests
S	{B, I, W}
n	{B, W, G, S}
р	{H, R}
m	{B, H, R}
r	{B, K, H}



- Probing:
 - 1. Set s: $C \leftarrow \emptyset \cup \{(s, n), (s, m), (s, r)\}$
 - 2. Set $n: C \leftarrow C \cup \{(n, m), (n, r)\}$
 - 3. Set p: $C \leftarrow C \cup \{(p, m), (p, r)\}$
 - 4. Set $m: C \leftarrow C \cup \{(m,r)\}$
- 8 (non-reflexive, non-symmetric) comparisons!
- Most candidates are the result of the long list B.

Augsten (Univ. Salzburg)

WS 2025/26

Prefix Signature

• Idea: Exploit token order to construct a signature that is based on a subset of tokens.

Filters for the Set Similarity Join Signatures for Overlap Similarity

Definition (Prefix Signature Scheme)

The prefix signature Pre(r) of a set r for overlap threshold t is constructed as follows:

- 1. Order the tokens of r by any fixed global^a order.
- 2. Each of the first |r| t + 1 tokens in the ordered set is a prefix signature.

^aglobal: the same order must be used for all sets

• Example:

Set n,
$$t = 2$$

G W

1. order (e.g., alphabetically)



2. take first 4 - 2 + 1 = 3 tokens B • $Pre(n) = \{B, G, S\}$

¹Implementation includes some optimizations compared to framework

Filters for the Set Similarity Join Signatures for Overlap Similarity

Prefix Signature Correctness I

Lemma

Pre is a signature scheme for overlap similarity, i.e.,

$$|r \cap s| \geqslant t \Rightarrow |\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| \geqslant 1$$

We show the contraposition, i.e.,

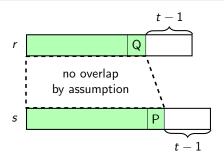
$$|\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| = 0 \Rightarrow |r \cap s| < t$$

Augsten (Univ. Salzburg)

WS 2025/26

Filters for the Set Similarity Join Signatures for Overlap Similarity

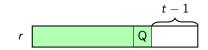
Prefix Signature Correctness II



- ullet Consider sets r, s; Q and P are the largest tokens in the respective prefixes, wlog. assume P < Q.
- Assume $Pre(r) \cap Pre(s) = \emptyset$. We bound $|r \cap s|$:
 - $|\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| = 0$ by assumption
 - $|(r \setminus \operatorname{Pre}(r)) \cap \operatorname{Pre}(s)| = 0$ as P < Q
 - $|r \cap (s \setminus \operatorname{Pre}(s))| \leq t 1$ as $|s \setminus \operatorname{Pre}(s)| = t 1$
 - Hence, $|r \cap s| < t$.

Filters for the Set Similarity Join Signatures for Overlap Similarity

Prefix Signature Correctness II





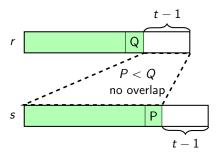
- Consider sets r, s; Q and P are the largest tokens in the respective prefixes, wlog. assume P < Q.
- Assume $Pre(r) \cap Pre(s) = \emptyset$. We bound $|r \cap s|$:
 - $|\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| = 0$ by assumption
 - $|(r \setminus Pre(r)) \cap Pre(s)| = 0$ as P < Q
 - $|r \cap (s \setminus \operatorname{Pre}(s))| \leq t 1$ as $|s \setminus \operatorname{Pre}(s)| = t 1$
 - Hence, $|r \cap s| < t$.

Similarity Search Augsten (Univ. Salzburg)

WS 2025/26

Filters for the Set Similarity Join Signatures for Overlap Similarity

Prefix Signature Correctness II

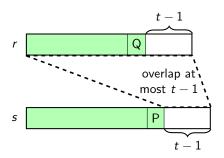


- \bullet Consider sets r, s; Q and P are the largest tokens in the respective prefixes, wlog. assume P < Q.
- Assume $Pre(r) \cap Pre(s) = \emptyset$. We bound $|r \cap s|$:
 - $|\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| = 0$ by assumption
 - $|(r \setminus \operatorname{Pre}(r)) \cap \operatorname{Pre}(s)| = 0$ as P < Q
 - $|r \cap (s \setminus Pre(s))| \le t 1$ as $|s \setminus Pre(s)| = t 1$
 - Hence, $|r \cap s| < t$.

Augsten (Univ. Salzburg)

Filters for the Set Similarity Join Signatures for Overlap Similarity

Prefix Signature Correctness II



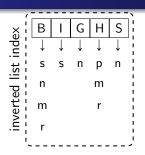
- Consider sets r, s; Q and P are the largest tokens in the respective prefixes, wlog. assume P < Q.
- Assume $Pre(r) \cap Pre(s) = \emptyset$. We bound $|r \cap s|$:
 - $|\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| = 0$ by assumption
 - $|(r \setminus Pre(r)) \cap Pre(s)| = 0$ as P < Q
 - $|r \cap (s \setminus \mathsf{Pre}(s))| \leq t 1$ as $|s \setminus \mathsf{Pre}(s)| = t 1$
 - Hence, $|r \cap s| < t$.

Filters for the Set Similarity Join Signatures for Overlap Similarity

WS 2025/26

Prefix Signature Example

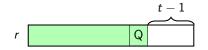
id	interests (ordered alphabetically)
S	{B, I, W}
n	{B, G, S, W}
р	{H, R}
m	{B, H, R}
r	{B, H, K}
	•



- Overlap threshold t = 2
- Indexing: all tokens except the last, alphabetical order
- Probing:
 - 1. Set s: $C \leftarrow \emptyset \cup \{(s,n),(s,m),(s,r)\}$
 - 2. Set $n: C \leftarrow C \cup \{(n, m), (n, r)\}$
 - 3. Set p: $C \leftarrow C \cup \{(p, m), (p, r)\}$
 - 4. Set $m: C \leftarrow C \cup \{(m,r)\}$
- still 8 comparisons!
- removing B from the prefix could help

Filters for the Set Similarity Join Signatures for Overlap Similarity

Prefix Signature Correctness II





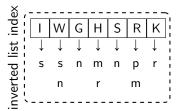
- Consider sets r, s; Q and P are the largest tokens in the respective prefixes, wlog. assume P < Q.
- Assume $Pre(r) \cap Pre(s) = \emptyset$. We bound $|r \cap s|$:
 - $|\operatorname{Pre}(r) \cap \operatorname{Pre}(s)| = 0$ by assumption
 - $|(r \setminus Pre(r)) \cap Pre(s)| = 0$ as P < Q
 - $|r \cap (s \setminus \mathsf{Pre}(s))| \le t 1$ as $|s \setminus \mathsf{Pre}(s)| = t 1$
 - Hence, $|r \cap s| < t$.

WS 2025/26

Filters for the Set Similarity Join Signatures for Overlap Similarity

Prefix Signature Example: Order Makes a Difference

id	interests (ordered by frequenc
S	{I, W, B}
n	{G, S, W, B}
p	{R, H}
m	{R, H, B}
r	{K, H, B}



- Overlap threshold t = 2
- Indexing: all tokens except the last, ordered by ascending frequency
- Probing:
 - 1. Set s: $C \leftarrow \emptyset \cup \{(s,n)\}$
 - 2. Set n: $C \leftarrow C \cup \emptyset$
 - 3. Set $p: C \leftarrow C \cup \{(p, m)\}$ 4. Set m: $C \leftarrow C \cup \{(m,r)\}$
- only 3 comparisons!
- Heuristic: Ordering by ascending token frequency reduces candidates

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Augsten (Univ. Salzburg)

Two Distributions

10^{0} zipf to ken by 10^{-3} 10^{-3} 10^{-3} uniform

Filters for the Set Similarity Join Signatures for Overlap Similarity

• Distribution: Real-world set-data often follow a zipfian distribution

tokens

 10^{1}

 10^{2}

• Skew: Some tokens appear frequently, a large number of tokens is uncommon. This favors the prefix signature.

Filters for the Set Similarity Join Signatures for Overlap Similarity

Augsten (Univ. Salzburg)

Augsten (Univ. Salzburg)

WS 2025/26

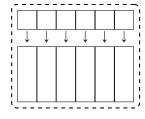
 10^{3}

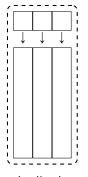
Impact of Universe Size

 10^{-4}

 10^{0}

- Identity and Prefix: individual tokens used as signatures
- Runtime: proportional to sum of all pairs in each list
- Problem: small universe size reduces filtering effectiveness





WS 2025/26

26 / 48

• Uniform distribution: halving the universe doubles the list lengths and runtime

Similarity Search

• Idea: use a more selective signature than individual tokens

Filters for the Set Similarity Join Signatures for Overlap Similarity

Demonstration

- Experiment: Identity signature vs. Prefix signature²
 - self-join with |R| = 100000
 - average set size 10
 - universe sizes |U| = 10000, uniform and zipfian distribution
 - overlap similarity with threshold $t \in \{4, 6, 8\}$
 - global order: ascending token frequency

Identity					Prefix			
dist.	t	#comp.	runtime [s]	dist.	t	#comp.	runtime [s]	
uni.	4	$5.0 \cdot 10^{7}$	0.187	uni.	4	$2.9 \cdot 10^{7}$	0.449	
	6	$5.0 \cdot 10^7$	0.186		6	$1.8 \cdot 10^7$	0.349	
	8	$5.0 \cdot 10^7$	0.186		8	$8.9 \cdot 10^6$	0.145	
zipf	4	$3.4 \cdot 10^{9}$	16.358	zipf	4	$3.1\cdot10^8$	3.935	
	6	$3.4 \cdot 10^{9}$	16.862		6	$6.2 \cdot 10^7$	0.873	
	8	$3.4 \cdot 10^9$	16.842		8	$1.2 \cdot 10^7$	0.197	

²Implementation includes some optimizations compared to framework

Similarity Search

WS 2025/26

Filters for the Set Similarity Join Signatures for Overlap Similarity

Subset Signature I

Augsten (Univ. Salzburg)

Lemma

$$|r \cap s| \geqslant t$$
 \Leftrightarrow

$$\exists p \subseteq U : |p| = t \land p \subseteq r \land p \subseteq s$$

• Similar sets have at least one common subset of size t. This proves correctness of the following signature:

Definition (Subsets Signature)

The subsets signature Sub(r) of a set r is defined as:

$$\mathsf{Sub}(r) = \{ p \subseteq r \mid |p| = t \}$$

for overlap threshold t.

Augsten (Univ. Salzburg) Similarity Search

Subset Signature II

• Sub is stronger than required for signature schemes. It also holds that

$$\mathsf{Sign}(r) \cap \mathsf{Sign}(s) \neq \emptyset \Rightarrow \mathsf{Sim}(r,s) \geqslant t.$$

Therefore, verification is not necessary.

- For set r and threshold t, we have $|\operatorname{Sub}(r)| = {r \choose t}$, growing very quickly depending on both |r| and t.
- Example: $n = \{B, W, G, S\}, t = 2$

$$Sub(n) = \{\{B, W\}, \{B, G\}, \{B, S\}\}$$
$$\{W, G\}, \{W, S\}, \{G, S\}\}$$

Augsten (Univ. Salzburg)

WS 2025/26

Filters for the Set Similarity Join Signatures for Hamming Distance

Outline

- Filters for the Set Similarity Join
 - Motivation
 - Signature-based Filtering
 - Signatures for Overlap Similarity
 - Signatures for Hamming Distance
- 2 Implementations of Set Similarity Joins
 - Other Similarity Functions
 - Table of Set Similarity Join Algorithms and their Signatures
- Conclusion

Filters for the Set Similarity Join Signatures for Overlap Similarity

Demonstration

• Experiment: Prefix signature vs. Subsets signature³

- self-join with |R| = 1000000
- average set size 6
- universe sizes $|U| \in \{1000, 10000\}$, uniform distribution
- overlap similarity with threshold $t \in \{3, 4, 5\}$

Prefix			Subsets		
U	t	runtime [s]	U	t	runtime [s]
1000	3	336	1000	3	25.6
	4	233		4	41.0
	5	138		5	42.8
10000	3	40.7	10000	3	32.0
	4	24.7		4	53.9
	5	14.9		5	66.2

- Sub outperforms Pre for small set sizes and small universes
- Pre scales better wrt. set size and threshold for large universes

Similarity Search

Filters for the Set Similarity Join Signatures for Hamming Distance

WS 2025/26

Partitioning I

Augsten (Univ. Salzburg)

Definition (Partition)

A partition P of universe U is a family of sets $P = \{p_1, \dots, p_n\}$ with the following properties:

- 1. Ø *∉ P*
- 2. $\bigcup_{p \in P} p = U$

Augsten (Univ. Salzburg)

3. For any $p_i, p_j, i \neq j$, we have $p_i \cap p_i = \emptyset$

Lemma

For any partition P of universe U and any two sets $r, s \subseteq U$, we have

$$Ham(r,s) = \sum_{p \in P} Ham(r \cap p, s \cap p)$$

Augsten (Univ. Salzburg) WS 2025/26 30 / 48 Similarity Search

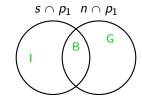
Similarity Search

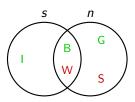
³Implementation includes some optimizations compared to framework

Partition Signature

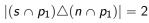
Partitioning II

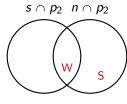
• Example: $P = \{\{B, G, H, I\}, \{K, R, S, W\}\} = \{p_1, p_2\}$





 $|s\triangle n|=3$





 $|(s \cap p_2) \triangle (n \cap p_2)| = 1$

Filters for the Set Similarity Join Signatures for Hamming Distance

Correctness of the Partition Signature

Lemma (Correctness of Par)

$$Ham(r,s) \leqslant t \Rightarrow Par(r) \cap Par(s) \neq \emptyset$$

We show the contraposition $Par(r) \cap Par(s) = \emptyset \Rightarrow Ham(r, s) > t$

Proof.

Assume $Par(r) \cap Par(s) = \emptyset$.

- For any $p_i \in P$, if $r \cap p_i \neq s \cap p_i$, then $Ham(r \cap p_i, s \cap p_i) \geqslant 1$.
- Hence, $Ham(r,s) = \sum_{p \in P} Ham(r \cap p, s \cap p) \geqslant |P| = t + 1 > t$.

Definition (Partition Signature)

The partition signature Par(r) of a set r is constructed as follows:

1. Fix any partition P of U into t+1 parts (for Hamming distance t)

Filters for the Set Similarity Join Signatures for Hamming Distance

- 2. $Par(r) = \{(r \cap p_i, i) \mid p_i \in P\}$
- Example: $P = \{\{B, G, H, I\}, \{K, R, S, W\}\} = \{p_1, p_2\}$
 - $Par(\{I, B, W\}) = \{(\{I, B\}, 1), (\{W\}, 2)\}$
 - $Par(\{S, W\}) = \{(\emptyset, 1), (\{S, W\}, 2)\}$

Augsten (Univ. Salzburg)

Filters for the Set Similarity Join Signatures for Hamming Distance

WS 2025/26

Demonstration

- Experiment: Prefix signature⁴ vs. Partition signature
 - self-join with |R| = 1000000
 - average set size 20
 - universe sizes $|U| \in \{1000, 10000\}$, uniform distribution
 - Hamming distance with threshold $t \in \{3, 4, 5\}$

	efix	Partition			
U	t	runtime [s]	U	t	runtime [s]
1000	3	326	1000	3	4
	4	451		4	23
	5	576		5	144
10000	3	34	10000	3	4
	4	46		4	23
	5	70		5	144

- Par outperforms Pre for large set sizes and small universes
- Par is less sensitive to universe size compared to Pre
- Pre works better for large universes and small set sizes

⁴Adapted to work with Hamming distance

Augsten (Univ. Salzburg)

WS 2025/26 35 / 48

Filters for the Set Similarity Join Signatures for Hamming Distance

The Empty Par Signature

• Consider partitioning $P = \{\{A, B, C\}, \{D, E, F\}, \{G, H, I\}\}, t = 2$:

	Α	В	С	D	Ε	F	G	Н	1
r	1	1					1	1	
s	1		1				1		1
и		1	1					1	1

- Although all pairs of sets are at Hamming distance 4, they all share the red signature.
- Hence, all pairs of sets are candidates!
- This can happen for small sets or heavily skewed distributions.
- More sophisticated partitioning and more flexible searching in each partition can remedy this problem.

Augsten (Univ. Salzburg)

WS 2025/26

Filters for the Set Similarity Join Signatures for Hamming Distance

Enumeration Signature

Definition (Enumeration Signature)

The enumeration signature En(r) of a set r for Hamming distance with threshold t is given by:

$$\mathsf{En}(r) = \{ p \subseteq r \mid |p| \geqslant |r| - t \}$$

• Example: t = 2

Augsten (Univ. Salzburg)

- $En(\{I, B, G\}) = \{\{I, B, G\}, \{I, B\}, \{I, G\}, \{B, G\}, \{I\}, \{B\}, \{G\}\}\}$
- $En(\{H,R\}) = \{\{H,R\},\{H\},\{R\},\emptyset\}$

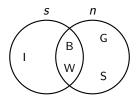
Lemma (Correctness of En)

$$Ham(r,s) \leqslant t \Rightarrow En(r) \cap En(s) \neq \emptyset$$

Filters for the Set Similarity Join Signatures for Hamming Distance

Enumeration: Idea

• How can we make s and n the same?



- By removing I from s and both S and G from n.
- Idea: By removing at most t tokens, all pairs of sets in Hamming distance t can be made equal.

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Filters for the Set Similarity Join Signatures for Hamming Distance

Correctness of the Enumeration Signature

Proof.

- Assume $Ham(r, s) = |(r \setminus s) \cup (s \setminus r)| \leq t$.
- Hence, $|r \setminus (r \setminus s)| \ge |r| t$ and $|s \setminus (s \setminus r)| \ge |s| t$
- Consider the set $r \setminus (r \setminus s) = r \cap s = s \setminus (s \setminus r)$
- As $r \cap s \subseteq r$ and $|r \cap s| \ge |r| t$, $r \cap s \in \text{En}(r)$.
- Similarly, $r \cap s \in \text{En}(s)$.
- So $r \cap s \in \operatorname{En}(r) \cap \operatorname{En}(s)$.

Similarity Search WS 2025/26 Augsten (Univ. Salzburg)

Demonstration

- Experiment: Partition signature vs. Enumeration signature⁵
 - self-join with |R| = 1000000
 - average set size $|r| \in \{8, 16\}$
 - universe size |U| = 10000, uniform distribution
 - Hamming distance with threshold $t \in \{1, 2, 3\}$

Partition				Enumeration			
r	t	runtime [s]	Ī	$\overline{r }$	t	runtime [s]	
8	1	4		8	1	12	
	2	141			2	42	
	3	1034			3	165	
16	1	1	1	.6	1	21	
	2	3			2	170	
	3	14			3	(out of memory)	

- En can outperform Par for small thresholds and set sizes
- For large thresholds and sets, En generates too many signatures

Augsten (Univ. Salzburg)

Similarity Search

Implementations of Set Similarity Joins Other Similarity Functions

WS 2025/26

WS 2025/26

Outline

- Filters for the Set Similarity Join
 - Motivation
 - Signature-based Filtering
 - Signatures for Overlap Similarity
 - Signatures for Hamming Distance
- 2 Implementations of Set Similarity Joins
 - Other Similarity Functions
 - Table of Set Similarity Join Algorithms and their Signatures
- Conclusion

Implementations of Set Similarity Joins

Implementations of Set Similarity Joins

Real implementations of set similarity join algorithms typically

- also support similarity functions other than overlap and Hamming
- use a combination of multiple signature schemes
- extend the algorithmic framework to optimize for their signature schemes
- use additional filters (e.g., based on set length or the positions of matching signatures)

Augsten (Univ. Salzburg)

Other Similarity Functions

- Normalization: often, normalized similarity functions are preferred
 - $r = \{A, B, C\}, s = \{A, B\}, u = \{A, B, C, D, E, F, G, H, I\}$

Implementations of Set Similarity Joins Other Similarity Functions

- The pair (r, u) has higher overlap than the pair (r, s)
- Still, (r, s) might appear more similar due to fewer different tokens
- Normalizations also consider set sizes and take values in [0, 1]
- Jaccard: $Jac(r,s) = \frac{|r \cap s|}{|r \cup s|} = \frac{|r \cap s|}{|r| + |s| |r \cap s|}$
- Dice: Dice $(r,s) = \frac{2|r \cap s|}{|r|+|s|}$
- Cosine: $Cos(r, s) = \frac{|r \cap s|}{\sqrt{|r| \cdot |s|}}$
- Example:

Augsten (Univ. Salzburg)

$$Jac(r,s) = \frac{|\{A,B\}|}{|\{A,B,C\}|} = \frac{2}{3}$$

$$Jac(r,u) = \frac{|\{A,B,C\}|}{|\{A,B,C,D,E,F,G,H,I\}|} = \frac{3}{9} = \frac{1}{3}$$

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

42 / 48

Similarity Search

43 / 48

⁵Implemented using an optimization called asymmetric signature scheme that avoids false positives.

Adapting the Prefix Signature for Jaccard

- The prefix signature operates with overlap similarity
- Idea: Bound minimum overlap s.t. two sets r, s can have $Jac(r, s) \ge t$

$$\frac{|r \cap s|}{|r| + |s| - |r \cap s|} \ge t$$

$$\Leftrightarrow \qquad |r \cap s| \ge t(|r| + |s| - |r \cap s|)$$

$$\Leftrightarrow \qquad |r \cap s| \ge \frac{t}{1 + t}(|r| + |s|) =: eqo_J(r, s)$$

- eqo, depends on the sizes of two sets. As we want to handle all possible size combinations, we have to get rid of one of them.
- Possible Solution: Assume minimal value |s| = 1, but better bounds are possible.

Augsten (Univ. Salzburg)

Similarity Search

Implementations of Set Similarity Joins Table of Set Similarity Join Algorithms and their Signatures

Outline

- Filters for the Set Similarity Join
 - Motivation
 - Signature-based Filtering
 - Signatures for Overlap Similarity
 - Signatures for Hamming Distance
- 2 Implementations of Set Similarity Joins
 - Other Similarity Functions
 - Table of Set Similarity Join Algorithms and their Signatures
- Conclusion

Length Bounds for Jaccard

Lemma

If $Jac(r, s) \ge t$, then $t|r| \le |s| \le \frac{|r|}{t}$.

Lemma

If $Jac(r, s) \ge t$, then

$$|r \cap s| \ge \operatorname{eqo}_{J}(r, s)$$

 $\ge \operatorname{eqo}_{J}(r, t|r|)$
 $= t|r|$

• Hence, for each set r use $\lceil t | r \rceil$ as the overlap and proceed as in Pre.

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Implementation	ons of Set Similarity Joins	Table of Set Similarity Join Algorithms and their Signatures
Algorithm	Signature ⁶	Remarks
AllPairs [BMS07] PPJoin [XWLY08]	Length+Pre Length+Pre	additional filter based on position of pre- fix matches
SkipJoin [WQL ⁺ 19]	Length+Pre	tighter length filter using prefix po- sitions, removes unmatchable entries from index, can leverage knowledge of set similarity "transitively"
SizeAware [DTL18]	Sub (small sets)+ Id (large sets)	· · · · · · · · · · · · · · · · · · ·
PartEnum [AGK06]	Length + Par + En	partitions sets into smaller subsets, enumerates in each partition
PartAlloc [DLWF15]	Length + Par + En	more flexible than PartEnum, has tighter filtering condition
GPH [QXW ⁺ 21]	Par + En	more flexible than PartAlloc, optimizes how partitions are chosen

⁶Refers to the closest signature discussed during the lecture; the listed algorithms often use a more efficient variation of the respective signatures.

Augsten (Univ. Salzburg) WS 2025/26 46 / 48 Similarity Search

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Conclusion

Summary

- Naive set similarity join inefficient due to large search space
- Signature-based filters speed up join:
 - Id and Pre: single tokens as signatures
 - Sub: all subsets of overlap size as signatures
 - Par: non-overlapping subsets as signatures
 - En: all subsets in Hamming distance as signatures
- Performance depends on dataset's characteristics

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

48 /

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

48 / 4

Dong Deng, Yufei Tao, and Guoliang Li.
Overlap set similarity joins with theoretical guarantees.
In Gautam Das, Christopher M. Jermaine, and Philip A. Bernstein, editors, Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018, pages 905–920. ACM, 2018.

Jianbin Qin, Chuan Xiao, Yaoshu Wang, Wei Wang, Xuemin Lin, Yoshiharu Ishikawa, and Guoren Wang.

Generalizing the pigeonhole principle for similarity search in hamming space.

IEEE Trans. Knowl. Data Eng., 33(2):489-505, 2021.

- Xubo Wang, Lu Qin, Xuemin Lin, Ying Zhang, and Lijun Chang. Leveraging set relations in exact and dynamic set similarity join. *VLDB J.*, 28(2):267–292, 2019.
- Chuan Xiao, Wei Wang, Xuemin Lin, and Jeffrey Xu Yu. Efficient similarity joins for near duplicate detection.

Arvind Arasu, Venkatesh Ganti, and Raghav Kaushik. Efficient exact set-similarity joins.

In Umeshwar Dayal, Kyu-Young Whang, David B. Lomet, Gustavo Alonso, Guy M. Lohman, Martin L. Kersten, Sang Kyun Cha, and Young-Kuk Kim, editors, *Proceedings of the 32nd International Conference on Very Large Data Bases, Seoul, Korea, September 12-15, 2006*, pages 918–929. ACM, 2006.

Roberto J. Bayardo, Yiming Ma, and Ramakrishnan Srikant. Scaling up all pairs similarity search.

In Carey L. Williamson, Mary Ellen Zurko, Peter F. Patel-Schneider, and Prashant J. Shenoy, editors, *Proceedings of the 16th International Conference on World Wide Web, WWW 2007, Banff, Alberta, Canada, May 8-12, 2007*, pages 131–140. ACM, 2007.

Dong Deng, Guoliang Li, He Wen, and Jianhua Feng. An efficient partition based method for exact set similarity joins. *Proc. VLDB Endow.*, 9(4):360–371, 2015.

In Jinpeng Huai, Robin Chen, Hsiao-Wuen Hon, Yunhao Liu, Wei-Ying Ma, Andrew Tomkins, and Xiaodong Zhang, editors, *Proceedings of the 17th International Conference on World Wide Web, WWW 2008, Beijing, China, April 21-25, 2008*, pages 131–140. ACM, 2008.

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 48 / 48 Augsten (Univ. Salzburg) Similarity Search WS 2025/26 48 / 48