

# Speed up Entity Resolution with Bit Arrays

## **Big Data Praktikum SS 17**

Moritz Engelmann  
Maik Fröbe

31.07.2017

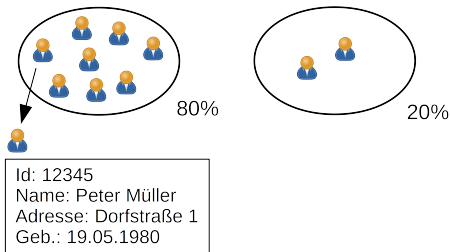
- 1 Einführung
- 2 Ansätze zur Entity-Resolution
  - Trivialer Ansatz
  - Sortier-Ansatz
  - Bit-Array-Ansatz
  - Vergleich
- 3 Optimierungen des Bit-Array-Ansatzes
  - Bit-Array als Filter
  - Filtern in 2 Phasen
- 4 Rückblick + Ausblick

# Einführung

# Problembeschreibung

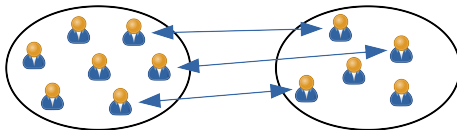
## • Eingabe

- 2 Mengen von Personen
- Verhältnis 80:20



## • Ziel

- finden ähnlicher Personen in beiden Datensätzen



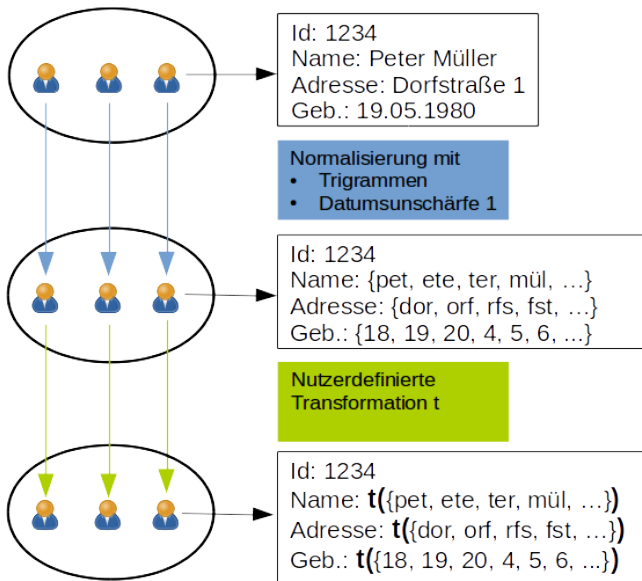
## • Anforderungen

- Bestimmung der Ähnlichkeit mit Jaccard-Index
- $jaccard(A, B) = \frac{|A \cap B|}{|A \cup B|}$

# Framework zur Entity-Resolution

- vollständig Parametrisierbar
- Modular
- Start der Entity-Resolution mit:
  - Transformation:  $\text{Person} \rightarrow V$
  - Ähnlichkeit:  $V \times V \rightarrow [0, 1]$
  - $n$  (Größe der  $n$ -Gramme)
  - Threshold
  - ...
- sequentieller Nested-Loop
  - Vollständige Berechnung der Ähnlichkeit für kartesisches Produkt

# Importphase



# Ansätze zur Entity-Resolution

## Trivialer Ansatz



## Transformation:



Id: 1234  
 Name: {pet, ete, ter, ...}  
 Adresse: {dor, orf, rfs, ...}  
 Geb.: {18, 19, 20, 4, 5, ...}

Transformation

Id: 1234  
 Name: {pet, ete, ter, ...}  
 Adresse: {dor, orf, rfs, ...}  
 Geb.: {18, 19, 20, 4, 5, ...}

## Ähnlichkeitsfunktion:



Id: 1234  
 Name: {pet, ete, ter, müll, üll, ...}  
 Adresse: {dor, orf, rfs, fst, ...}  
 Geb.: {18, 19, 20, 4, 5, 6, ...}



Id: 1234  
 Name: {pet, ete, ter, mue, ...}  
 Adresse: {dor, orf, rfs, fst, ...}  
 Geb.: {18, 19, 20, 4, 5, 6, ...}



$$\text{sim} = \text{avg}(\text{sim}_{\text{Jacc}}(\text{Name}), \text{sim}_{\text{Jacc}}(\text{Adresse}), \text{sim}_{\text{Jacc}}(\text{Geb.}))$$

## Sortier-Ansatz

- Analog zu Sort-Merge-Verbund<sup>1</sup>
  - Während Import: Sortiere Mengen
  - Während ER: Berechne Kardinalität der Schnittmenge in  $\mathcal{O}(n)$ 
    - Schritthaltende Traversierung der sortierten Mengen

---

<sup>1</sup>Siehe Vorlesung Implementierung von Datenbanksystemen

## Bit-Array-Ansatz

# Einschub: Bloom-Filter

pet

ete

ter

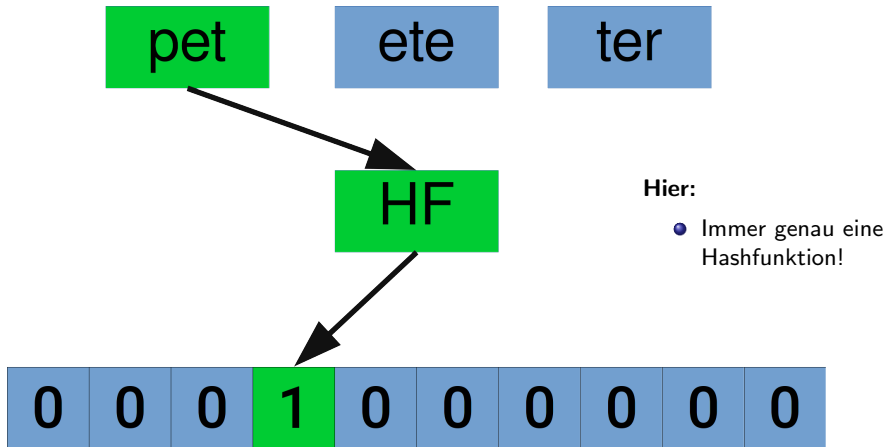
HF

Hier:

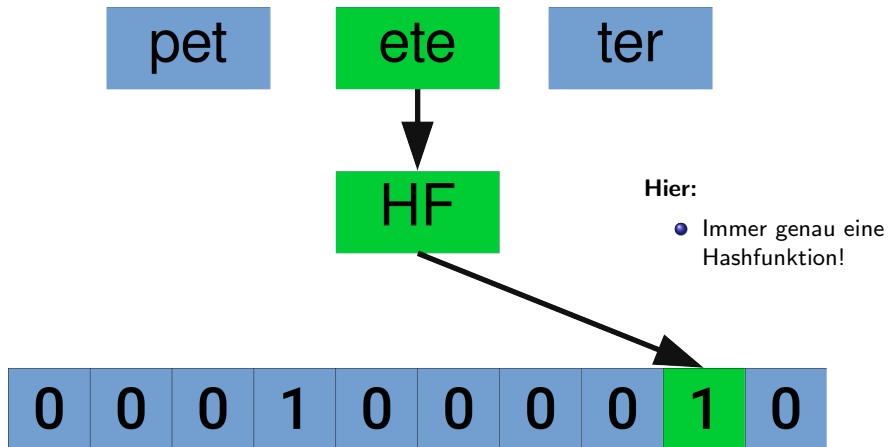
- Immer genau eine Hashfunktion!

0 0 0 0 0 0 0 0 0 0

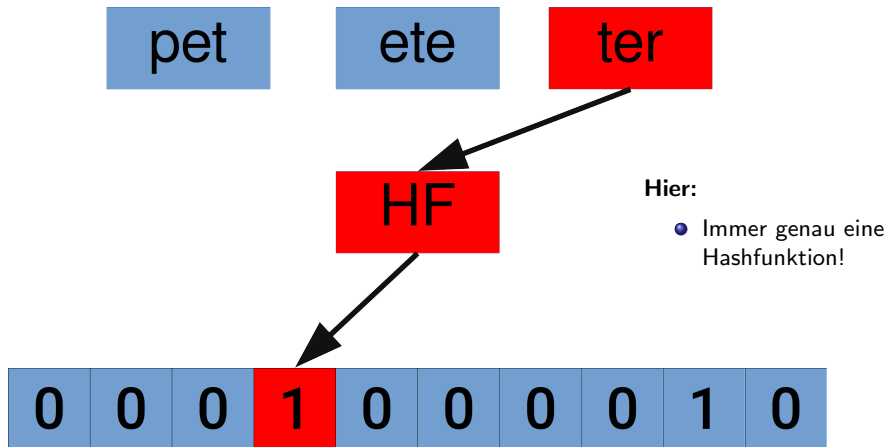
# Einschub: Bloom-Filter



# Einschub: Bloom-Filter



# Einschub: Bloom-Filter





## Transformation:



Id: 1234  
 Name: {pet, ete, ter, ...}  
 Adresse: {dor, orf, rfs, ...}  
 Geb.: {18, 19, 20, 4, 5, ...}

Transformation

Id: 1234  
 Name: 

0	0	0	1	0	1	0	0	1	0
---	---	---	---	---	---	---	---	---	---

  
 Adresse: 

0	1	0	0	1	0	0	0	0	1
---	---	---	---	---	---	---	---	---	---

  
 Geb.: 

0	0	1	1	0	0	1	0	1	0
---	---	---	---	---	---	---	---	---	---

## Ähnlichkeitsfunktion:



Id: 1234  
 Name: 

0	0	0	1	0	1	0	0	1	0
---	---	---	---	---	---	---	---	---	---

  
 Adresse: 

0	1	0	0	1	0	0	0	0	1
---	---	---	---	---	---	---	---	---	---

  
 Geb.: 

0	0	1	1	0	0	1	0	1	0
---	---	---	---	---	---	---	---	---	---



Id: 1234  
 Name: 

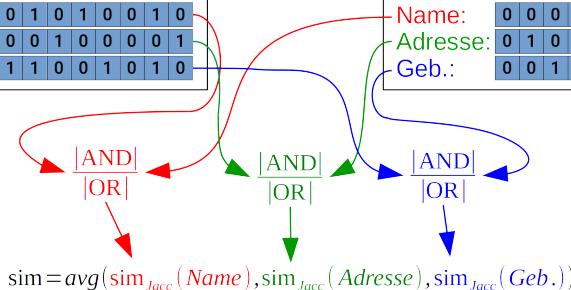
0	0	0	1	0	1	0	1	1	0
---	---	---	---	---	---	---	---	---	---

  
 Adresse: 

0	1	0	1	0	0	0	0	0	1
---	---	---	---	---	---	---	---	---	---

  
 Geb.: 

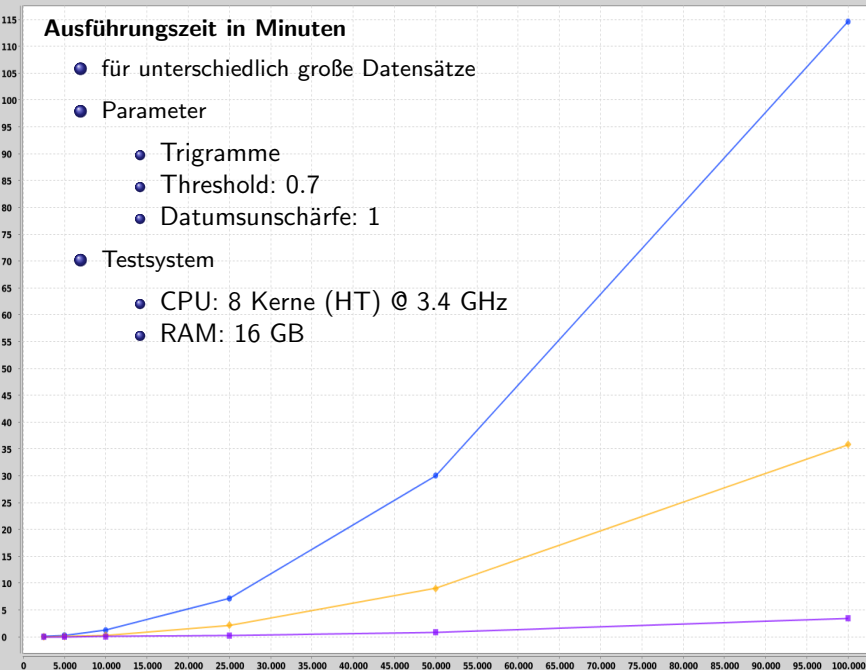
0	0	1	1	0	1	1	0	1	0
---	---	---	---	---	---	---	---	---	---



# Vergleich

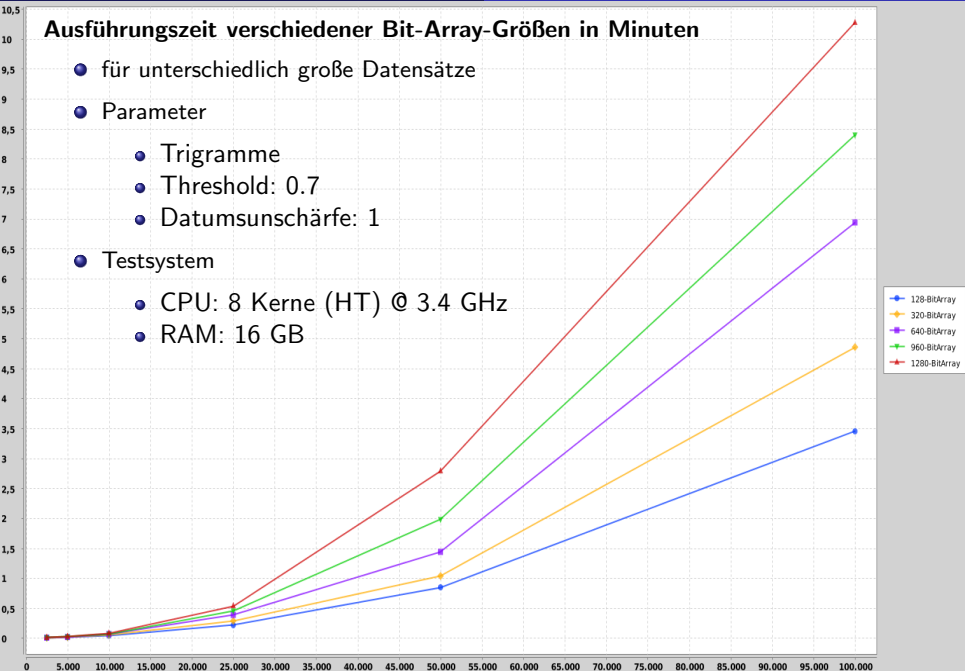
## Ausführungszeit in Minuten

- für unterschiedlich große Datensätze
- Parameter
  - Trigramme
  - Threshold: 0.7
  - Datumsunschärfe: 1
- Testsystem
  - CPU: 8 Kerne (HT) @ 3.4 GHz
  - RAM: 16 GB



## Ausführungszeit verschiedener Bit-Array-Größen in Minuten

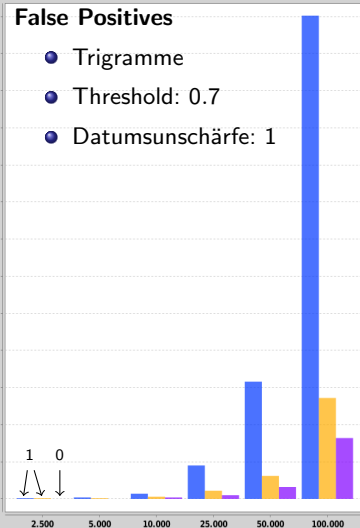
- für unterschiedlich große Datensätze
- Parameter
  - Trigramme
  - Threshold: 0.7
  - Datumsunschärfe: 1
- Testsystem
  - CPU: 8 Kerne (HT) @ 3.4 GHz
  - RAM: 16 GB



## Absolute Werte:

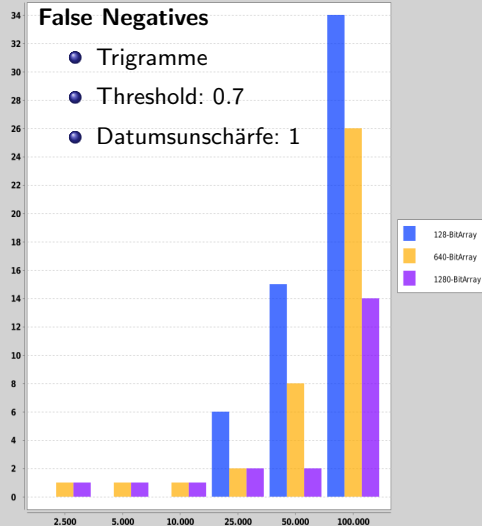
### False Positives

- Trigramme
- Threshold: 0.7
- Datumsunschärfe: 1

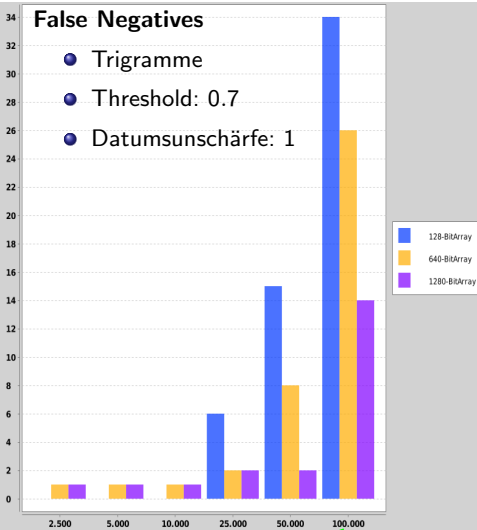
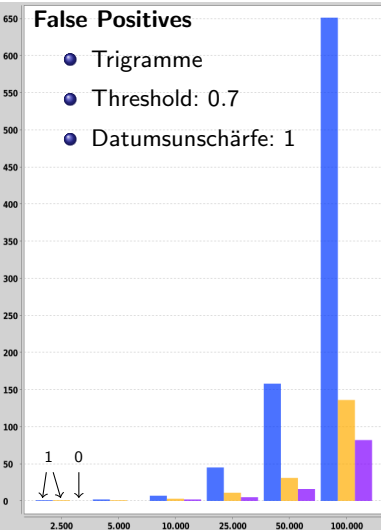


### False Negatives

- Trigramme
- Threshold: 0.7
- Datumsunschärfe: 1



## Absolute Werte:



	128-Bit	640-Bit	1280-Bit
Precision	0.955	0.990	0.994
Recall	0.998	0.998	0.999

# Optimierungen des Bit-Array-Ansatzes

## Bit-Array als Filter



- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

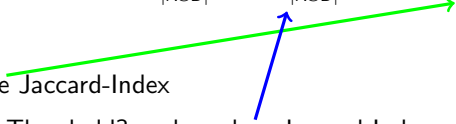
- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

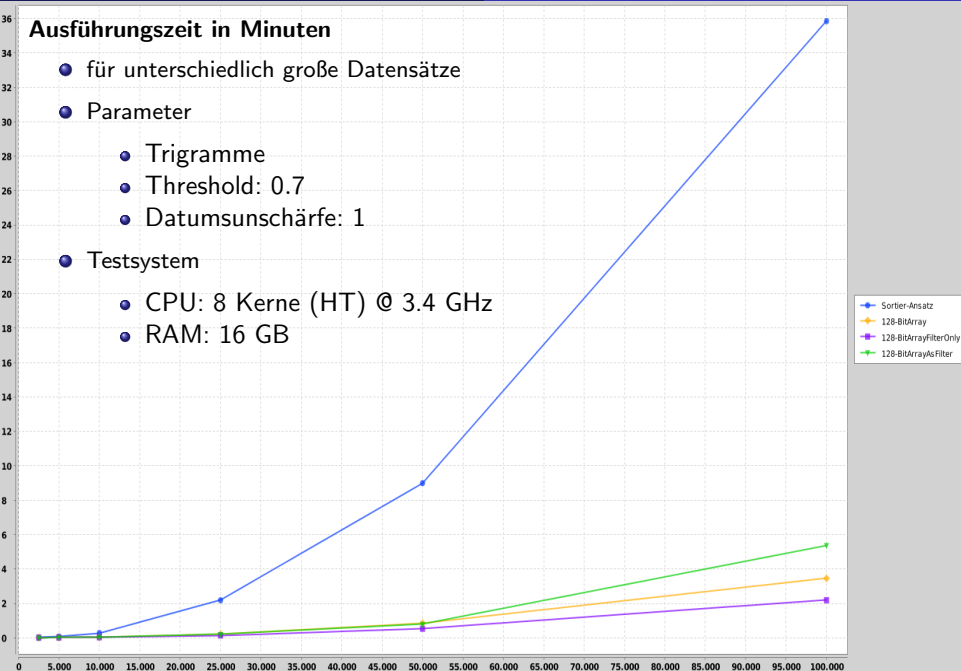


- Zwischenschritt:
  - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
- Idee:
  - $A, B$  - Mengen
  - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
  - Es gilt:
    - $|A_F| \leq |A|$
    - $|A_F \vee B_F| \leq |A \cup B|$
    - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
- Vorgehen:
  - Schätze Jaccard-Index
  - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index

- Zwischenschritt:
    - Obere Schranke des Jaccard-Index mit Bit-Arrays effizient bestimmen
  - Idee:
    - $A, B$  - Mengen
    - $A_F = \text{bloom}(A)$ ,  $B_F = \text{bloom}(B)$  - Bit-Arrays der Mengen
    - Es gilt:
      - $|A_F| \leq |A|$
      - $|A_F \vee B_F| \leq |A \cup B|$
      - $\text{jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A| + |B| - |A \cup B|}{|A \cup B|} \leq \frac{|A| + |B| - |A_F \vee B_F|}{|A_F \vee B_F|}$
  - Vorgehen:
    - Schätze Jaccard-Index
    - Größer Threshold?  $\Rightarrow$  berechne Jaccard-Index
- 

## Ausführungszeit in Minuten

- für unterschiedlich große Datensätze
- Parameter
  - Trigramme
  - Threshold: 0.7
  - Datumsunschärfe: 1
- Testsystem
  - CPU: 8 Kerne (HT) @ 3.4 GHz
  - RAM: 16 GB



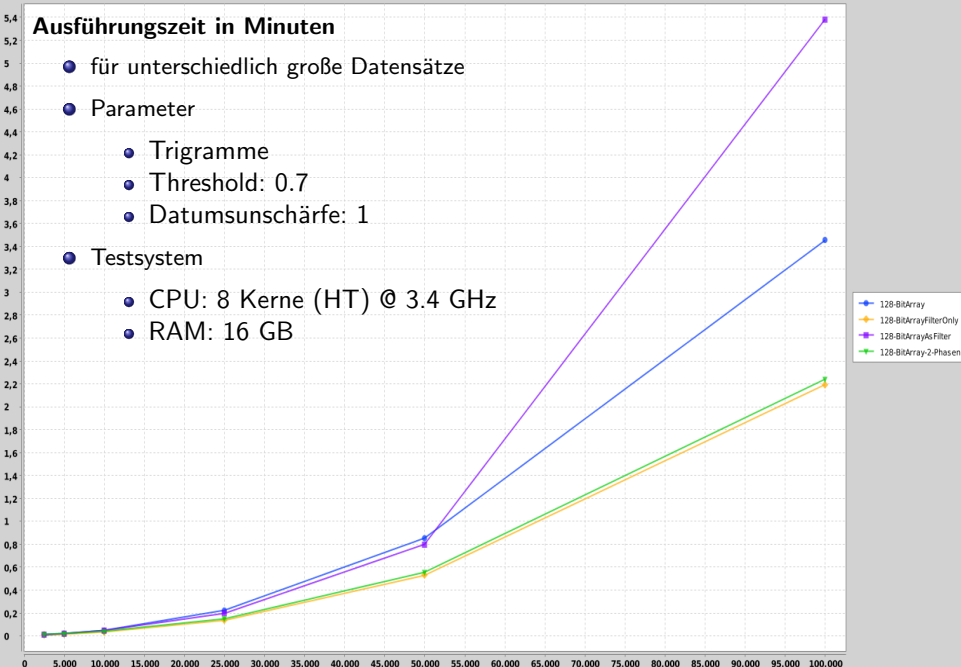
## Filtern in 2 Phasen

# Filtern in 2 Phasen

- Phase 1
  - ER mit Bit-Array-Filter-Only
- Phase 2
  - Eingabe: Id-Paare der Kandidaten aus Phase 1
  - Import (Normalisierung + Transformation) für Sortier-Ansatz aus
  - ER nur für Kandidaten-Paare

## Ausführungszeit in Minuten

- für unterschiedlich große Datensätze
- Parameter
  - Trigramme
  - Threshold: 0.7
  - Datumsunschärfe: 1
- Testsystem
  - CPU: 8 Kerne (HT) @ 3.4 GHz
  - RAM: 16 GB



## Rückblick + Ausblick

- Erfahrungen:
  - Abstraktion + große Datenstrukturen sind teuer
- Tipps für unbekannte Datensätze:
  - Untersuche Datensatz mit „Phase 1“
  - Schrittweise Anpassung der Parameter
    - Ziel: sehr hohe Selektivität trotz kleinem Bit-Array
  - Abschließend: ER mit „Phase 1 + 2“
- mögliche Nächste Schritte:
  - Parallelisierung
  - Ein Partitionierter Bit-Array
  - Vergleich mit weiteren ER-Ansätzen
  - Integration unterschiedlicher ER-Ansätze als „Phase 2“