



**Hasso
Plattner
Institut**

IT Systems Engineering | Universität Potsdam

Apriori Algorithmus

Endpräsentation

Stefan George, Felix Leupold

Gliederung

2

- Wiederholung Apriori
- Erweiterung: Parallelisierung
 - Parallele Programmierung in Python
 - Parallelisierungsszenarien
 - Implementierung
 - Ergebnisse/Benchmarks
- Usecase
 - Lift
 - Conviction
 - Ergebnisse

Wiederholung Apriori

3

- Algorithmus zum Auffinden von Assoziationsregeln
- Wichtige Parameter: Support und Confidence
- Apriori Überlegung:

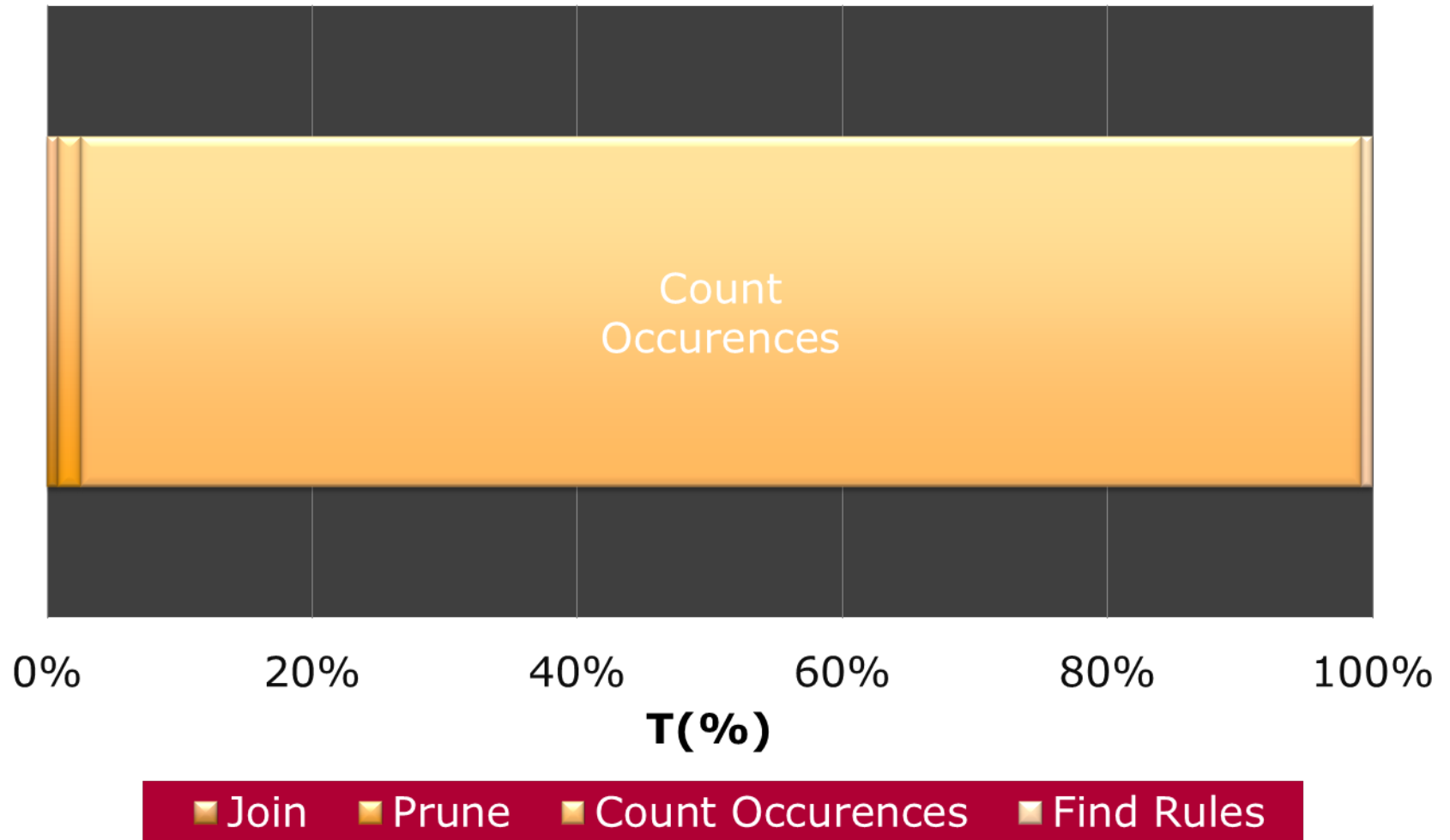
Wenn eine Teilmenge einer Menge M klein ist, dann ist die Menge M auch auf jeden Fall klein

- Kandidatengenerierung iterativ: Join und Prune
- Regeln aus Itemsets mit ausreichend Support generieren

$$A \rightarrow (L - A) \quad \Leftrightarrow \quad \frac{\text{sup}(L)}{\text{sup}(A)} > \min \text{Conf}$$

Motivation für Parallelisierung

4



Parallelisierung in Python

5

- Multiprocessing vs. Multithreading
- Threads haben „shared memory“
- Deutlich mehr overhead bei Context Switch zwischen Prozessen
- Global Interpreter Lock in Python
 - Nur ein Thread kann gleichzeitig Bytecode ausführen
 - Keine Parallelisierung bei CPU lastigem Code
- Für uns kommt nur Multiprocessing in Frage

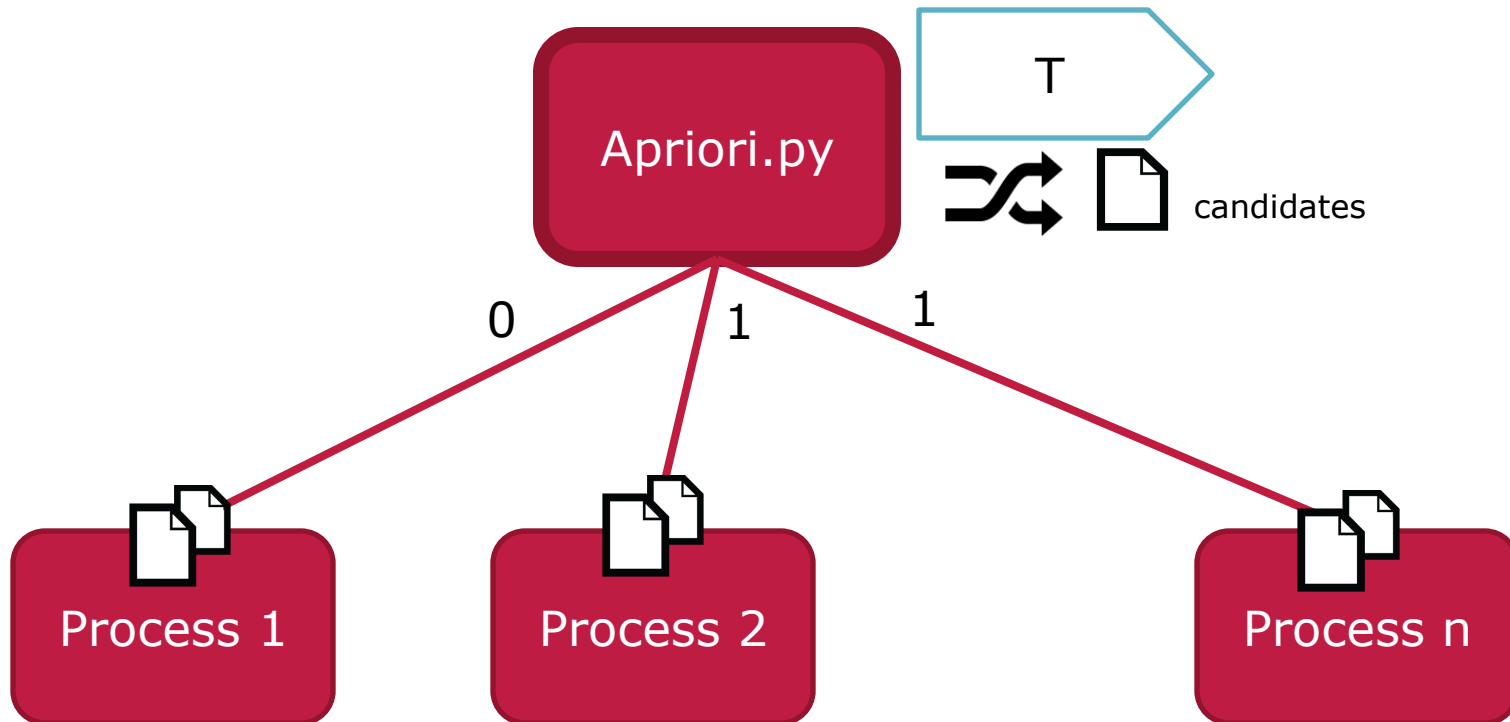
Parallelisierungsszenarien (1/3)

6



Parallelisierungsszenarien (1/3)

7



Analyse (1/3)

8

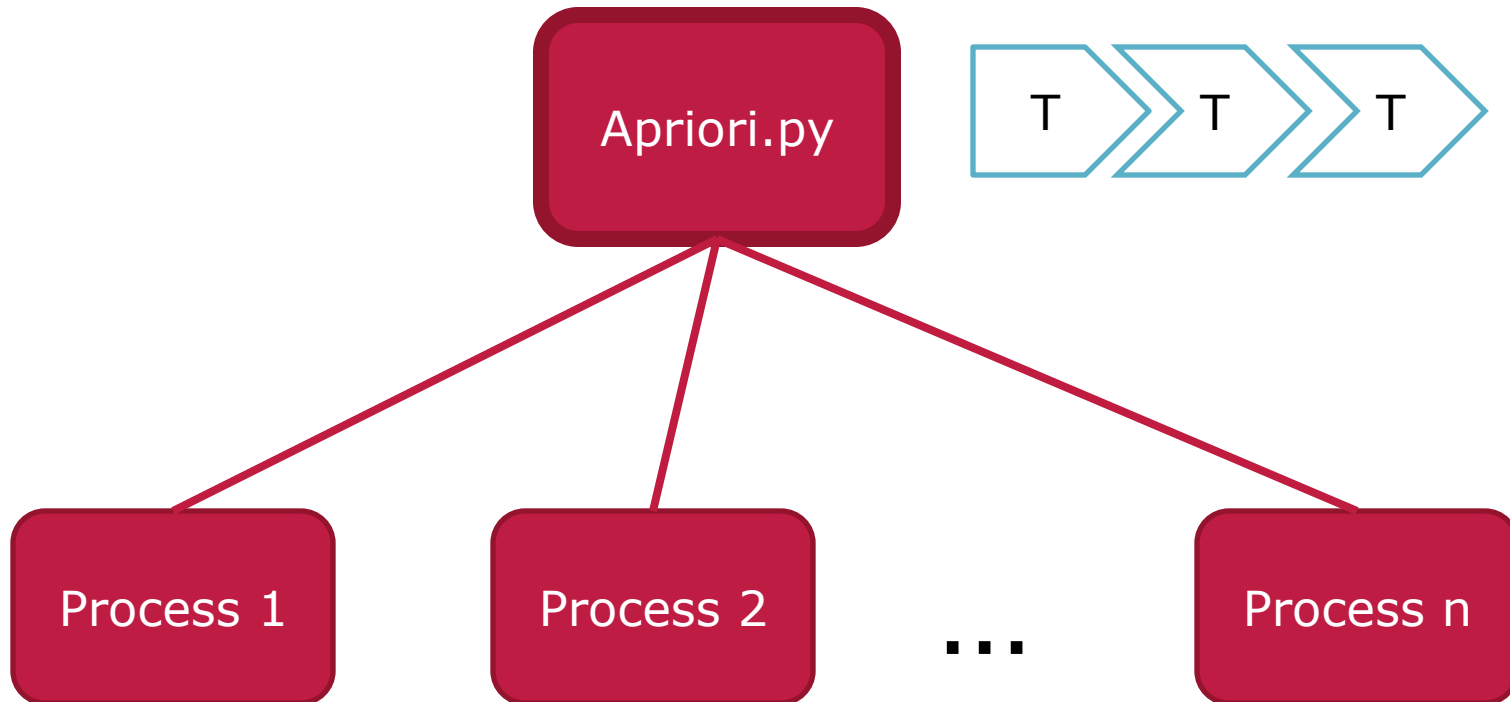
- Kein Initialisierungsaufwand 0
- Datenübertragung in Schritt k: $|D| * |C_k|$ 5Mio * 400
- Ergebnisgröße: $|C_k|$ 400

- Gesamtkosten = $\sum_k (|D| + 1) * |C_k|$ 2Mrd

- Parallelisierung muss grob granularer sein.

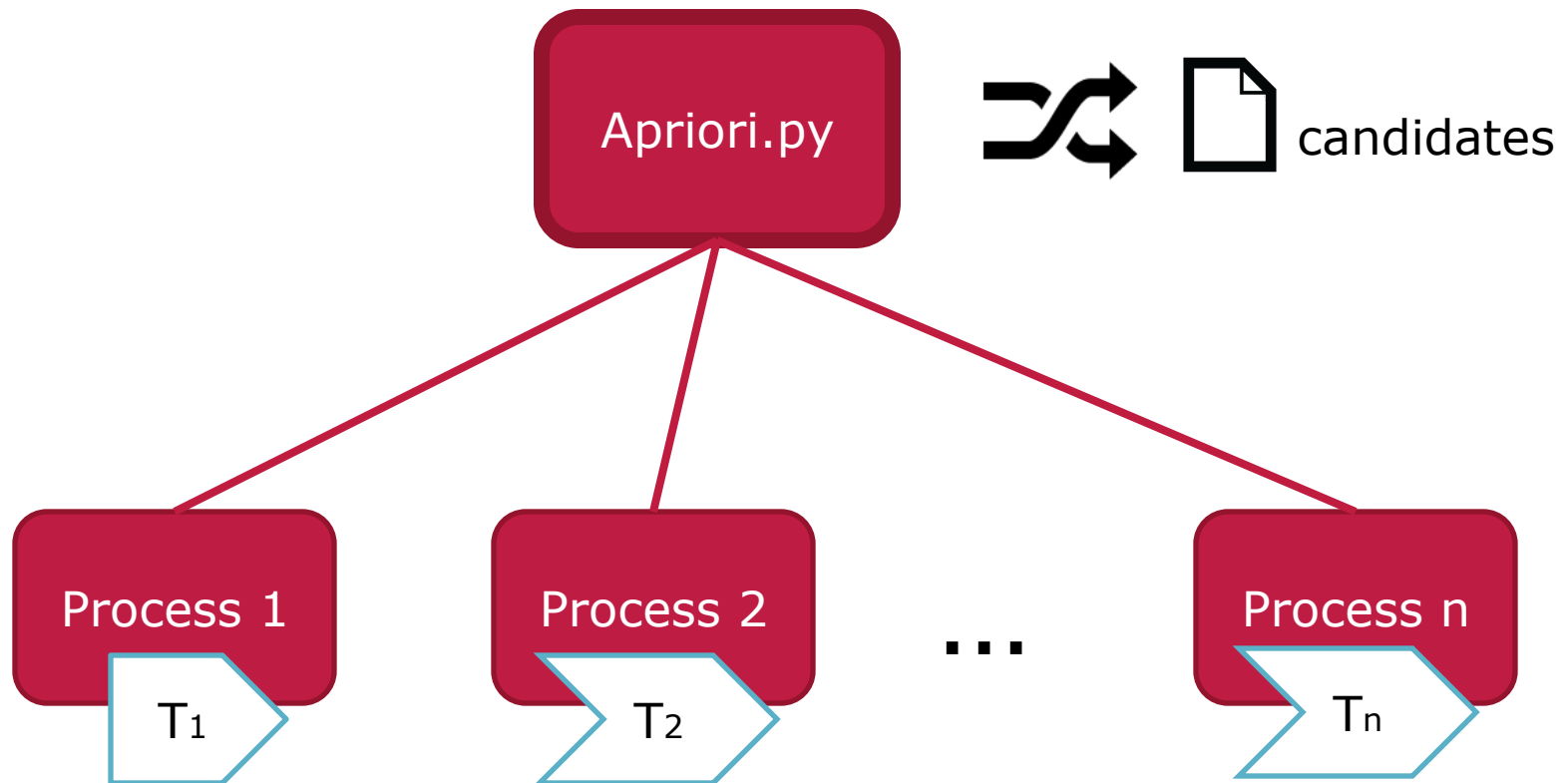
Parallelisierungsszenarien (2/3)

9



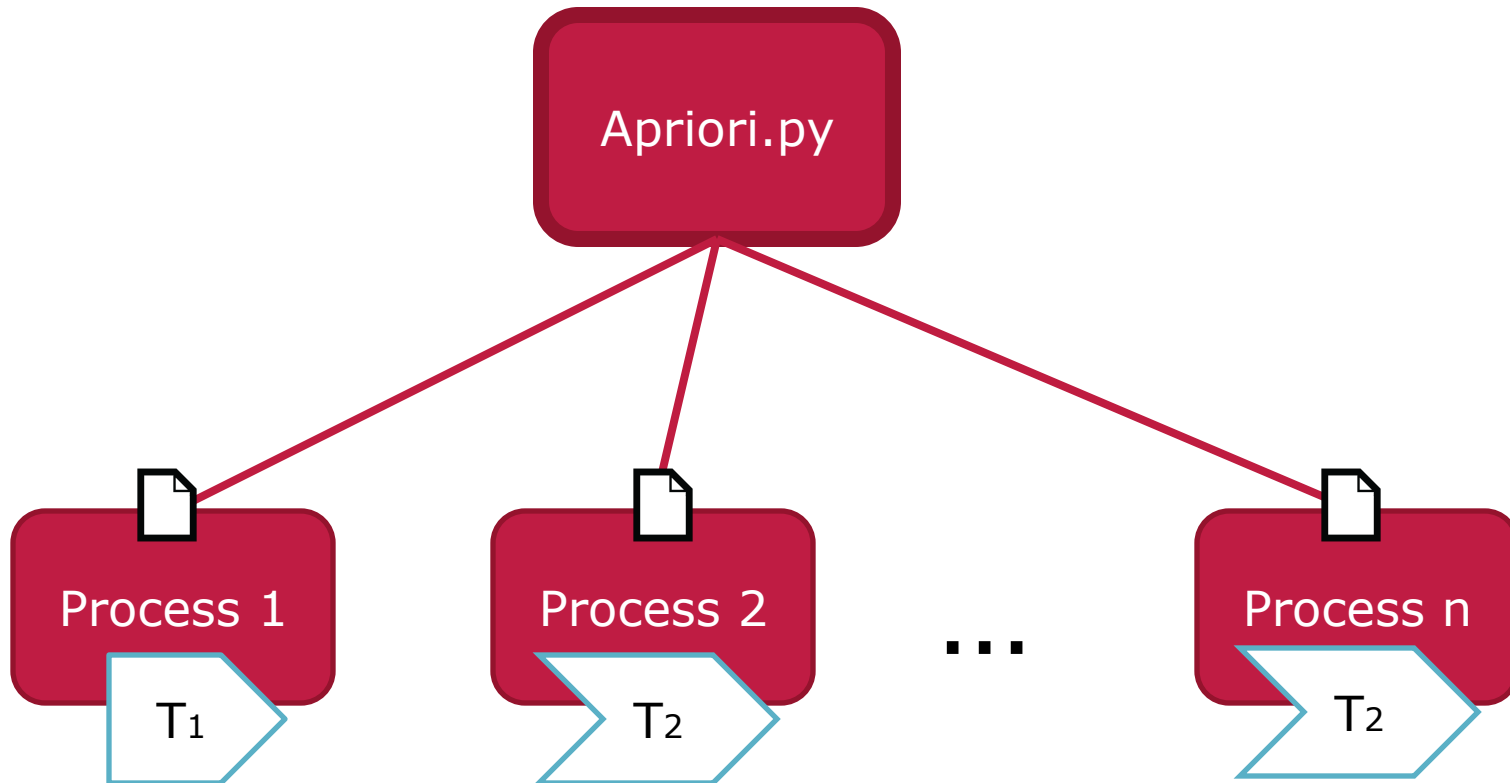
Parallelisierungsszenarien (2/3)

10



Parallelisierungsszenarien (2/3)

11



Analyse (2/3)

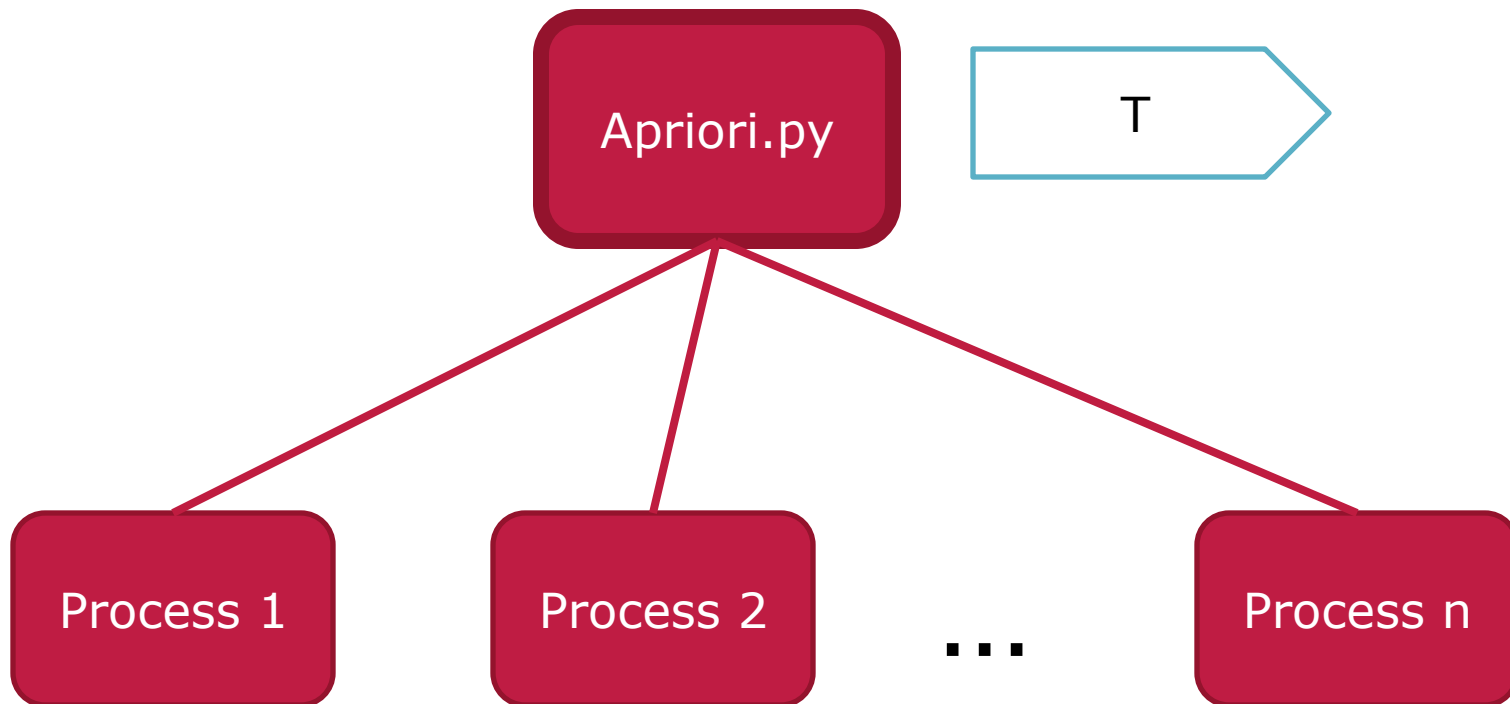
12

- Initialisierungsaufwand: $|D|$ 5Mio
- Aufwand in Schritt k: $n * |C_k|$ 4 * 400
- Ergebnisgröße: $n * |C_k|$ 4 * 400

- Gesamtkosten: $|D| + \sum_k 2n |C_k|$ 5Mio

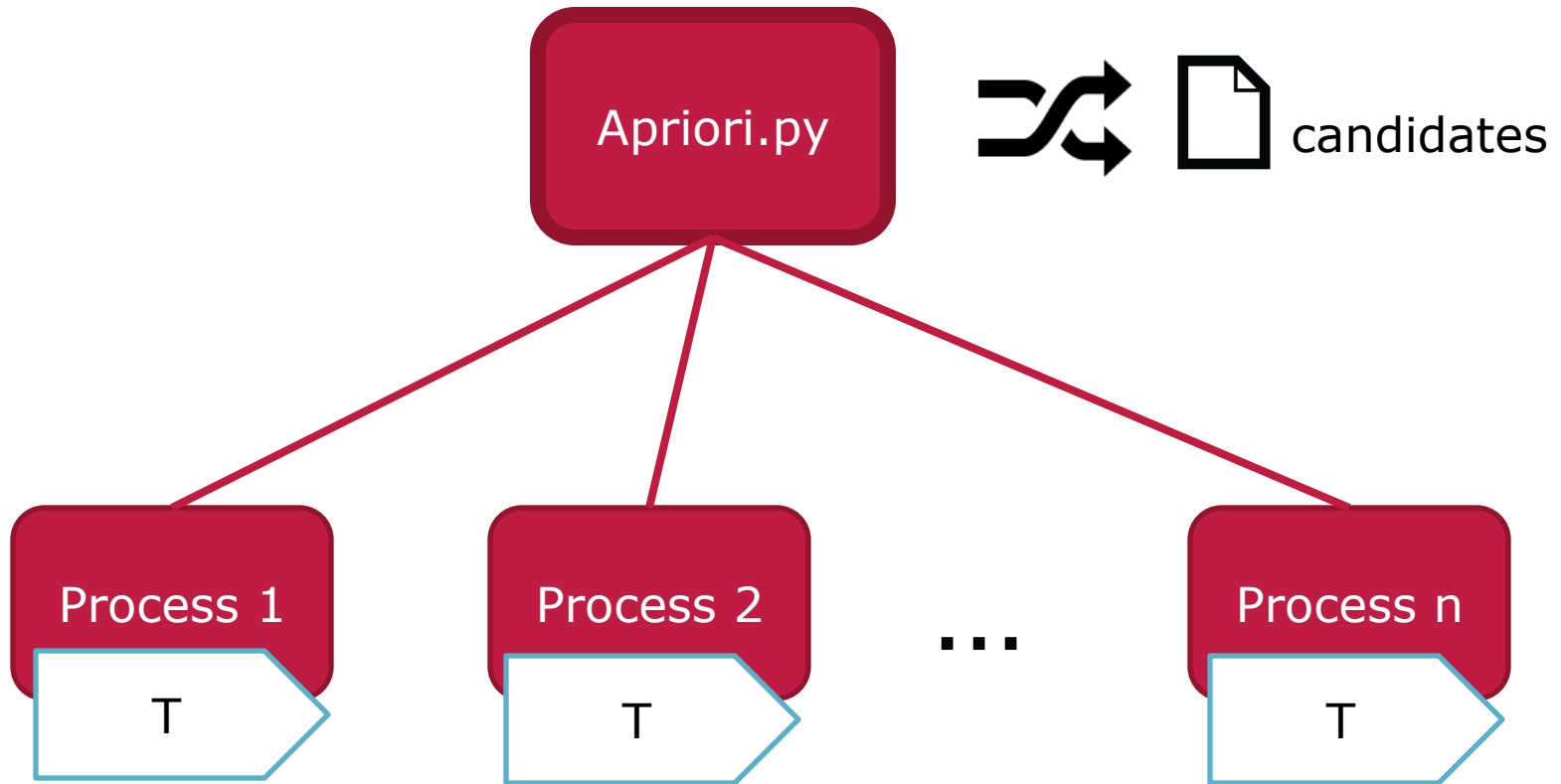
Parallelisierungsszenarien (3/3)

13



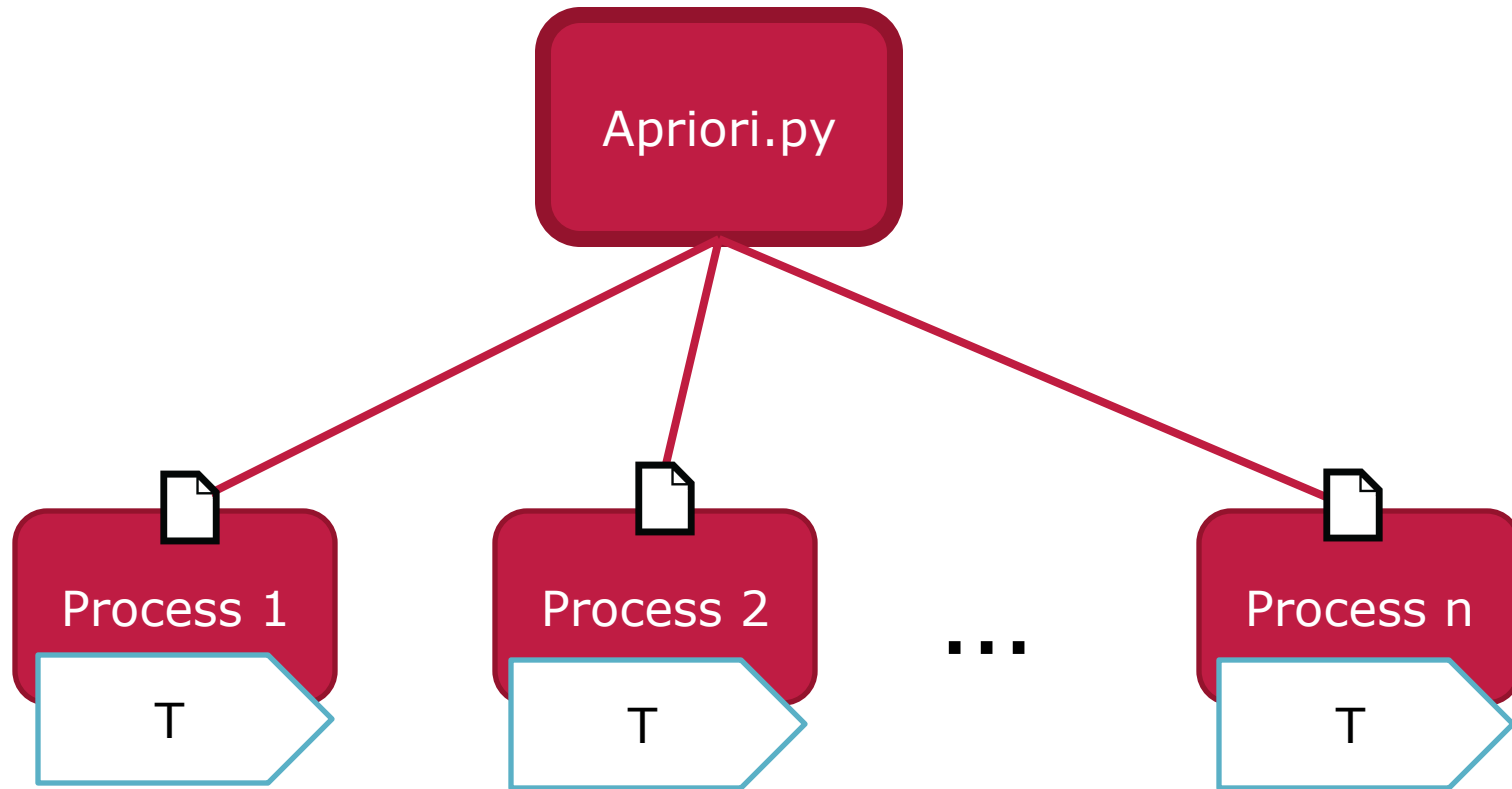
Parallelisierungsszenarien (3/3)

14



Parallelisierungsszenarien (3/3)

15



Analyse (3/3)

16

- Initialisierungsaufwand: $n * |D|$

4 * 5Mio

- Aufwand in Schritt k: $|C_k|$

400

- Ergebnisgröße: $|C_k|$

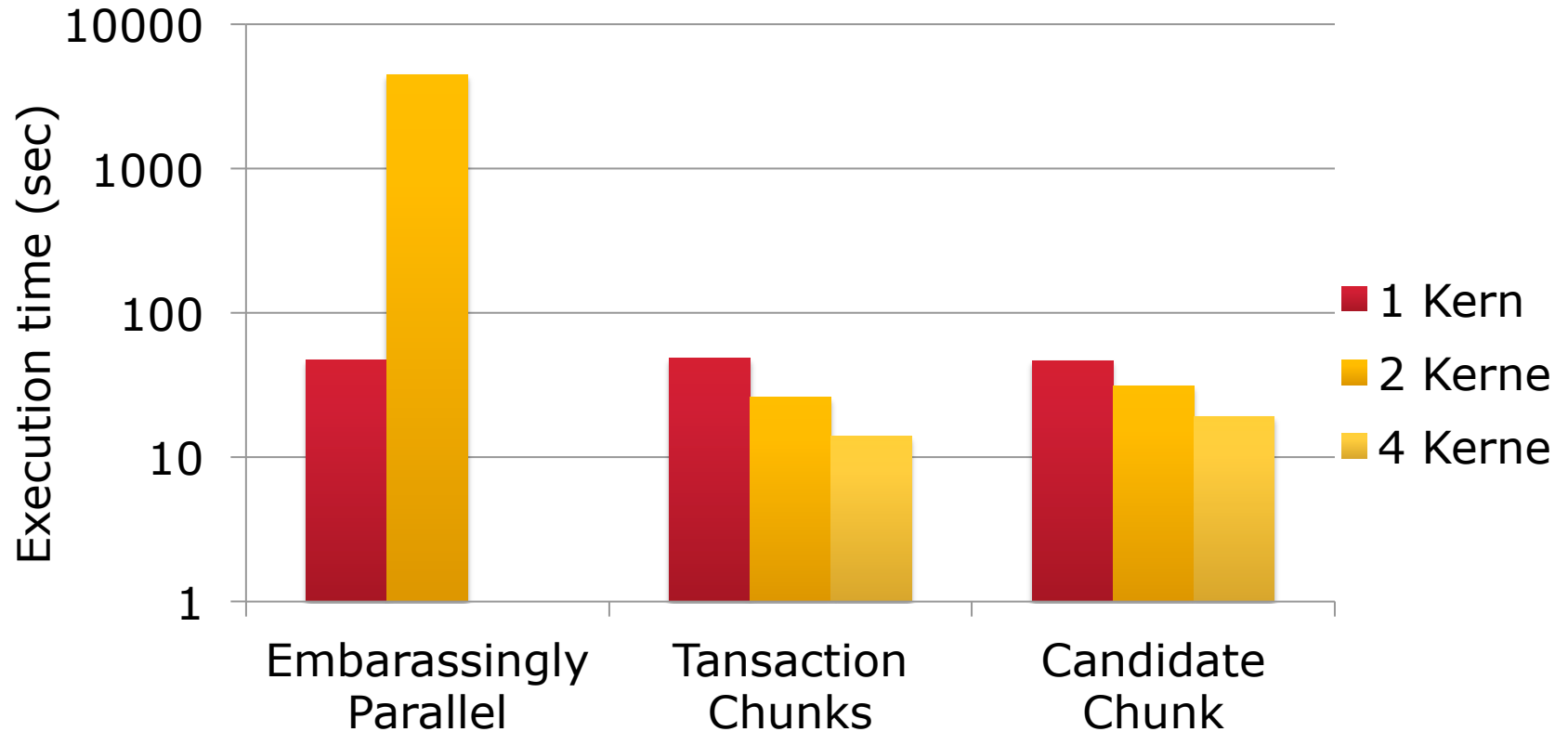
400

- Gesamtkosten: $n|D| + \sum_k 2|C_k|$

20Mio

Theorie vs. Praxis

17



$$\sum_k (|D| + 1) * |C_k|$$

$$|D| + \sum_k 2n |C_k|$$

$$n |D| + \sum_k 2 |C_k|$$

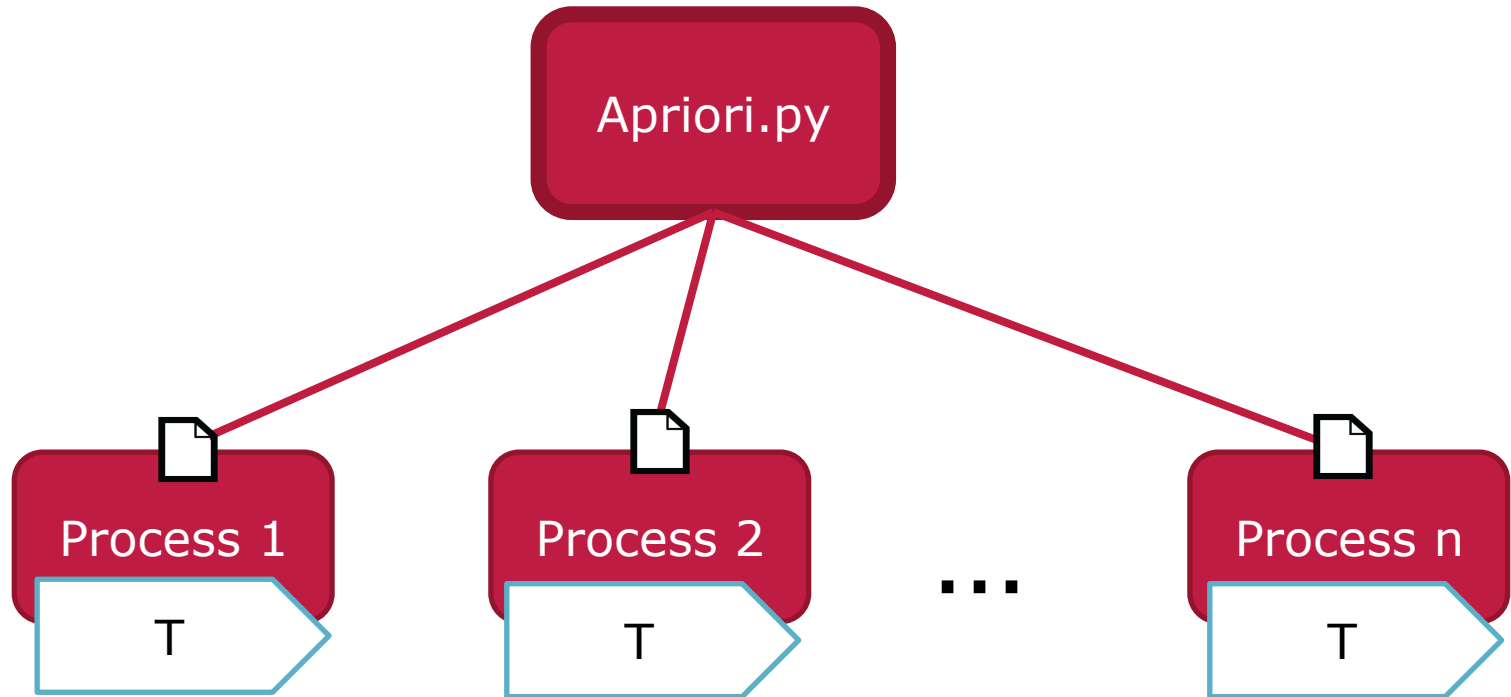
Implementierung (1)

18

- Map/Reduce
 - Map(candidates, transactions) → list(candidate, count)
 - Reduce(candidate, list (counts)) → list(candidates)

Implementierung (2)

19



- Update der Transaktionen
 - Würde Rückübertragung aller Transaktionen benötigen
 - Daher nicht mehr AprioriTID sondern Apriori

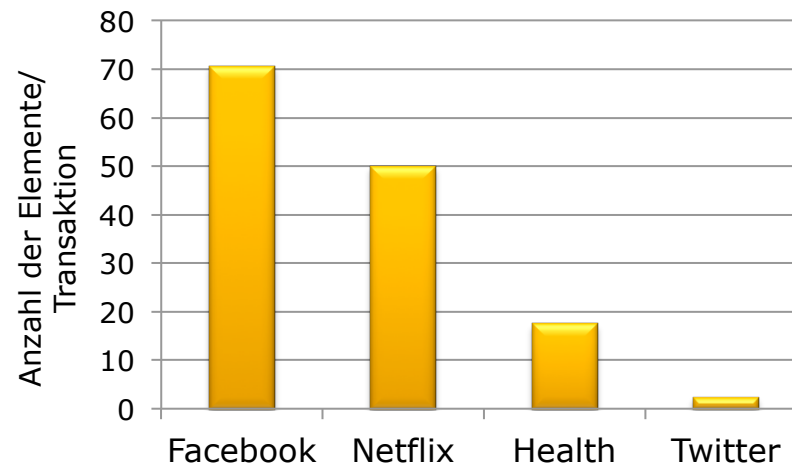
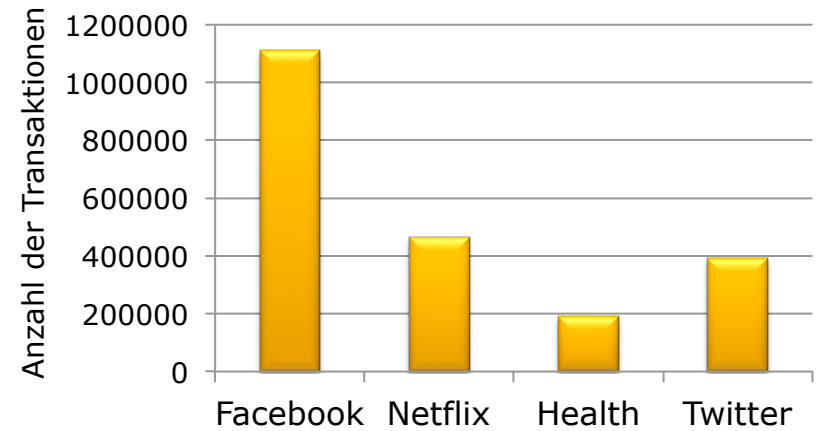
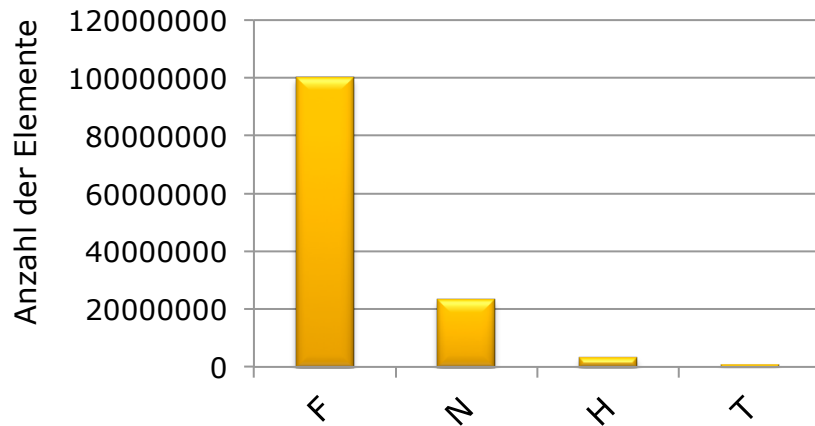
Fazit & Ideas

20

- Parallelisierung in anderen Programmiersprachen könnte mit Threads implementiert werden
 - Weniger Kommunikationsoverhead
 - Geringere Context Switch Kosten
 - Aber ggf. Wartezeiten wegen locking
- Parallelisierung auf Hadoop
- Hoher Data Parallelism
 - CUDA/OpenCL

Datensets

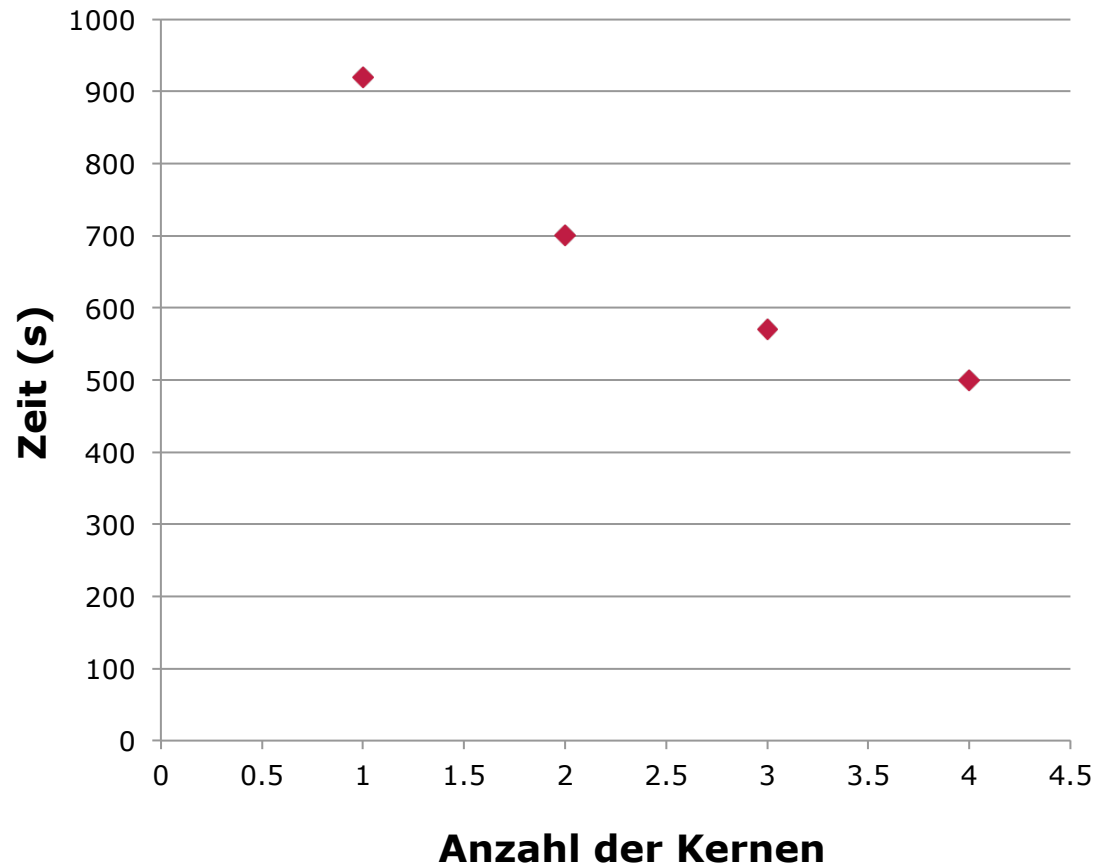
21



Benchmark (Netflix)

22

- $|D| = 200.000$
- $AVG(|T|) = 52,6$
- $N = 10.232.001$
- 40% speedup



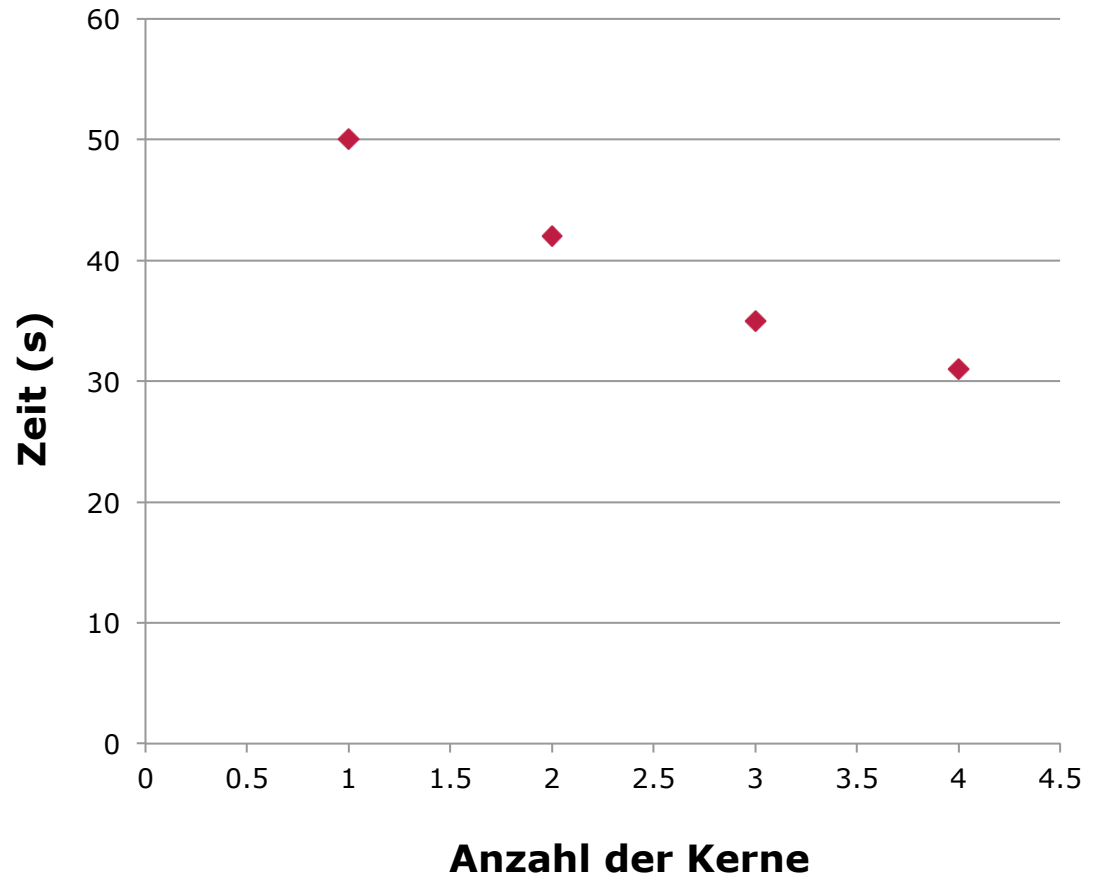
Hardware:

- Core i7-2600 Quad-Core
- 16GB RAM

Benchmark (Twitter)

23

- $|D| = 200.000$
- $AVG(|T|) = 2,4$
- $N = 480.000$
- 40% speedup



Hardware:

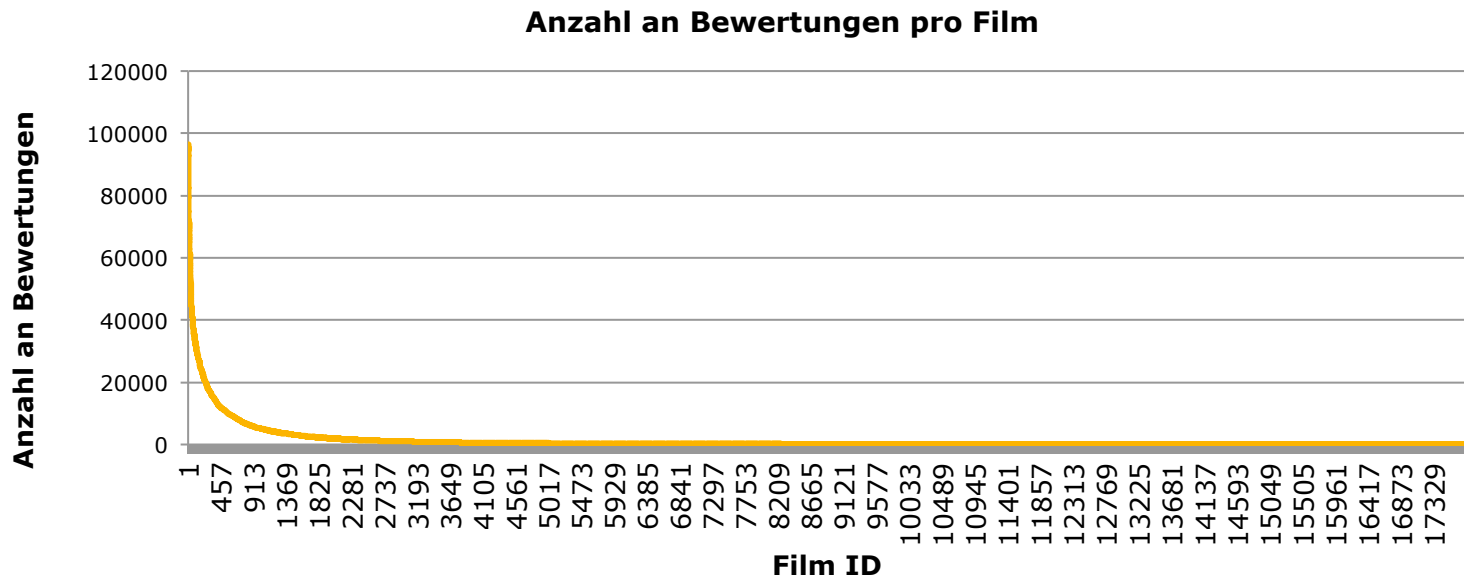
- Core i7-2600 Quad-Core
- 16GB RAM

Auffinden interessanter Regeln

24



Pirates of the Caribbean -> LoR II





$$\text{lift}(X \rightarrow Y) = \frac{\text{sup}(X \cup Y)}{\text{sup}(X) * \text{sup}(Y)}$$

$$\text{lift}(\text{MONSTERS, INC.} \rightarrow \text{NEMO}) = \frac{5.78\%}{7.2\% * 18\%} = 4.46$$

Errechnete Werte: 3,7 – 7,5

Conviction

26

$$\text{conv}(X \rightarrow Y) = \frac{1 - \text{sup}(Y)}{1 - \text{conf}(X \rightarrow Y)}$$

$$\text{conv}(\text{Monsters, Inc.} \rightarrow \text{Nemo}) = \frac{1 - 0.18}{1 - 0.77} = 3.55$$


Errechnete Werte: 3,5 - 25.32

Ergebnisse

27

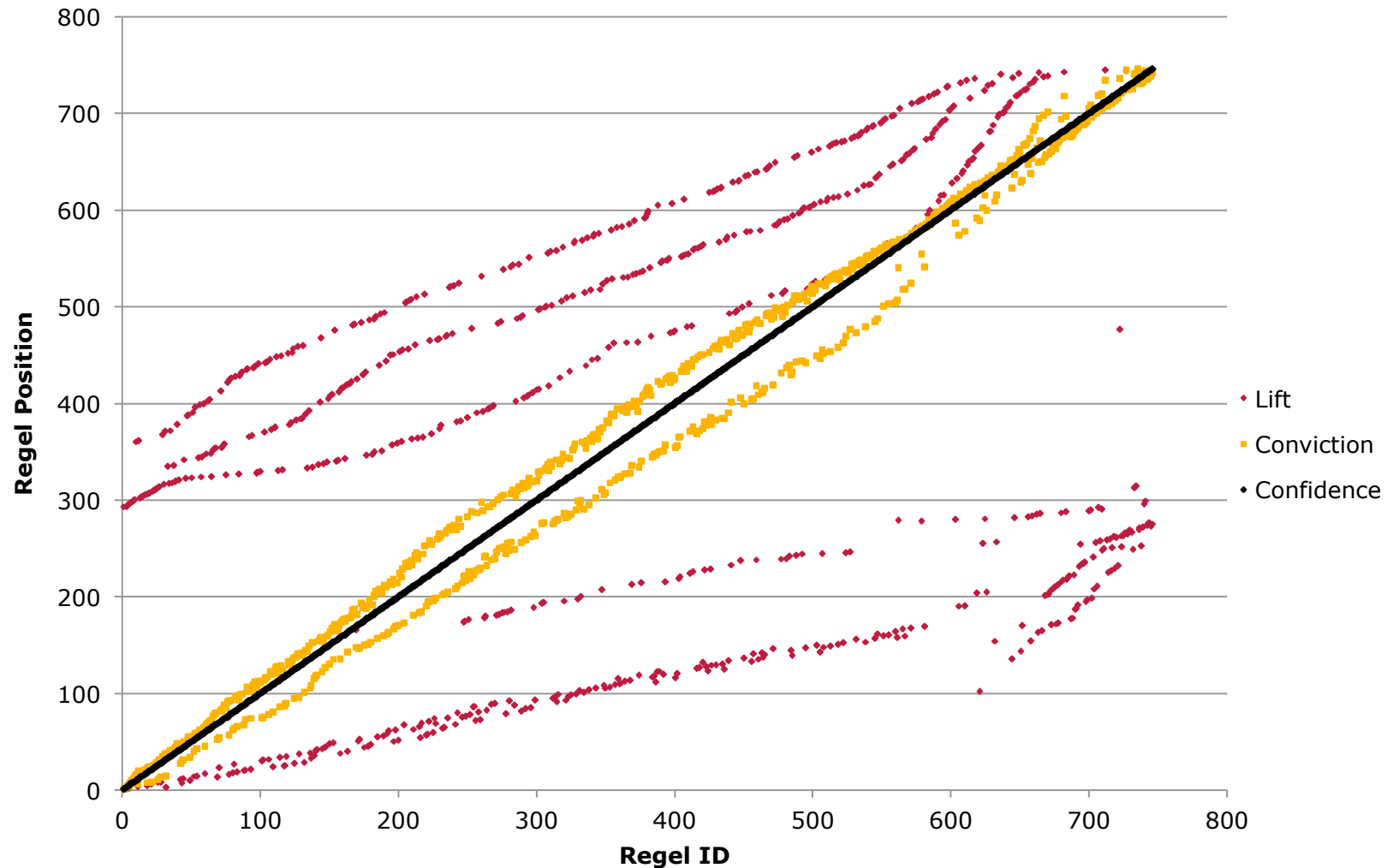


Support: 5,3 %
 Confidence: 96,5 %
 Lift: 7,5
 Conviction: 25,3

Support: 5%
 Confidence: 75%
 746 Regeln

Einfluss von Lift und Conviction

28



Vergleich zu IMDb

29

People who liked this also liked... [Learn more](#)



The Lord of the Rings: The Return of the King (2003)
 PG-13 Action | Adventure | Drama ...
 ★★★★★★ 8.9/10

Aragorn leads the World of Men against Sauron's army to draw the dark lord's gaze from Frodo and Sam who are on the doorstep of Mount Doom with the One Ring.

Director: Peter Jackson
Stars: Elijah Wood and Viggo Mort...

[Add to Watchlist](#)
[Next »](#)

◀ Prev 6 Next 6 ▶

Stichprobe (N=30):

21/30 = 70% stimmen mit IMDb überein

- Nachfolger
- Gleiches Genre

Quellen

30

- [Sergey Brin, Rajeev Motwani, Jeffrey D. Ullman, and Shalom Tsur. Dynamic itemset counting and implication rules for market basket data.](#)
- [P. Becuzzi, M. Coppola and M. Vanneschi. Mining of Association Rules in Very Large Databases: A Structured Parallel Approach*](#)
- Ali Tarhini: Parallel Apriori algorithm for frequent pattern mining. <http://alitarhini.wordpress.com/2011/02/26/parallel-apriori-algorithm-for-frequent-pattern-mining/> [Stand 3. Juli 2012]
- Anuradha.T, Satya Pasad R, S.N.Tirumalarao. Parallelizing Apriori on Dual Core using OpenMP. International Journal of Computer Applications, Volume 43
- R. Agrawal, R. Srikant: "Fast Algorithms for Mining Association Rules", Proc. of the 20th Int'l Conference on Very Large Databases, Santiago, Chile, Sept. 1994