

tf/idf computation

Florian Thomas, Christian Reß
Map/Reduce Algorithms on Hadoop
6. Juli 2009



tf/idf computation

- **Was ist tf/idf?**
- **Verschiedene Implementierungen**
- **Map/Reduce-Aufbau**
- **Implementierungsbesonderheiten**
- **Analyse**
- **Fazit**

Was ist tf/idf?

term frequency

$$\mathbf{tf}_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

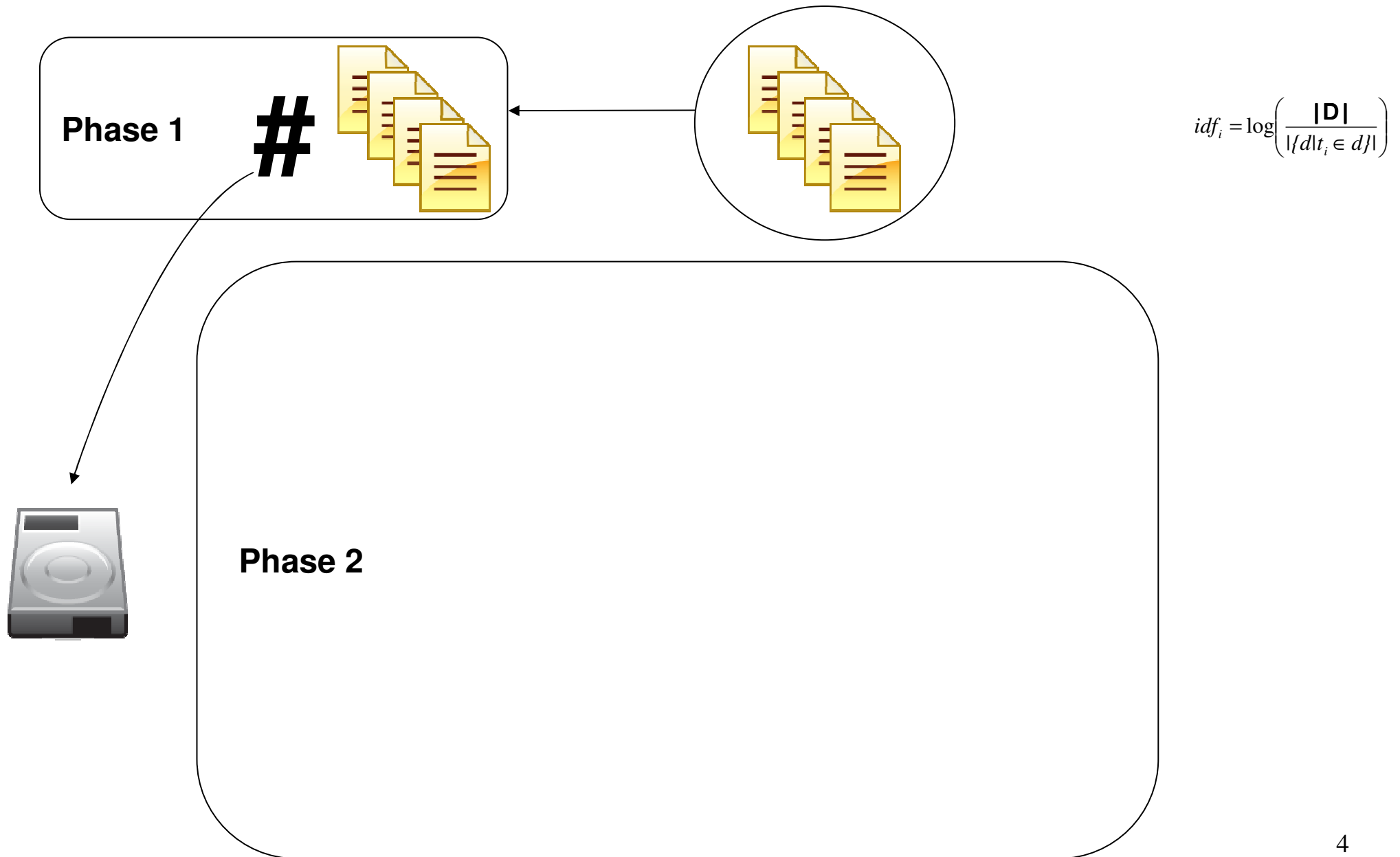
inverse document frequency

$$\mathbf{idf}_i = \log \left(\frac{|D|}{|\{d | t_i \in d\}|} \right)$$

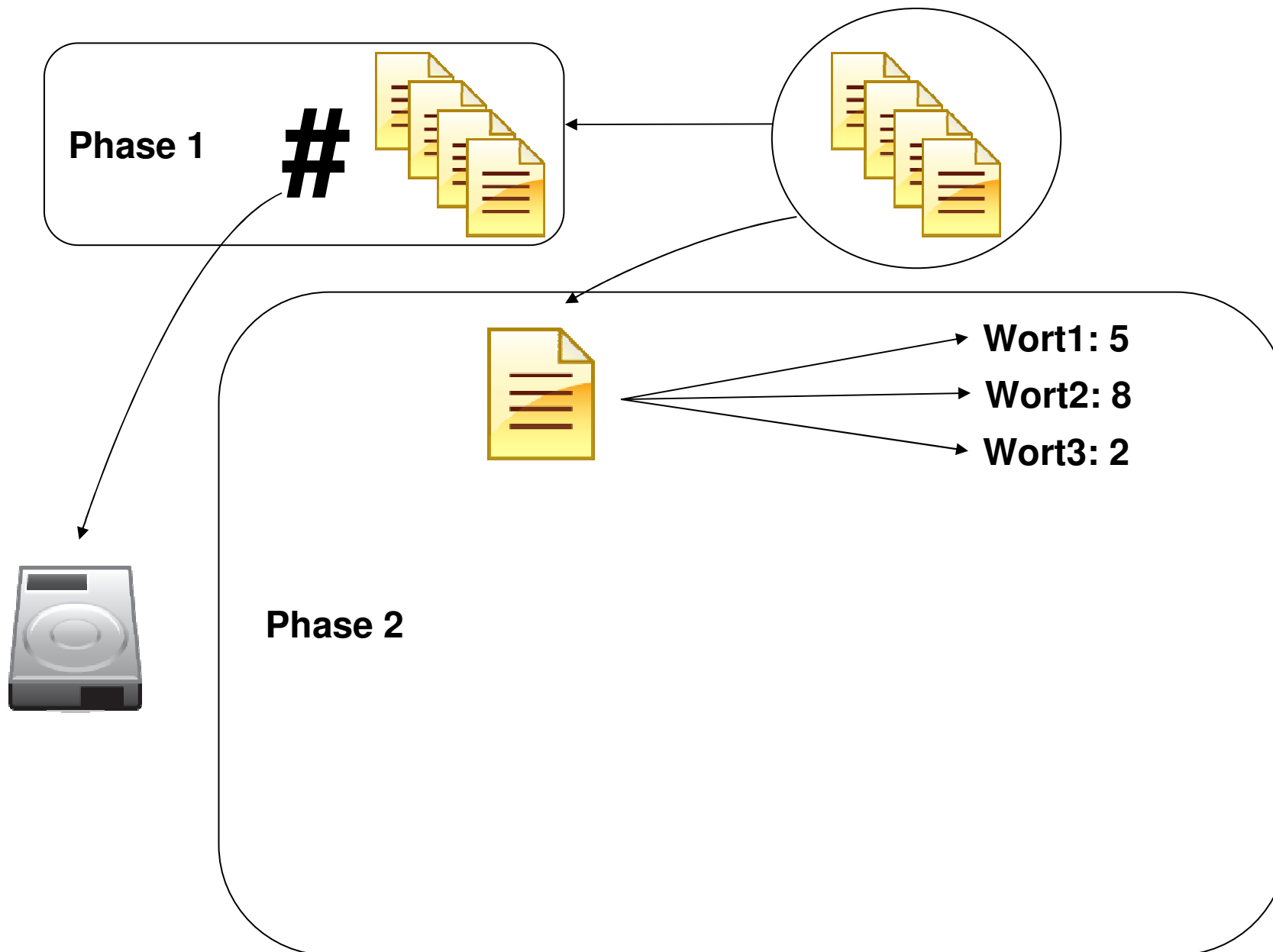
tf/idf

$$(\mathbf{tf} - \mathbf{idf})_{i,j} = tf_{i,j} \cdot idf_i$$

Implementierungen – 2 Phasen



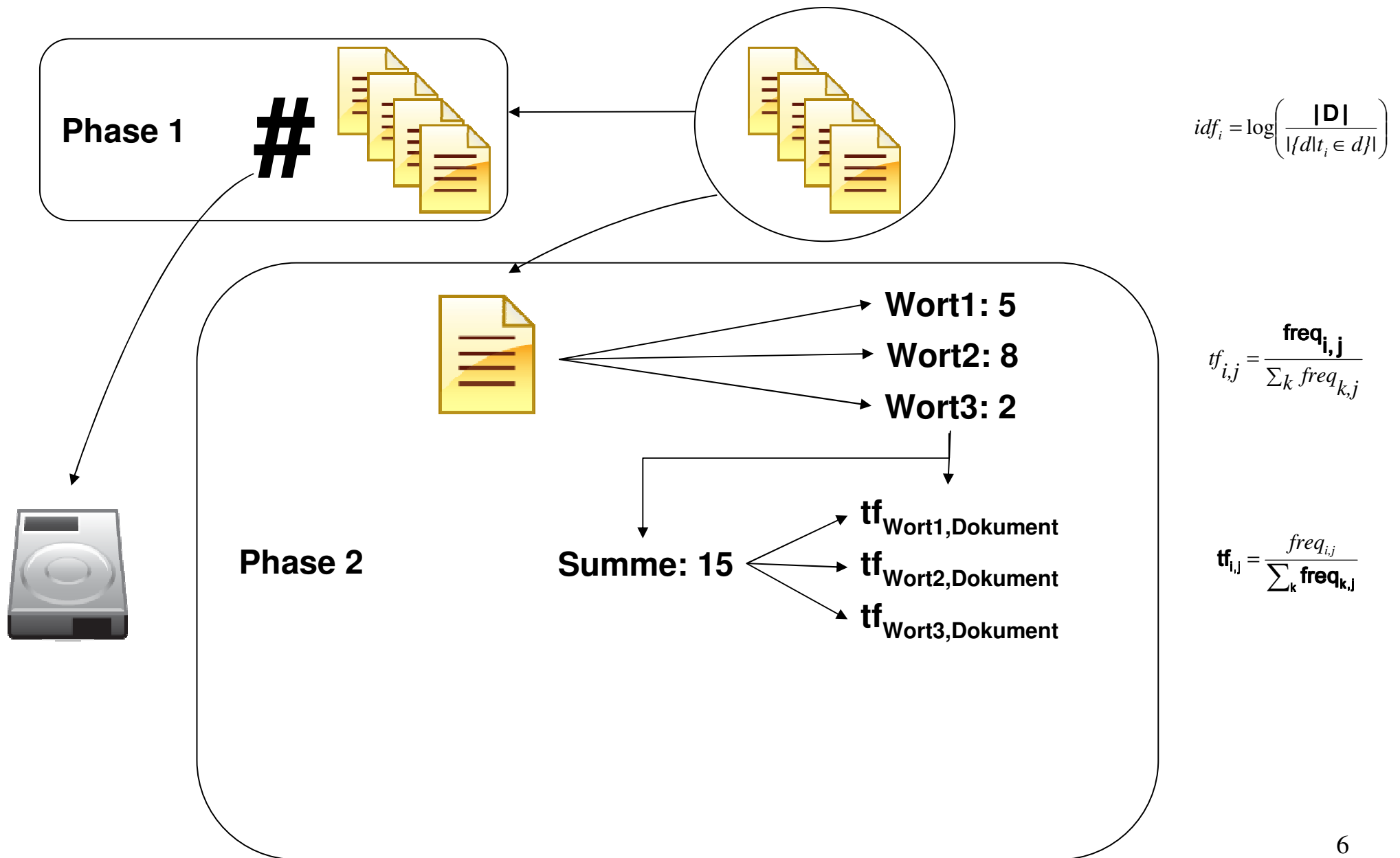
Implementierungen – 2 Phasen



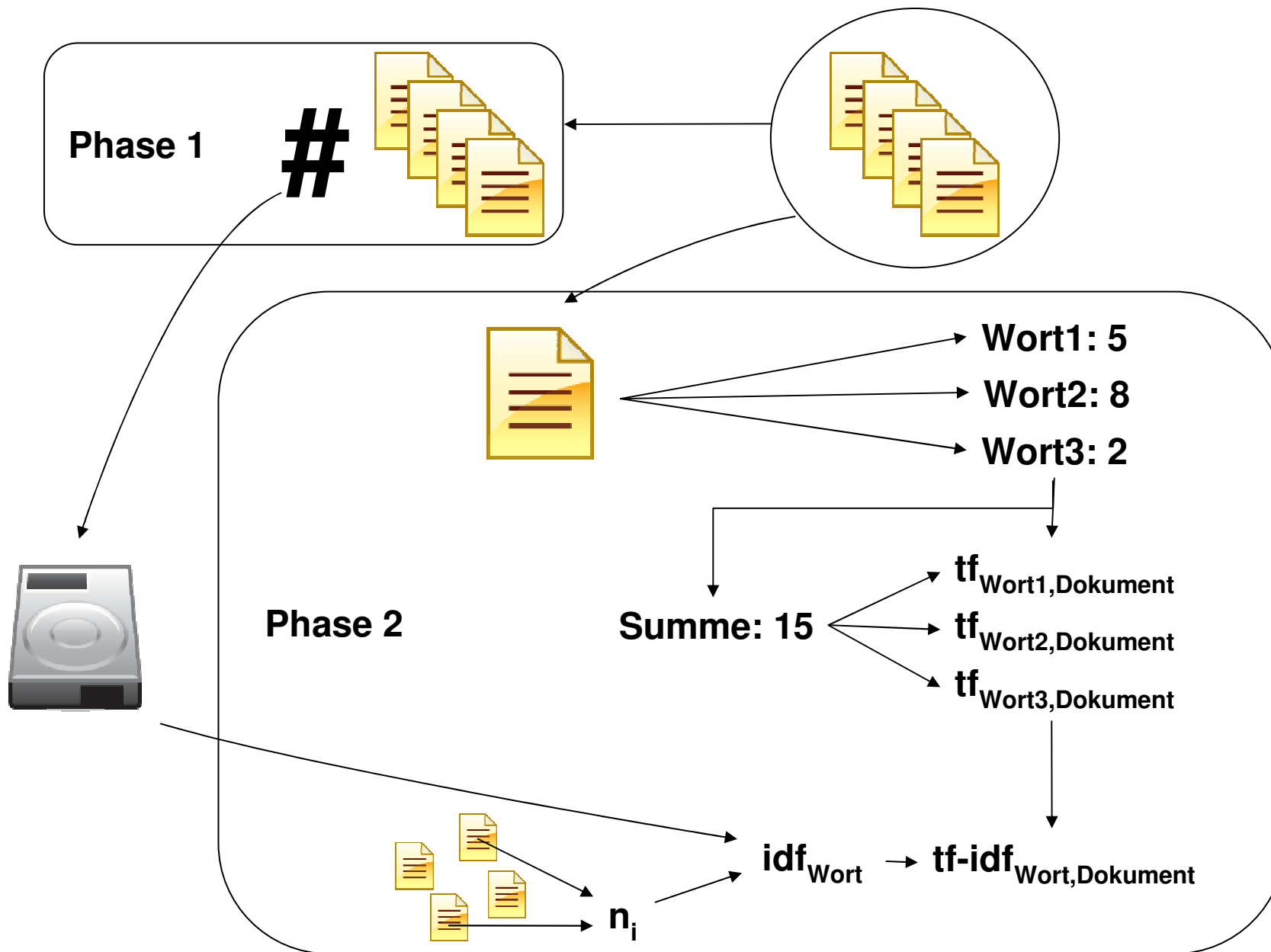
$$idf_i = \log \left(\frac{|D|}{|\{d | t_i \in d\}|} \right)$$

$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

Implementierungen – 2 Phasen



Implementierungen – 2 Phasen



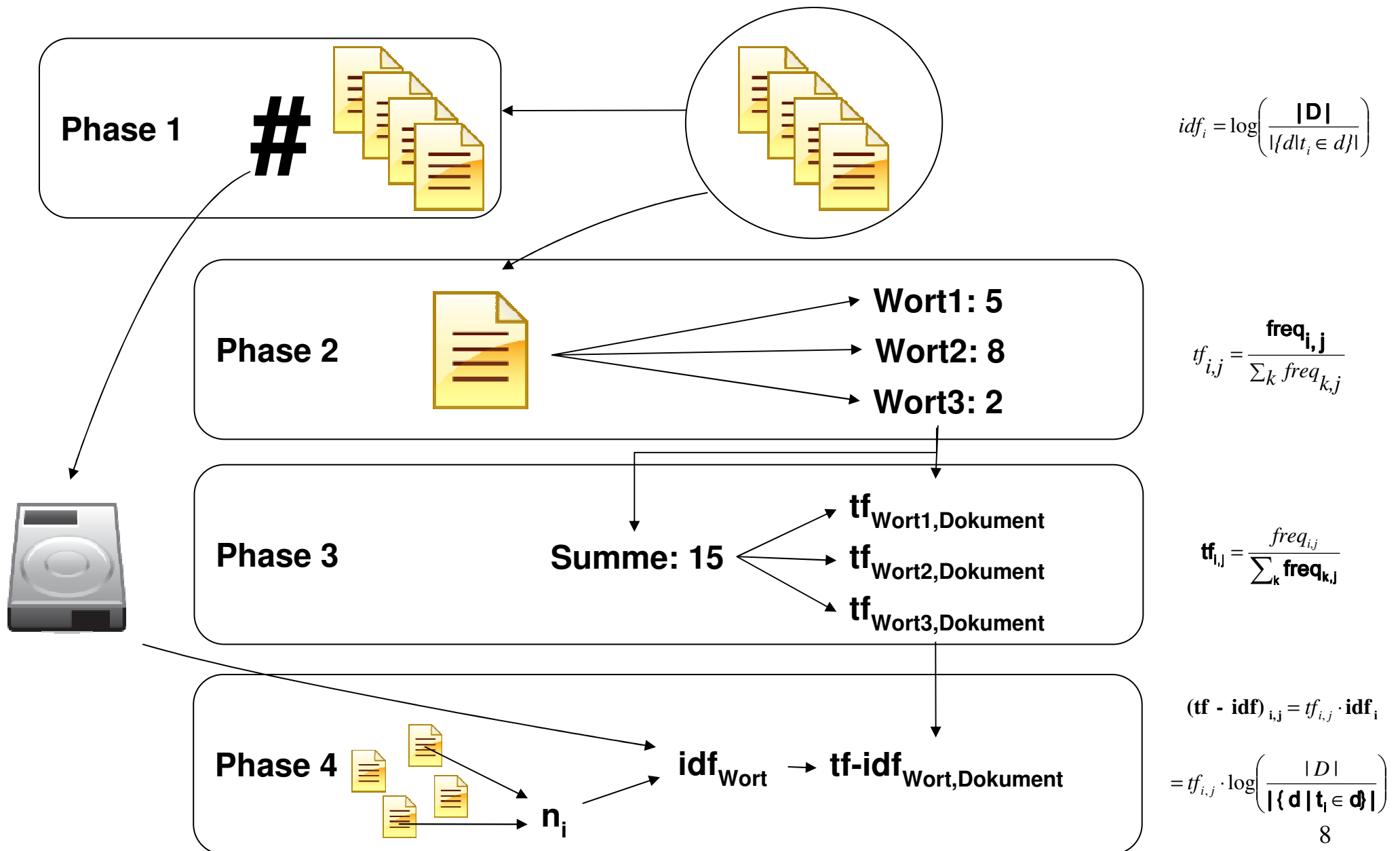
$$idf_i = \log \left(\frac{|D|}{|\{d | t_i \in d\}|} \right)$$

$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

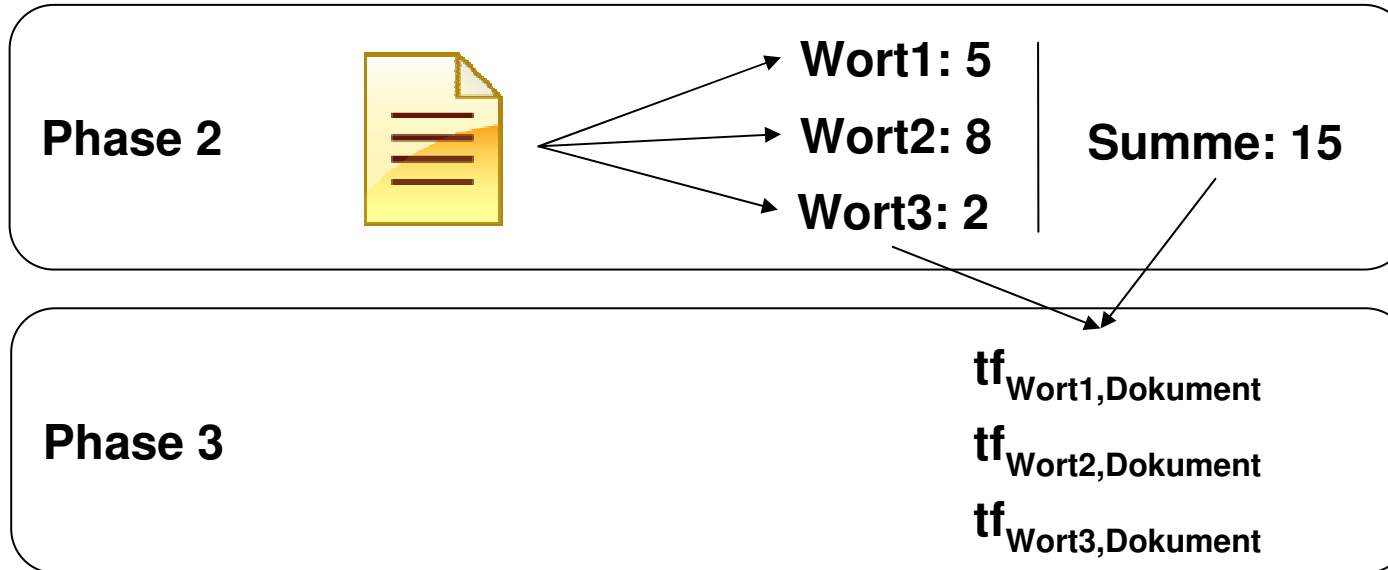
$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

$$\begin{aligned} (tf-idf)_{i,j} &= tf_{i,j} \cdot idf_i \\ &= tf_{i,j} \cdot \log \left(\frac{|D|}{|\{d | t_i \in d\}|} \right) \end{aligned}$$

Implementierungen – 4 Phasen



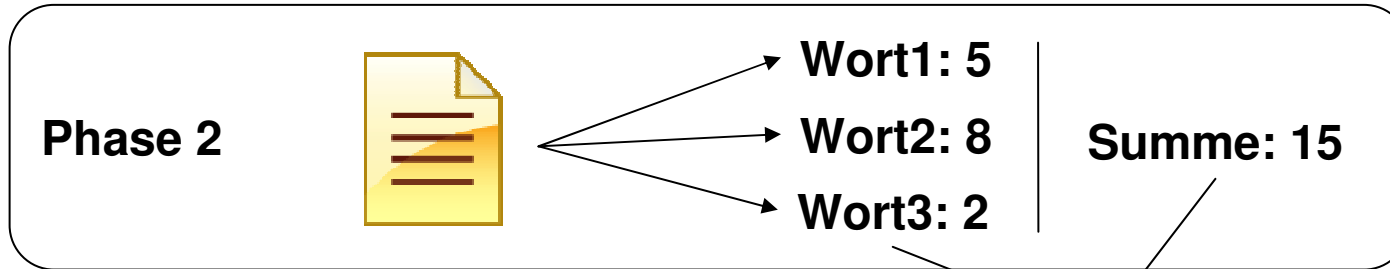
Implementierungen – 5 Phasen



$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

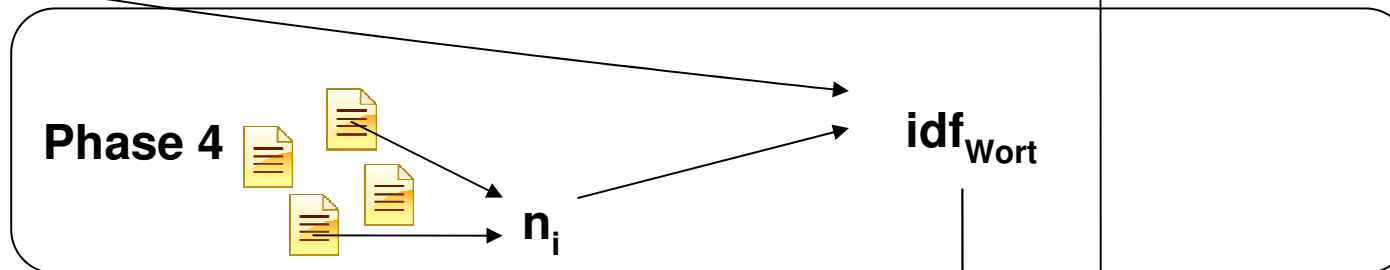
Implementierungen – 5 Phasen



$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$



$$tf_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$



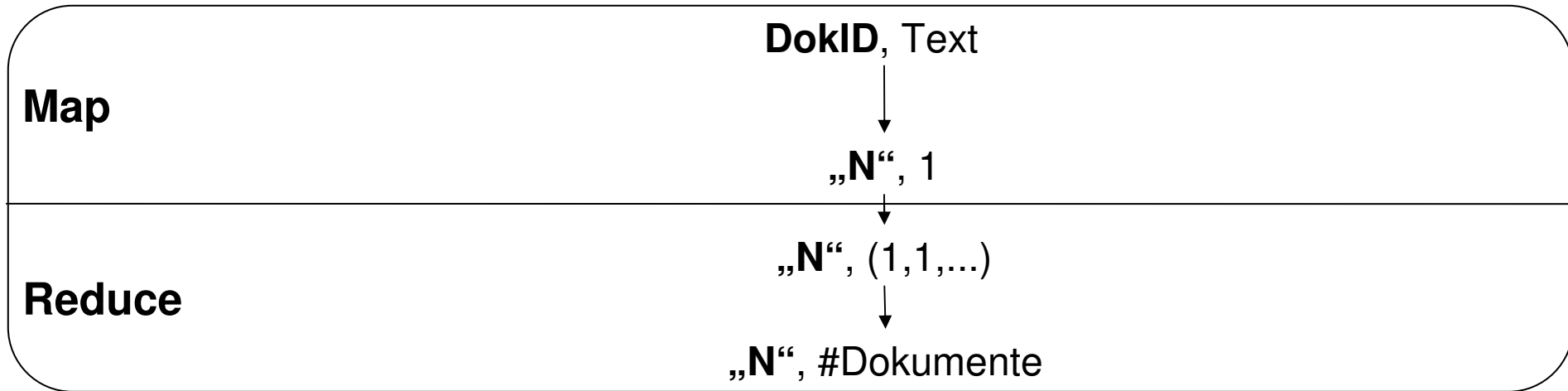
$$idf_i = \log \left(\frac{|D|}{|\{d | t_i \in d\}|} \right)$$



$$(tf-idf)_{i,j} = tf_{i,j} \cdot idf_i$$

Map/Reduce-Aufbau – 5 Phasen

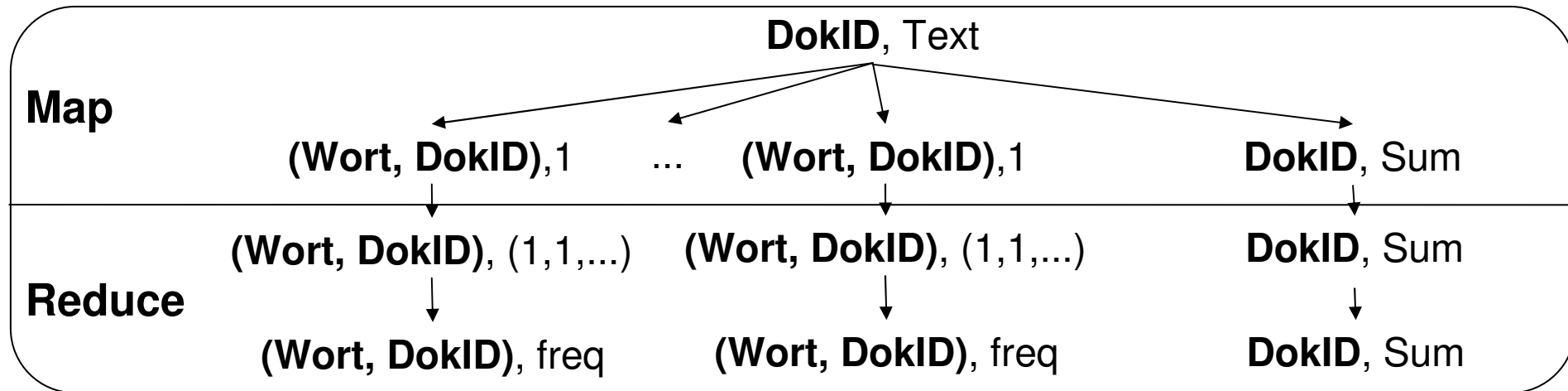
Phase 1: Dokumente zählen



$$idf_i = \log \left(\frac{|D|}{|\{d | t_i \in d\}|} \right)$$

Map/Reduce-Aufbau – 5 Phasen

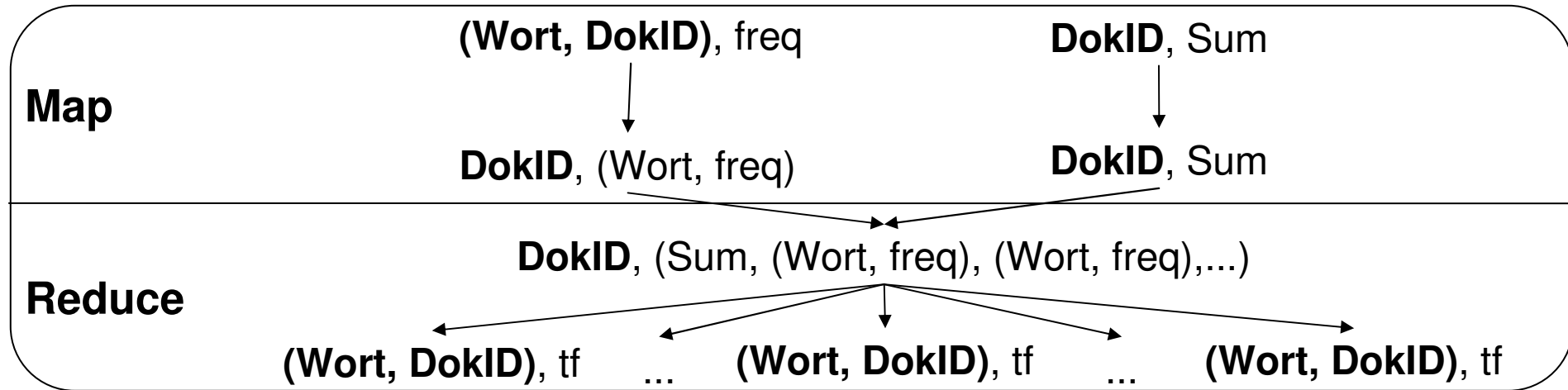
Phase 2: Wörter zählen in jedem Dokument



$$tf_{i,j} = \frac{\text{freq}_{i,j}}{\sum_k \text{freq}_{k,j}}$$

Map/Reduce-Aufbau – 5 Phasen

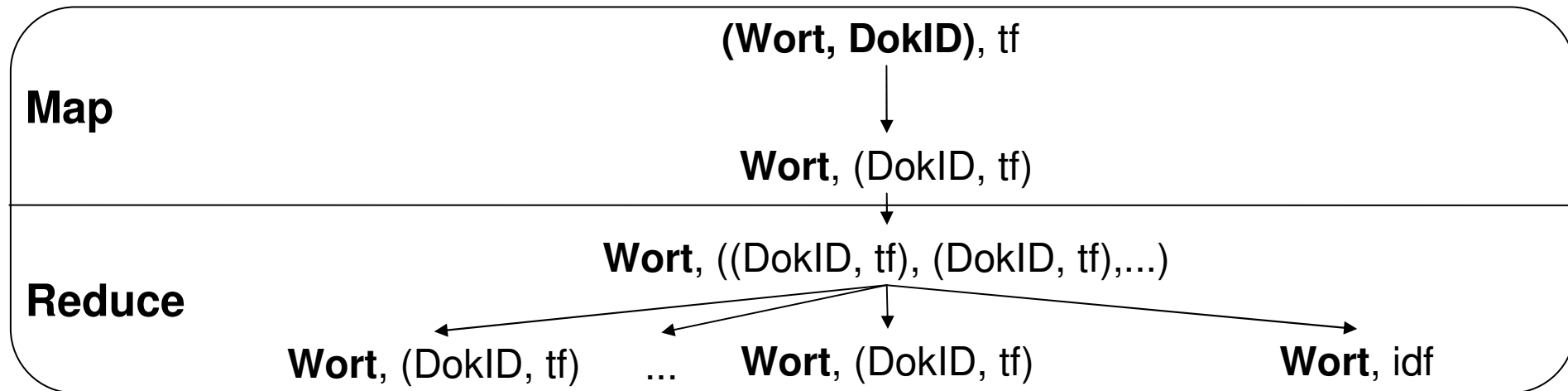
Phase 3: tf berechnen



$$\mathbf{tf}_{i,j} = \frac{freq_{i,j}}{\sum_k freq_{k,j}}$$

Map/Reduce-Aufbau – 5 Phasen

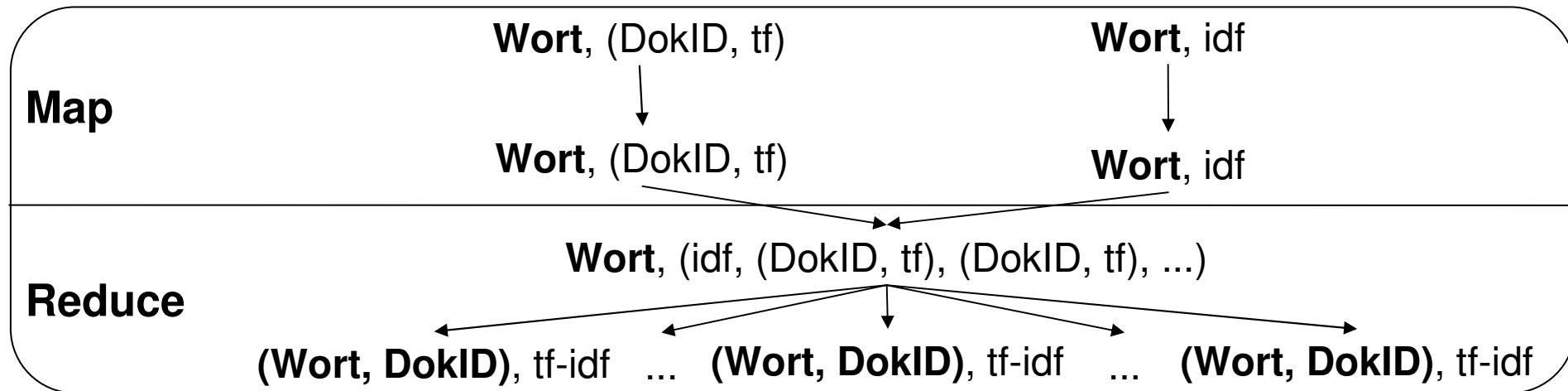
Phase 4: #Dokumente_{Wort} bestimmen, idf berechnen



$$\text{idf}_i = \log \left(\frac{|D|}{|\{d \mid t_i \in d\}|} \right)$$

Map/Reduce-Aufbau – 5 Phasen

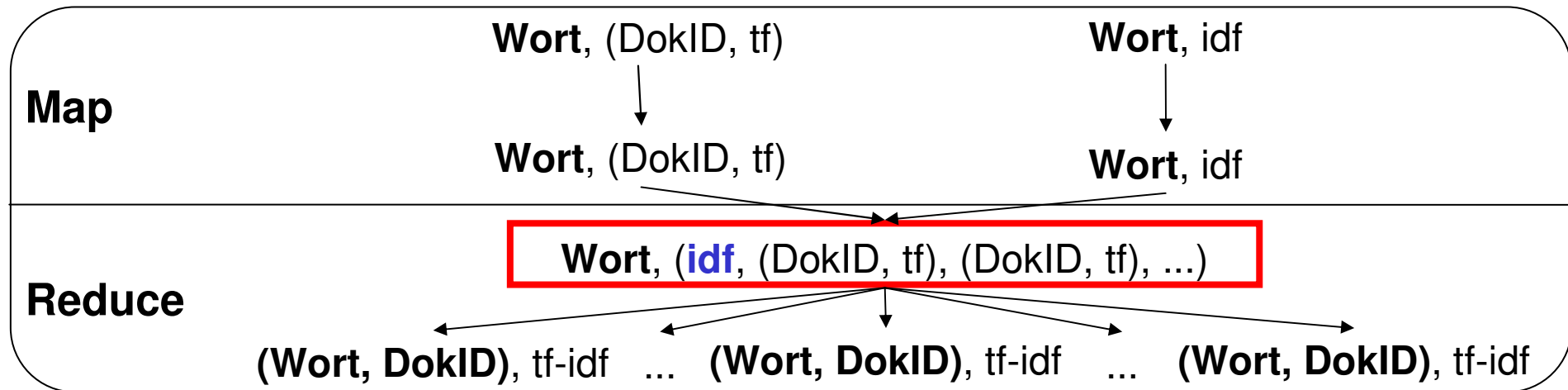
Phase 5: tf-idf berechnen



$$(\mathbf{tf - idf})_{i,j} = tf_{i,j} \cdot idf_i$$

Implementierung – Secondary Sort

Motivation: Phase 5



Normalerweise: **Wort**, ((DokID, tf), (DokID, tf), ..., **idf**, ...)

Implementierung – Secondary Sort

Hadoop – normales Verhalten

Wort	(DokID, tf)
mental	(1, 0.4)
mental	(3, 0.9)
mental	(5, 0.3)
mental	(0, 0.7)
mental	1.8
mental	(2, 0.6)
retardation	(8, 0.5)
retardation	5.1
retardation	(6, 0.1)
retardation	(9, 0.4)
retardation	(4, 0.8)

Implementierung – Secondary Sort

Hadoop – normales Verhalten

Wort	(DokID, tf)	
mental	(1, 0.4)	
mental	(3, 0.9)	
mental	(5, 0.3)	
mental	(0, 0.7)	→ mental , ((1, 0.4), (3, 0.9), (5, 0.3), (0, 0.7), 1.8 , (2, 0.6))
mental	1.8	
mental	(2, 0.6)	
retardation	(8, 0.5)	
retardation	5.1	
retardation	(6, 0.1)	→ retardation , ((8, 0.5), 5.1 , (6, 0.1), (9, 0.4), (4, 0.8))
retardation	(9, 0.4)	
retardation	(4, 0.8)	

Implementierung – Secondary Sort

Hadoop – gewünschtes Verhalten

Wort **(DokID, tf)**

mental 1 (1, 0.4)

mental 1 (3, 0.9)

mental 1 (5, 0.3)

mental 1 (0, 0.7)

mental 0 1.8

mental 1 (2, 0.6)

retardation 1 (8, 0.5)

retardation 0 5.1

retardation 1 (6, 0.1)

retardation 1 (9, 0.4)

retardation 1 (4, 0.8)

- **Anhängen eines Index-Wertes an den Schlüssel**

Implementierung – Secondary Sort

Hadoop – gewünschtes Verhalten

Wort (DokID, tf)

mental 0 1.8

mental 1 (1, 0.4)

mental 1 (3, 0.9)

mental 1 (5, 0.3)

mental 1 (0, 0.7)

mental 1 (2, 0.6)

retardation 0 5.1

retardation 1 (8, 0.5)

retardation 1 (6, 0.1)

retardation 1 (9, 0.4)

retardation 1 (4, 0.8)

- **Anhängen eines Index-Wertes an den Schlüssel**

Implementierung – Secondary Sort

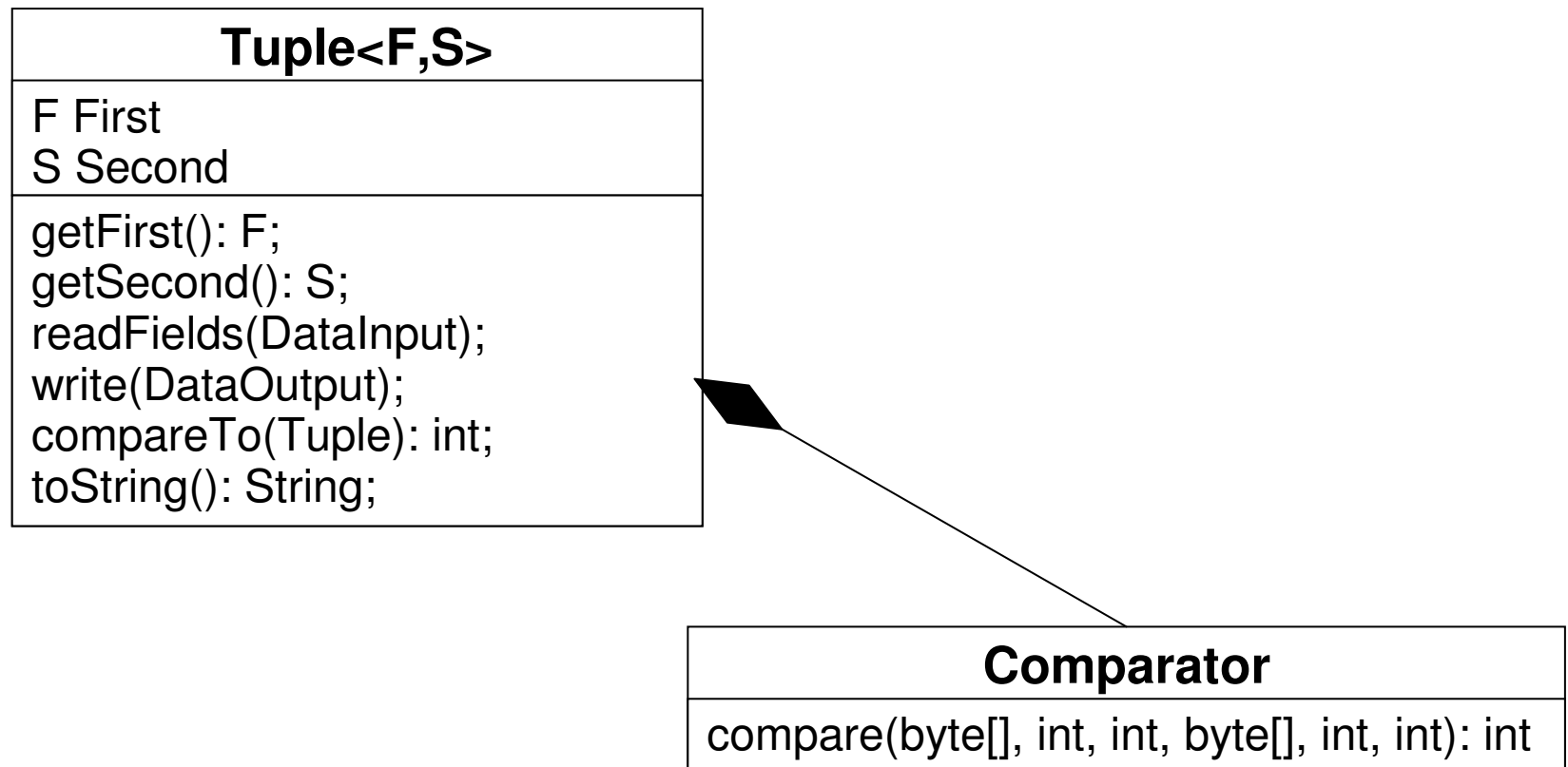
Hadoop – gewünschtes Verhalten

Wort	(DokID, tf)	
mental 0	1.8	
mental 1	(1, 0.4)	
mental 1	(3, 0.9)	→ mental , (1.8, (1, 0.4), (3, 0.9), (5, 0.3), (0, 0.7), (2, 0.6))
mental 1	(5, 0.3)	
mental 1	(0, 0.7)	
mental 1	(2, 0.6)	
<hr/>		
retardation 0 5.1		
retardation 1	(8, 0.5)	
retardation 1	(6, 0.1)	→ retardation , (5.1, (8, 0.5), (6, 0.1), (9, 0.4), (4, 0.8))
retardation 1	(9, 0.4)	
retardation 1	(4, 0.8)	
<hr/>		

- Anhängen eines Index-Wertes an den Schlüssel
- Anpassen der Gruppierung

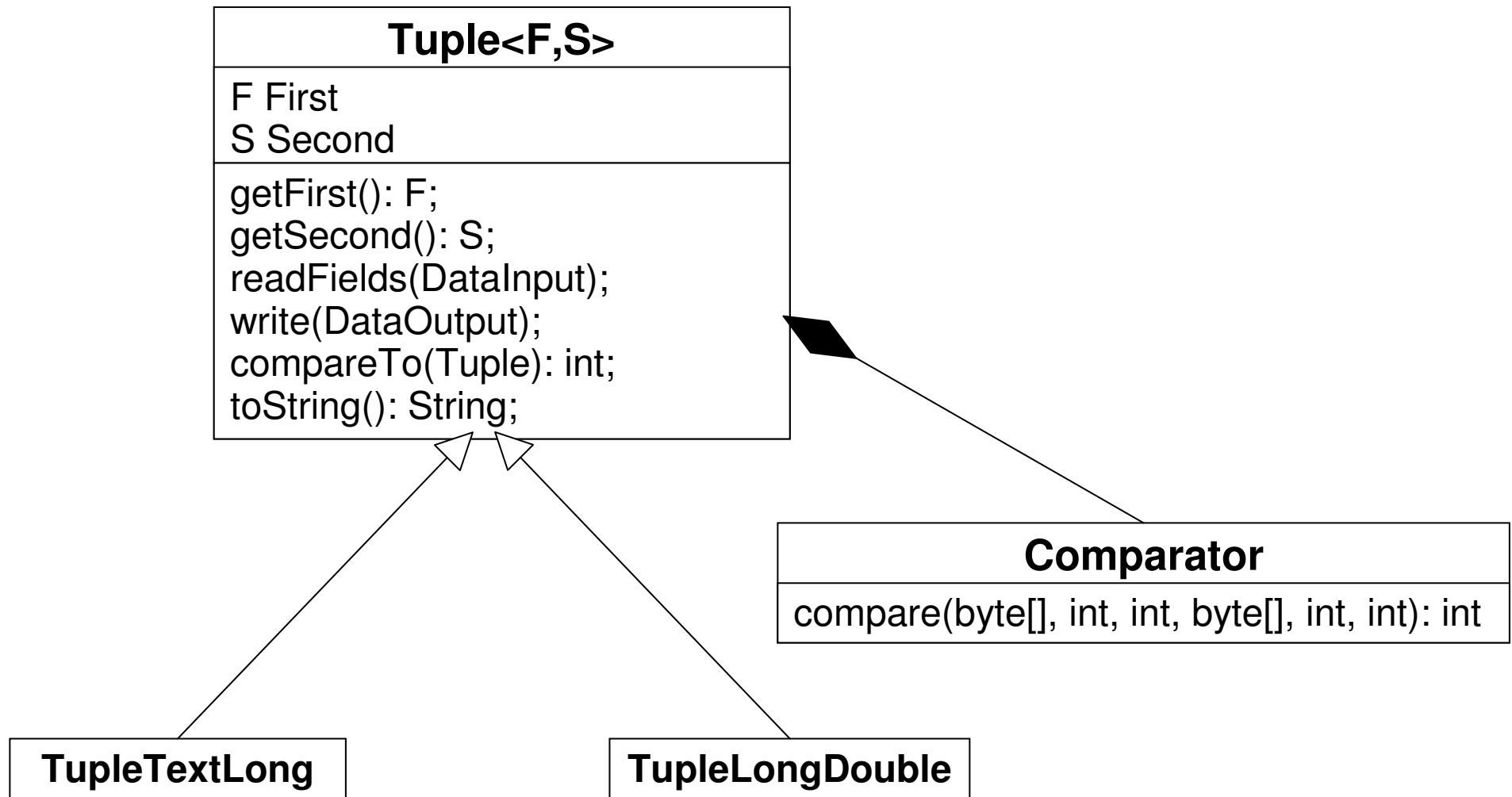
Implementierung – Eigene Datentypen

Tupel



Implementierung – Eigene Datentypen

Tupel



Speicherauslastung – 2 Phasen

- Phase 1
 - Map: konstant
 - Reduce: konstant – 1 long
- Phase 2
 - Map: 1 Dokument, pro Wort ein float
 - Reduce: $|\{\text{Wort} \in \text{Korpus}\}| * (\text{Wort}, \text{Long}, \text{Double})$
 - Abschätzung Obergrenze:

$$|\{\text{Wort} \in \text{Korpus}\}| = |\text{Korpus}|$$

$$1'000'000 * (10+8+8 \text{ Byte}) \approx 25 \text{ MiB}$$

$$3 \text{ GiB RAM} \approx 3,2 \text{ Mrd. Dokumente}$$

Speicherauslastung – 4 Phasen

- Phase 1
 - Map: konstant
 - Reduce: konstant – 1 long
- Phase 2
 - Map: 1 Dokument, 1 int
 - Reduce: konstant, 1 int
- Phase 3
 - Map: 1 Tupel (String, Long, Double)
 - Reduce: { (Wort, tf) } eines Dokuments
- Phase 4
 - Map: 1 Tupel (String, Long, Double)
 - Reduce: | {Wort \in Korpus} | * (String, Double)

Speicherauslastung – 5 Phasen

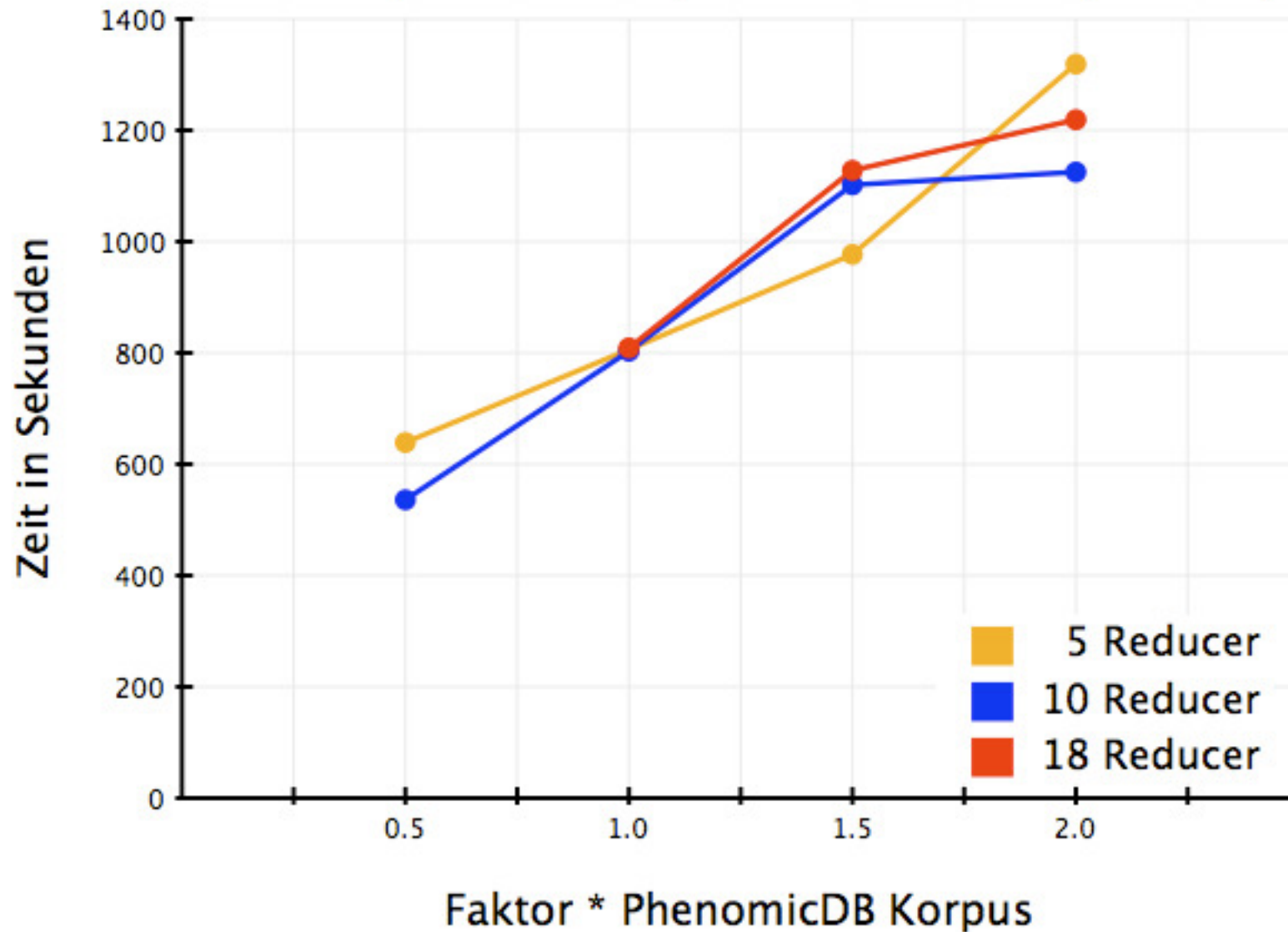
- Phase 1
 - Map: konstant
 - Reduce: konstant – 1 long
- Phase 2
 - Map: 1 Dokument, 1 int
 - Reduce: konstant, 1 int
- Phase 3
 - Map: 1 Tupel (String, Long, Long)
 - Reduce: 1 Tupel (String, Long, Long)
- Phase 4 & Phase 5
 - Map: 1 Tupel (String, Long, Double)
 - Reduce: 1 Tupel (String, Long, Double)

Obergrenze – 5 Phasen

- durch Phase 2 beschränkt
 - 1 Wort => max. INT_MAX Vorkommen im Korpus
 - 1 Dokument => max. LONG_MAX Worte
 - max. LONG_MAX Dokumente
- Abschätzung: $10^{19} * 10^{19} * 10\text{Byte} = 10^{39}$ Byte
- => unbegrenzt (10^{80} Atome im Universum...)

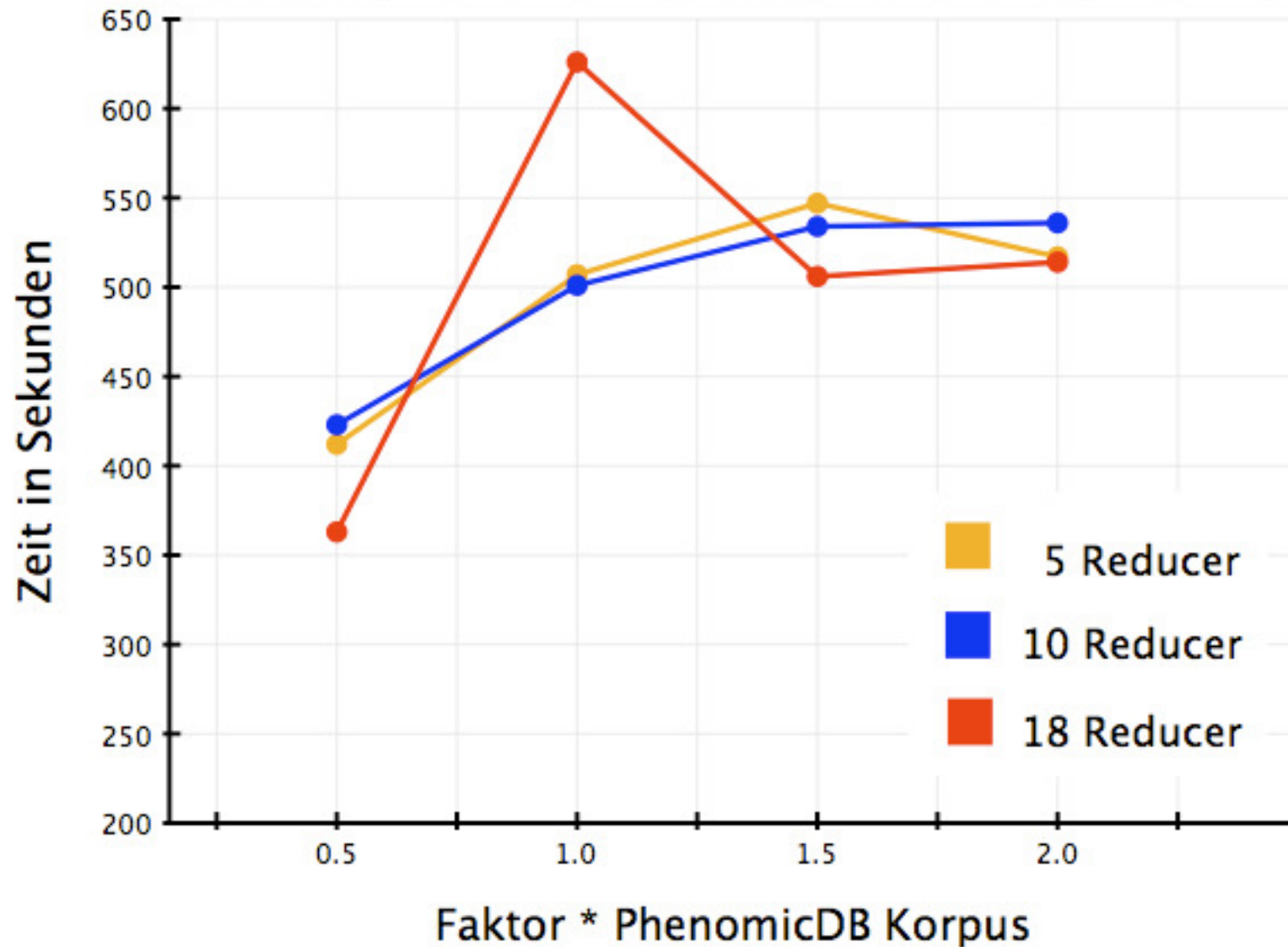
Laufzeit

4-Phasen Variante, PhenomicDB (320'000 Dokumente, 110 MiB)

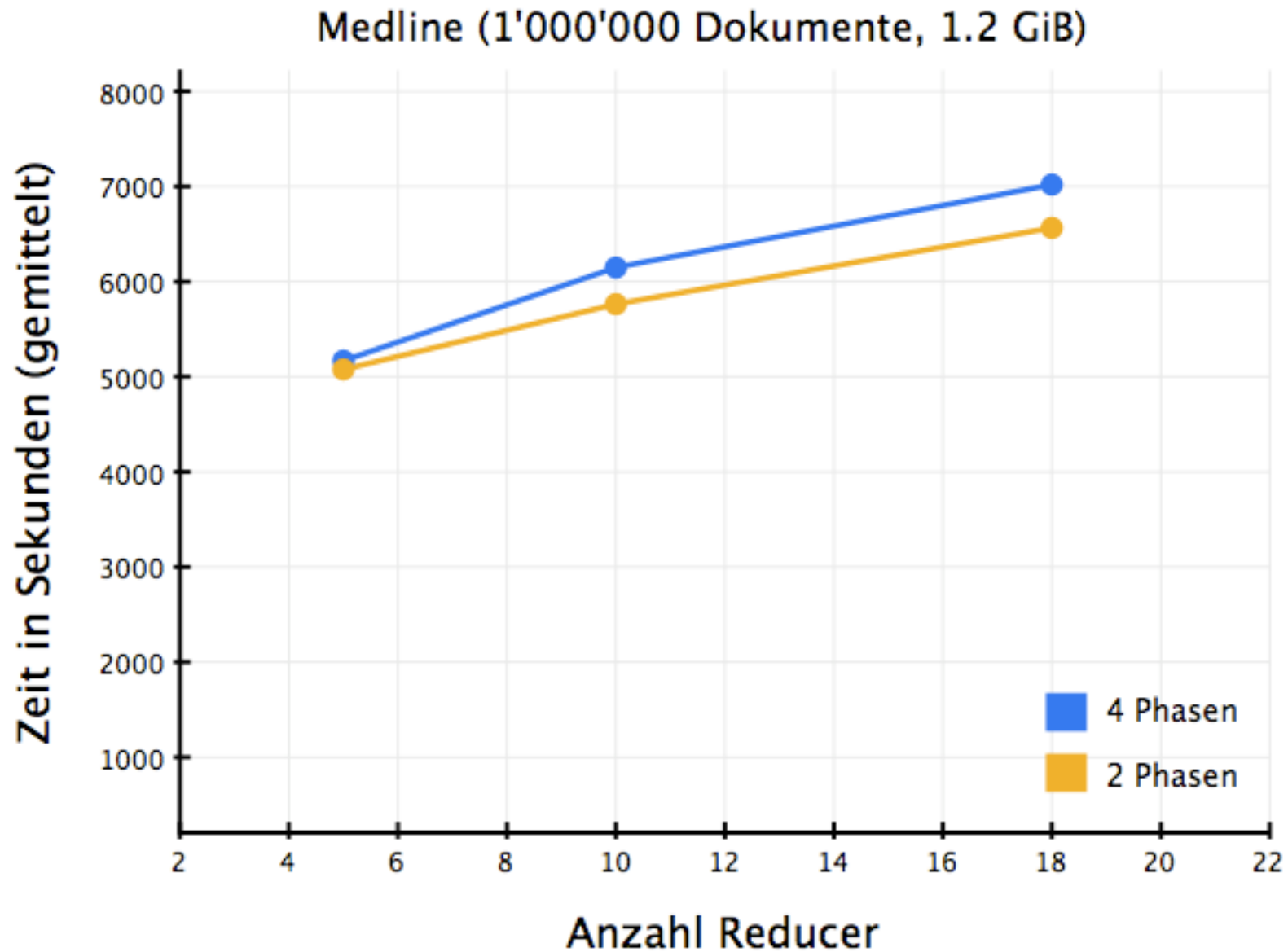


Laufzeit

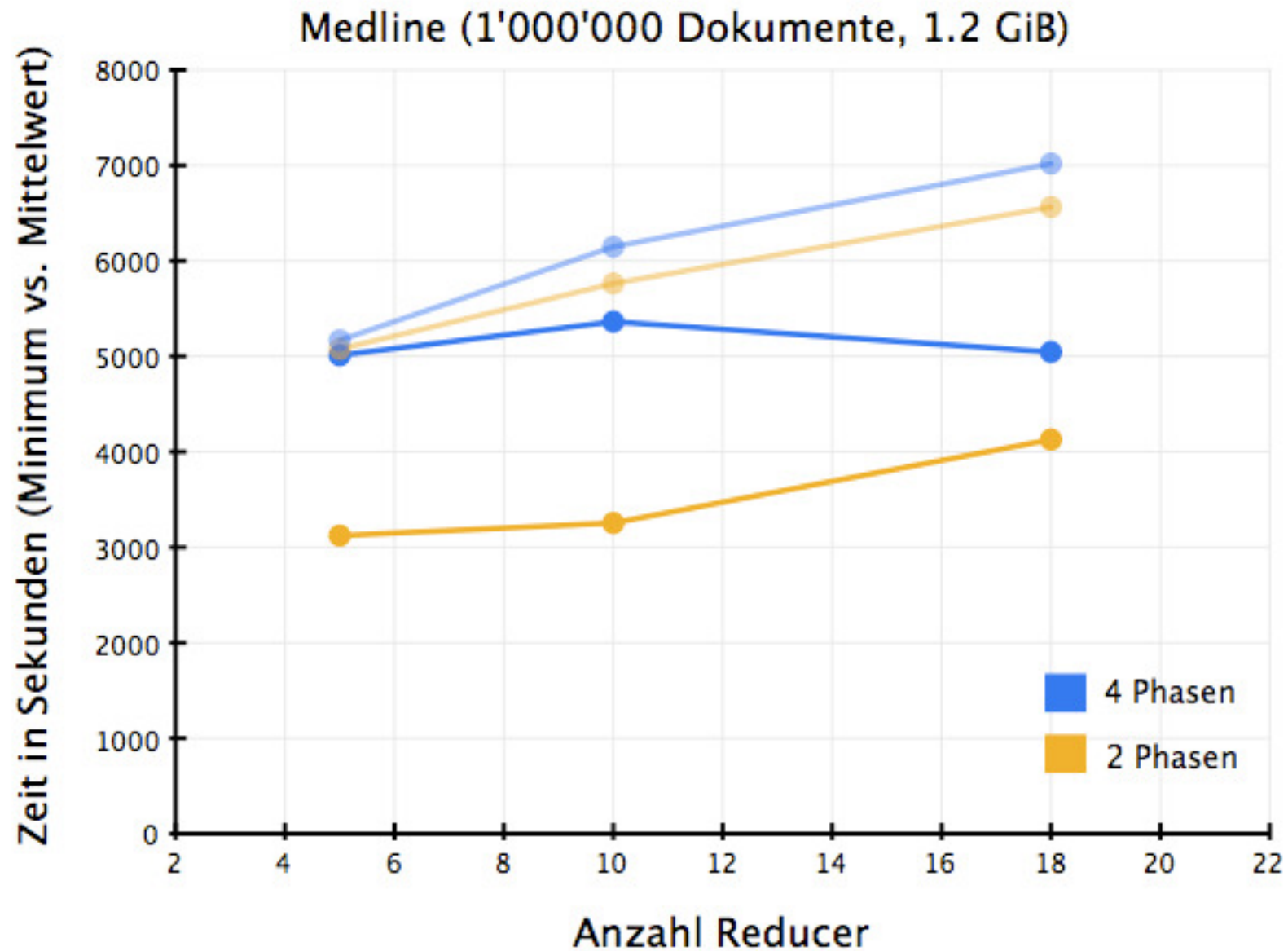
2-Phasen Variante, PhenomicDB (320'000 Dokumente, 110 MiB)



Laufzeit

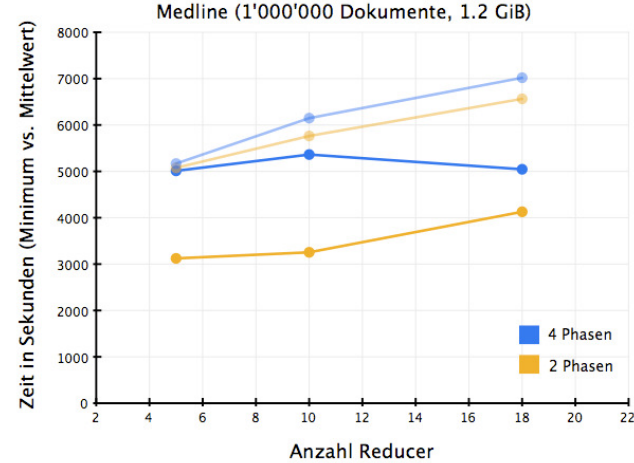
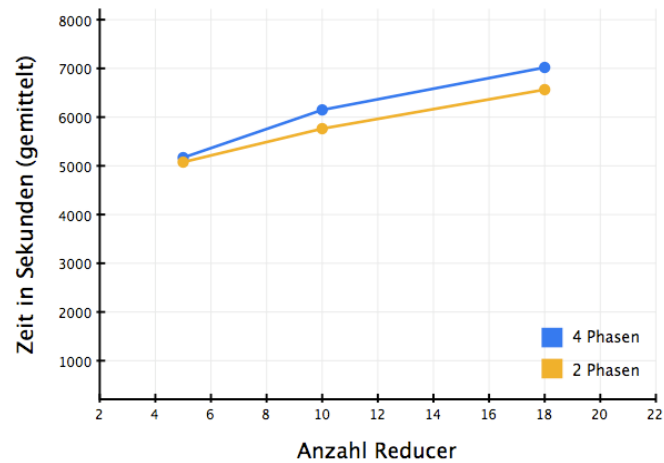
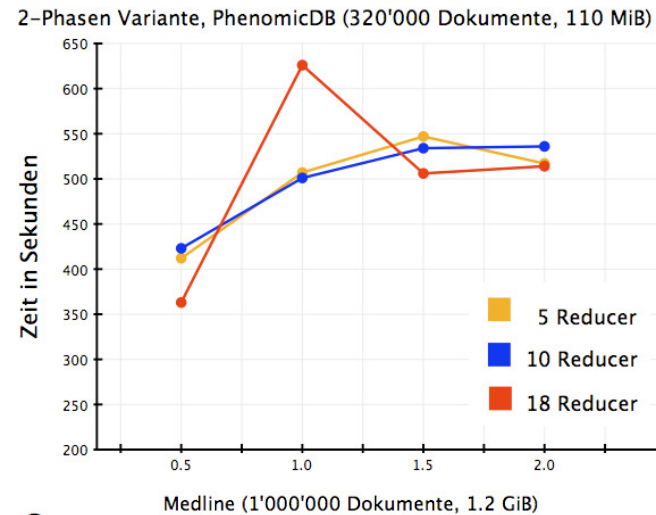
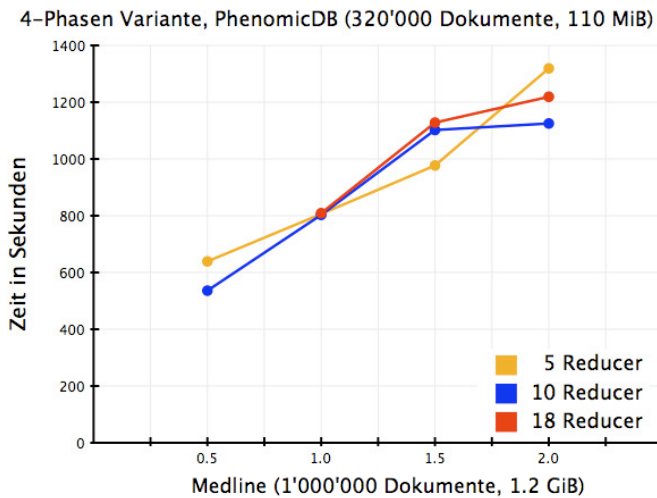


Laufzeit



Fazit

- beide Varianten bei größeren Korpora anwendbar, max. $O(N)$
- Vergleich bei großen Korpora nötig (Wikipedia!)



Quellen

- Hadoop API Documentation
<http://hadoop.apache.org/core/docs/r0.20.0/api/>
- Yahoo! Hadoop Tutorial
<http://developer.yahoo.com/hadoop/tutorial/>
- Wikipedia
<http://en.wikipedia.org/wiki/Tf-idf>
- Introduction to Information Retrieval
<http://nlp.stanford.edu/IR-book/html/htmledition/tf-idf-weighting-1.html>