

# Read as Needed: Building WiSER, a Flash-Optimized Search Engine

Jun He<sup>1</sup>, Kan Wu, Sudarsun Kannan<sup>#</sup>,  
Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau

Department of Computer Sciences, University of Wisconsin - Madison

<sup>#</sup> Department of Computer Science, Rutgers University



<sup>1</sup>*Jun He is now at Google.*

# SSDs are Fast



■ ■ ■

- Sequential read: 3.5GB/s
- Random read: 500,000 IOPS
- 0.17 dollar/GB

*\*Not an endorsement for any brand*

# **Many applications/systems have been optimized for SSDs**

## **Key-value stores**

- RocksDb, Wisckey, ...

## **Graph stores**

- FlashGraph, Mosaic, ...

## **File systems**

- SFS, F2FS, ...

**...**

# Search engines are overlooked

Search engines are important and widely used

DB-Engines Ranking	
Feb 2020	DBMS
1.	Oracle +
2.	MySQL +
3.	Microsoft SQL Server +
4.	PostgreSQL +
5.	MongoDB +
6.	IBM Db2 +
7.	Elasticsearch +
8.	Redis +
9.	Microsoft Access
10.	SQLite +

- Wikipedia
- Github
- Uber
- ...

# **Search engines are challenging for storage systems**

## **Low data latency**

- queries are interactive

## **High data throughput**

- engines retrieve info from a large amount of data

## **High scalability**

- data grows over time

# Just use more RAM?

## **RAM is critical in existing search engines**

- “your RAM will limit all other resources”  
— an Elasticsearch user

## **OK at small scale**

- total cost is low

## **Cost prohibitive at large scale**

- data grows fast
- may waste bandwidth: rarely read and process 100GB/s



Can search engines perform well with  
a **small** memory and a **fast** SSD?

RAM

Dataset

Working set

A diagram illustrating memory hierarchy. At the top is a small black rounded square labeled 'RAM'. Below it is a horizontal bar representing the 'Dataset'. The bar is divided into three sections: a small light gray section on the left, a larger black section in the middle labeled 'Working set', and a medium-sized light gray section on the right.

# Yes

## **Built an engine that performs well on small memory**

- could outperform other engine with entire dataset in memory

## **Our philosophy: read as needed**

- use small memory
- read data from SSDs as needed
- do not attempt to cache data in memory
- attempt to read data from SSDs efficiently

## **SSDs have high read bandwidth and IOPS**

- read an entire drive (1TB) in 5 minutes
- higher bandwidth is possible with RAID

## **Computation/network is more likely to be the bottleneck**



# **To optimize a read-as-needed system, we must:**

## **Reduce read amplification**

- SSD bandwidth is still limited

## **Hide I/O latency**

- SSD latency is still high

## **Use large requests to improve device efficiency**

- SSDs favor large requests

# **We apply four techniques to build a read-as-needed search engine**

## **1. Grouping data by term**

- reduce read amplification
- use large requests

## **2. Two-way Cost-aware Bloom Filters**

- reduce read amplification

## **3. Adaptive Prefetching**

- hide I/O latency

## **4. Trade Disk Space for I/O**

- reduce read amplificationm

# Our Contributions

## **Four techniques to improve search engine performance**

- small memory, large dataset, fast SSDs

## **An open-source flash-optimized search engine from scratch (WiSER)**

- 11,000 lines of C++

## **A benchmark for evaluating search engines**

- based on Wikipedia
- a variety of queries for stressing the system

## **Better performance than the state-of-the-art Elasticsearch**

- increase end-to-end query throughput by up to 2.7x
- reduce end-to-end latency by up to 16x

# Outline

Overview

**Techniques**

**Evaluation**

**Final Thoughts**

## **Techniques**

- 1. Cross-stage Data Grouping**
- 2. Two-way Cost-aware Bloom Filters**
- 3. Adaptive Prefetching**
- 4. Trade Disk Space for I/O**

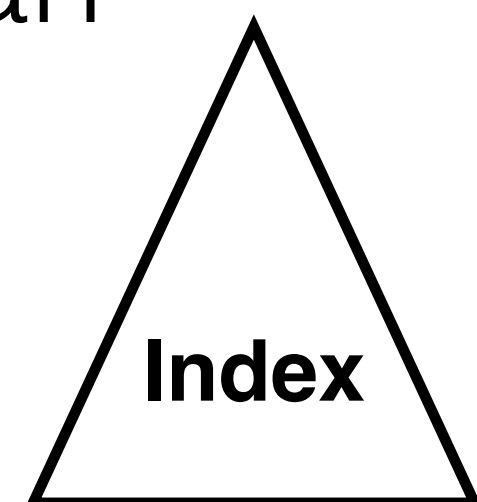
**Guided by the Wikipedia dataset**

## **Techniques**

- 1. Cross-stage Data Grouping**
2. Two-way Cost-aware Bloom Filters
3. Adaptive Prefetching
4. Trade Disk Space for I/O

# Background: stages of a search query

'tugman'



Doc IDs

2, 7, 10



Term  
Frequencies



Term  
Positions



Term Byte  
Offsets



Blair **Tugman** is  
an American  
mixed martial  
artist ...

2

Blair **Tugman** is  
an American  
mixed martial  
artist ...

Byte Byte  
6 11

Blair **Tugman** is  
an American  
mixed martial  
artist ...



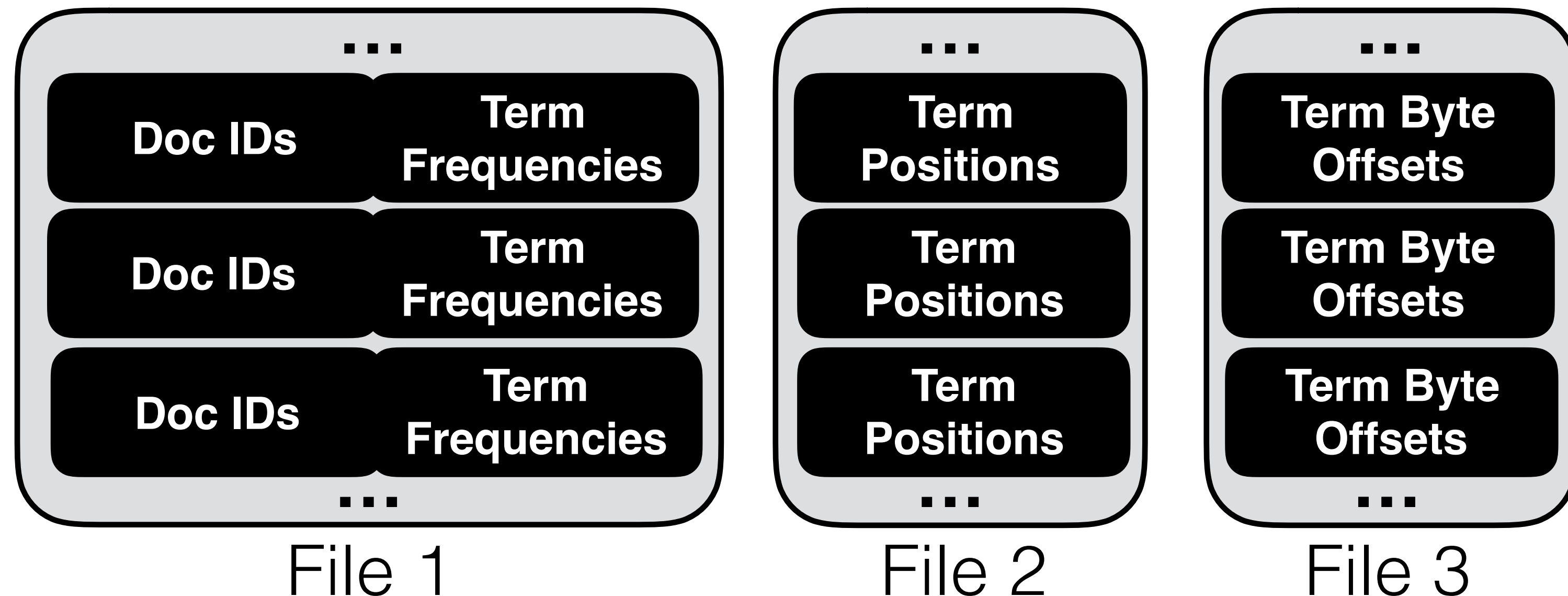


# Data in Elasticsearch is grouped by stage

(previous term)

‘tugman’

(next term)



# How much I/O is needed for 'tugman'?

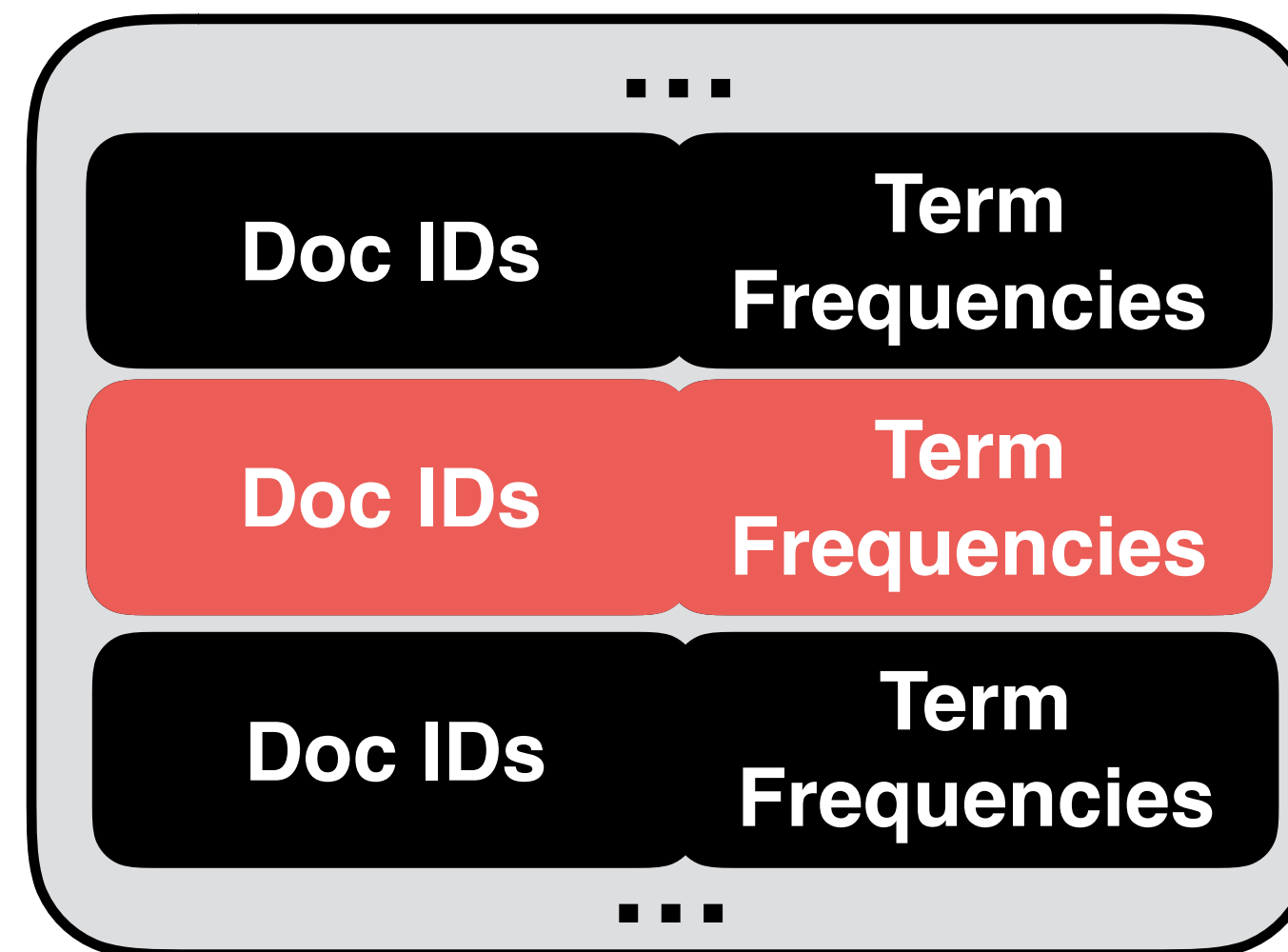
I/O Count: 

CPU

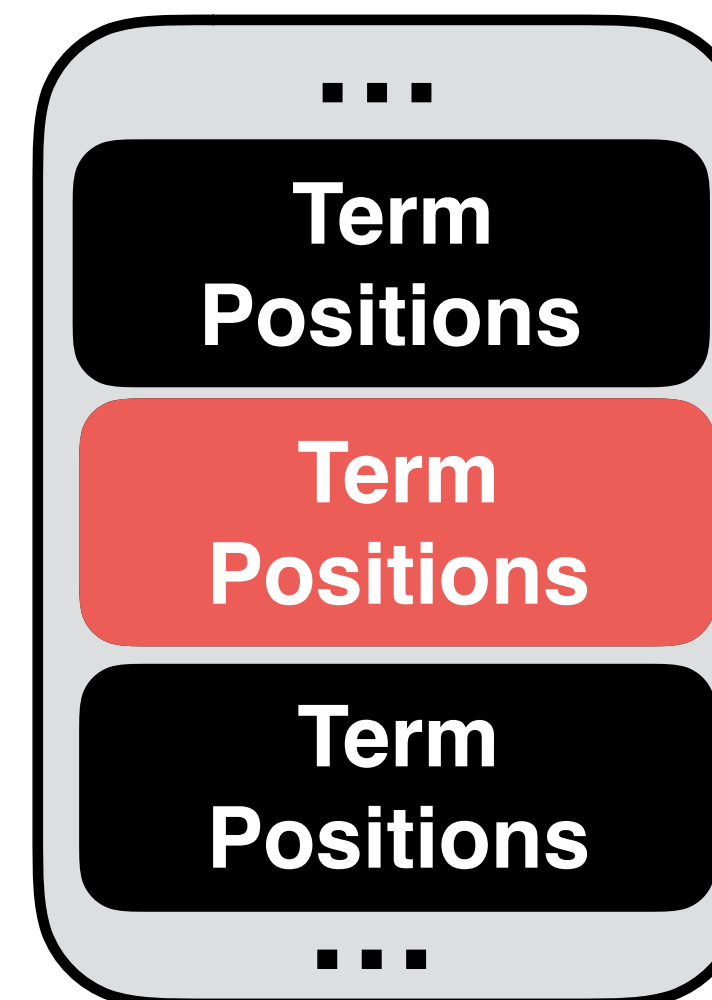
(previous term)

'tugman'

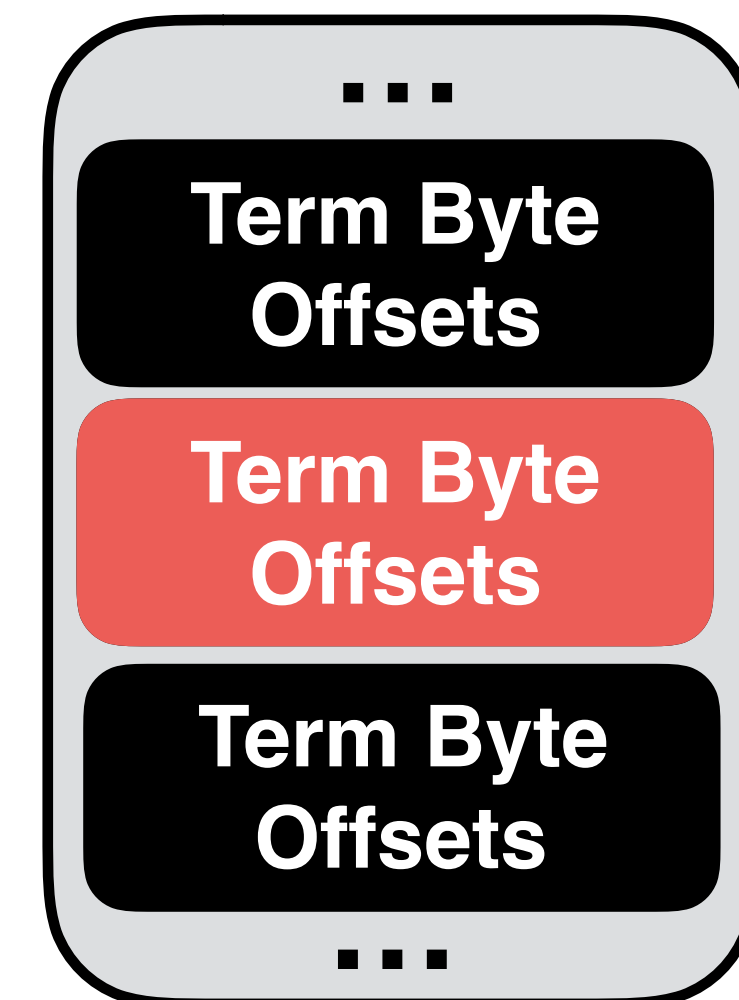
(next term)



File 1



File 2



File 3

# Transferred data is often wasted

**I/O Count: 3, Size: 12 KB**

**Data of 'tugman' can fit in 4KB**

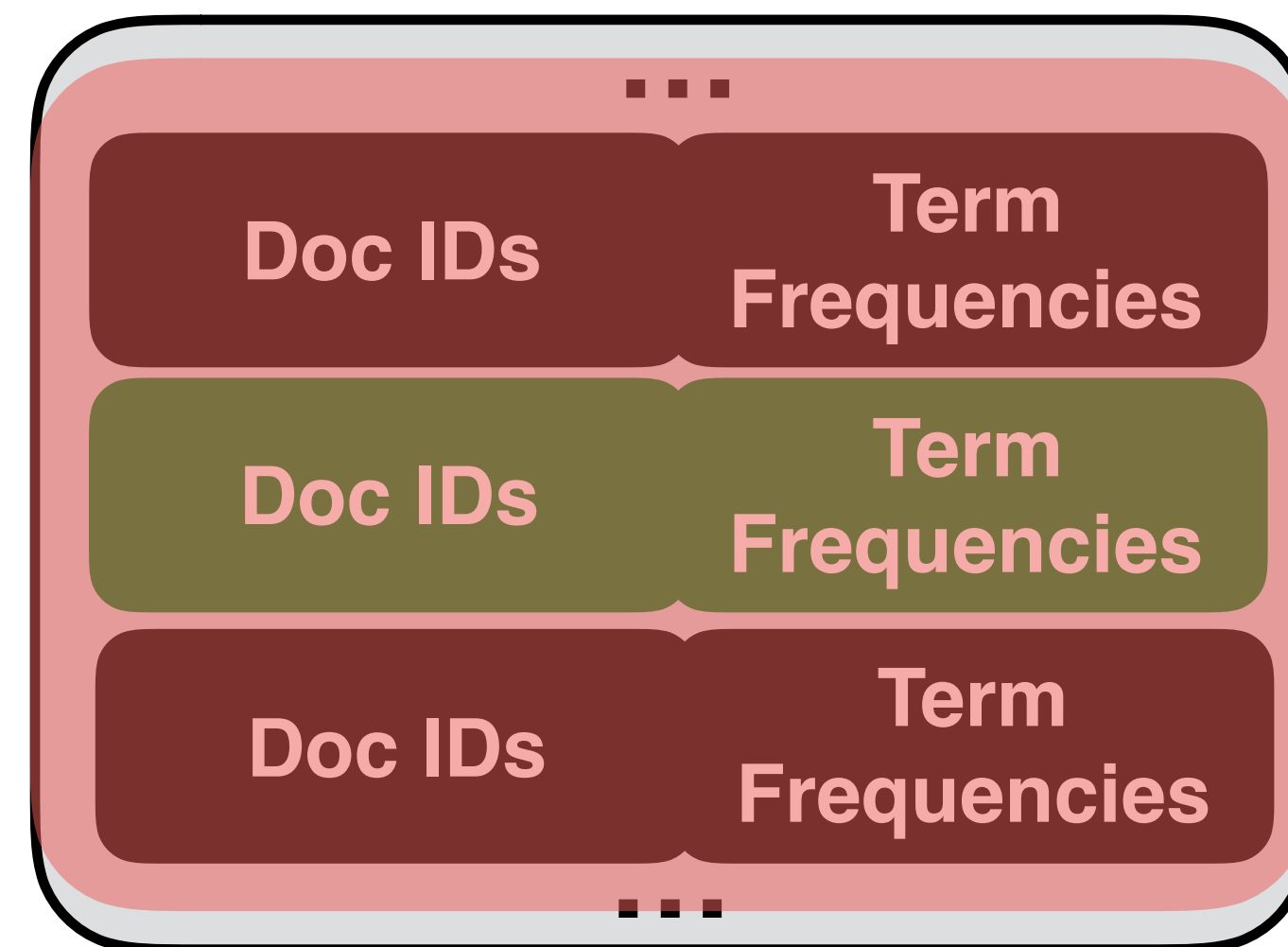
**Small term is common**

- 99% of Wikipedia terms can fit in 4 KB

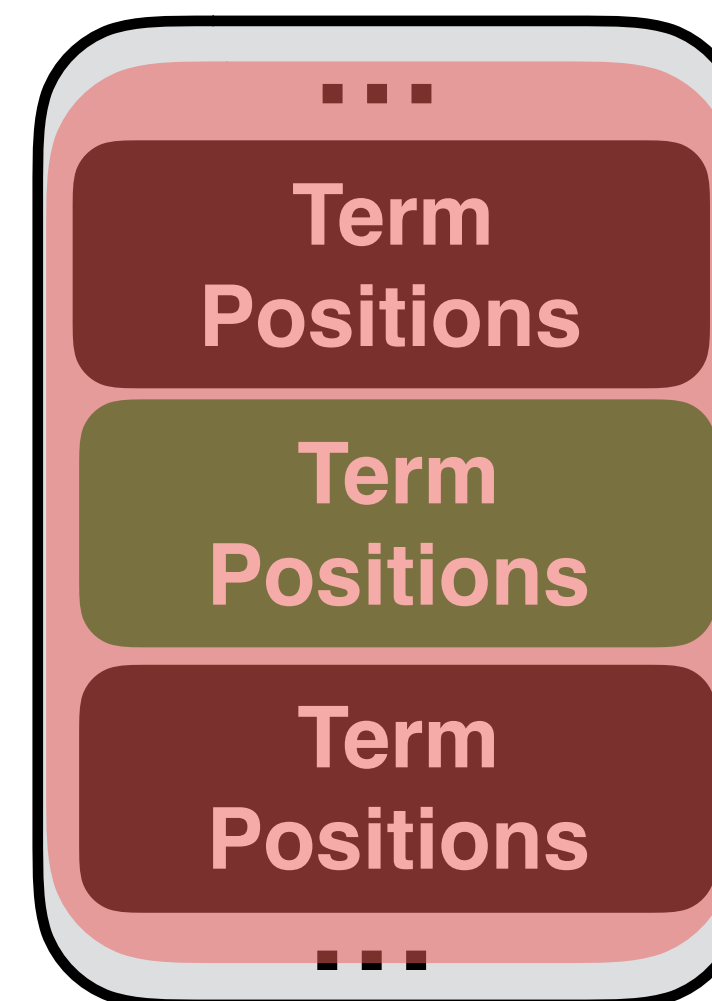
(previous term)

'tugman'

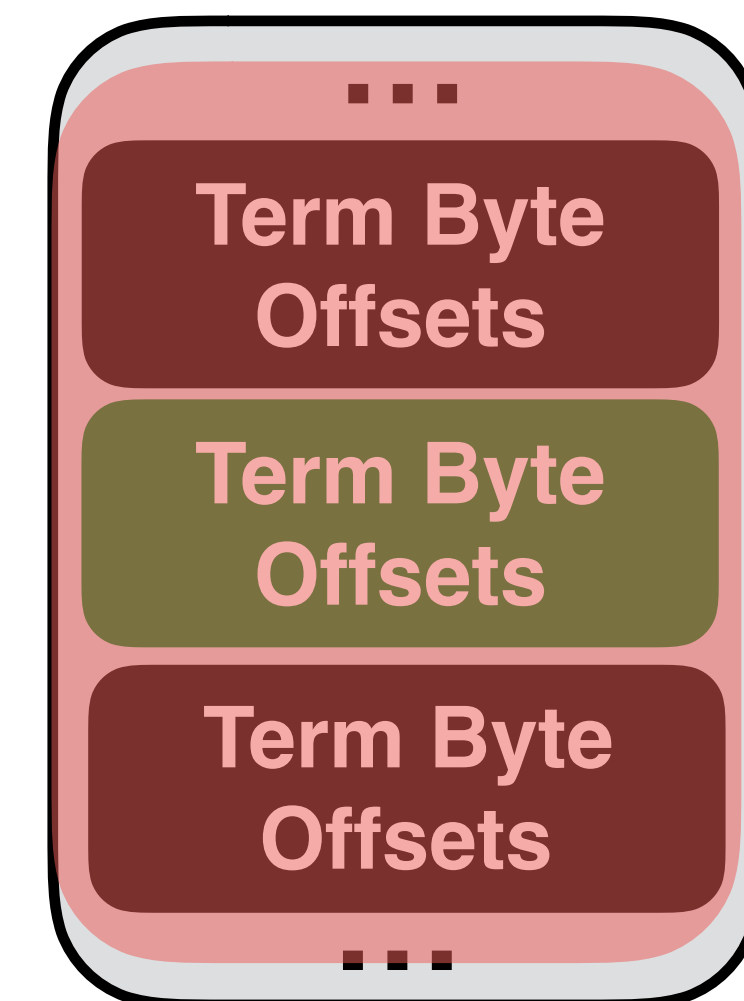
(next term)



File 1



File 2



File 3

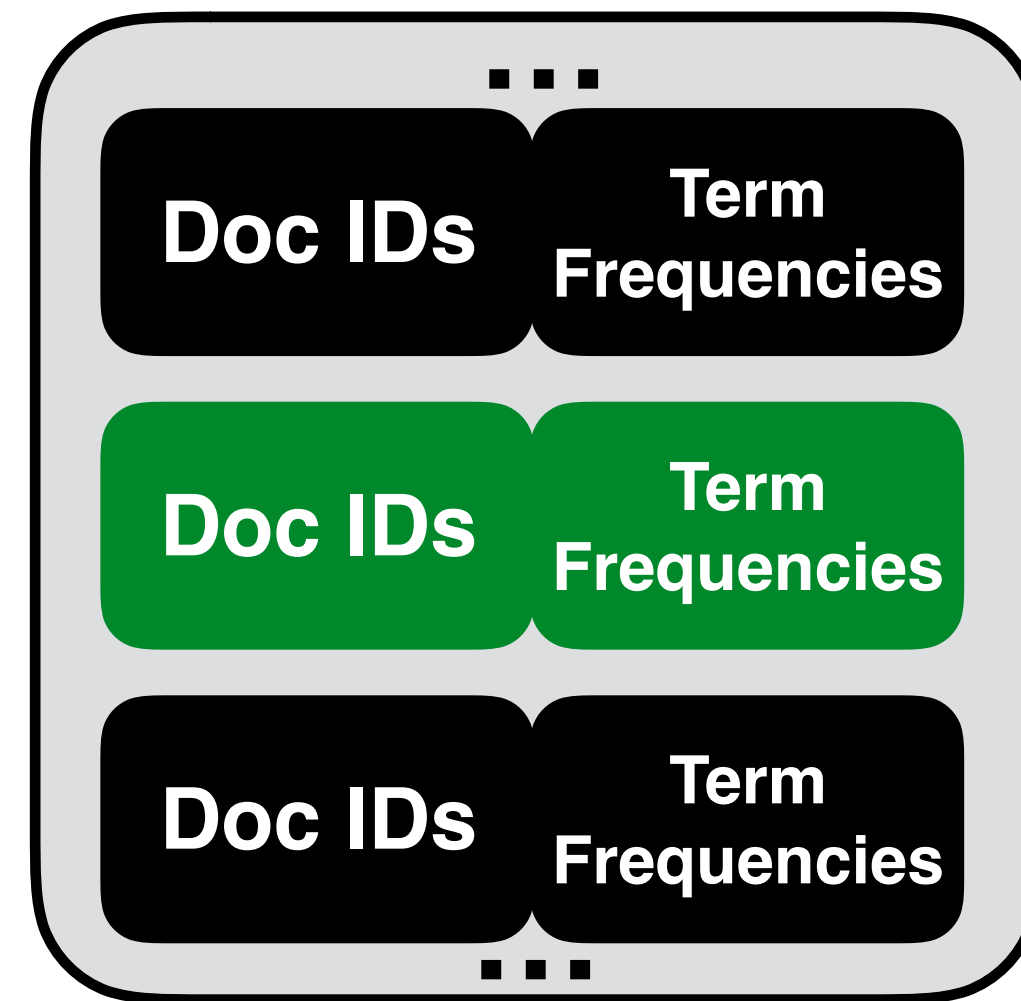
**~30 bytes**

# WiSER groups data by term

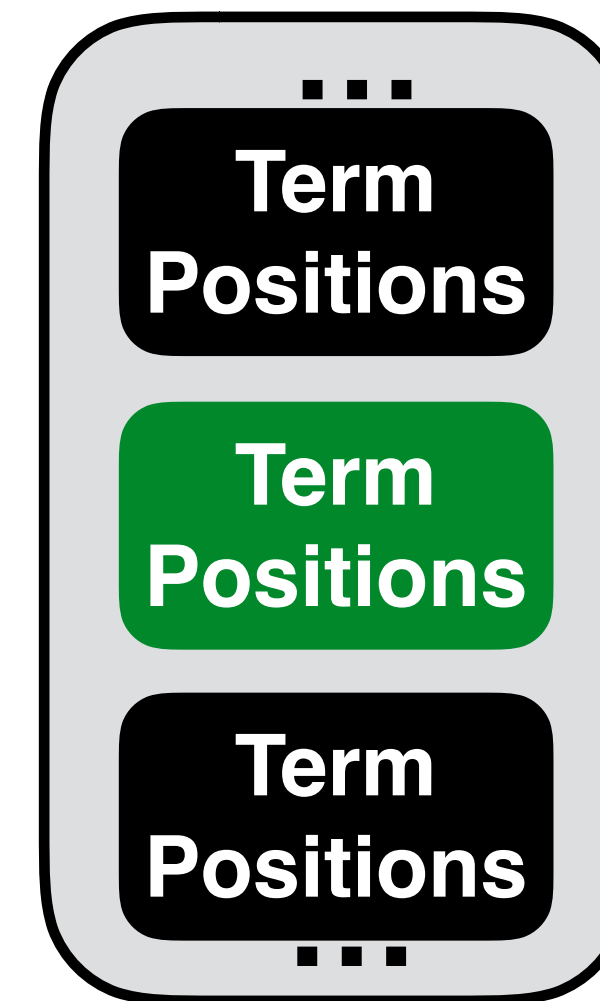
(previous term)

‘tugman’

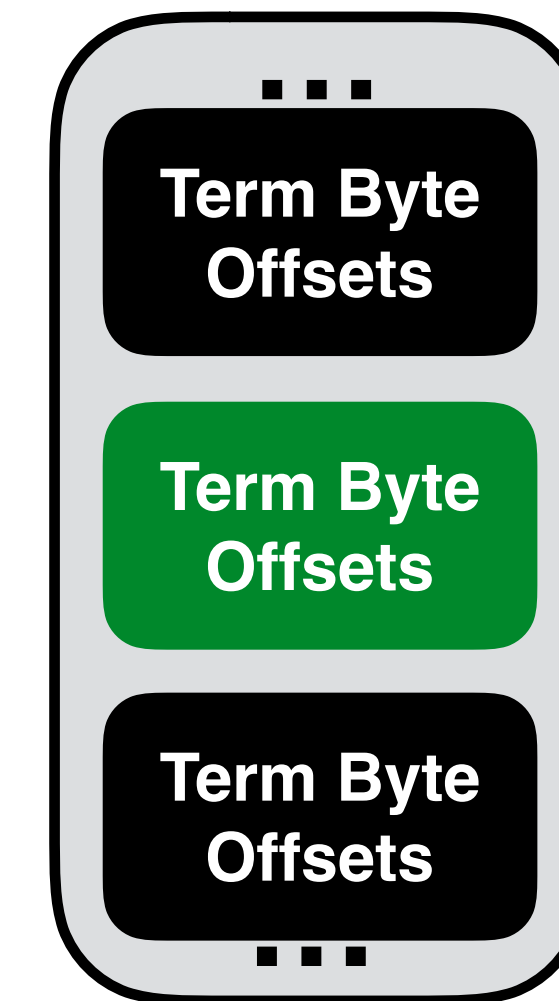
(next term)



File 1

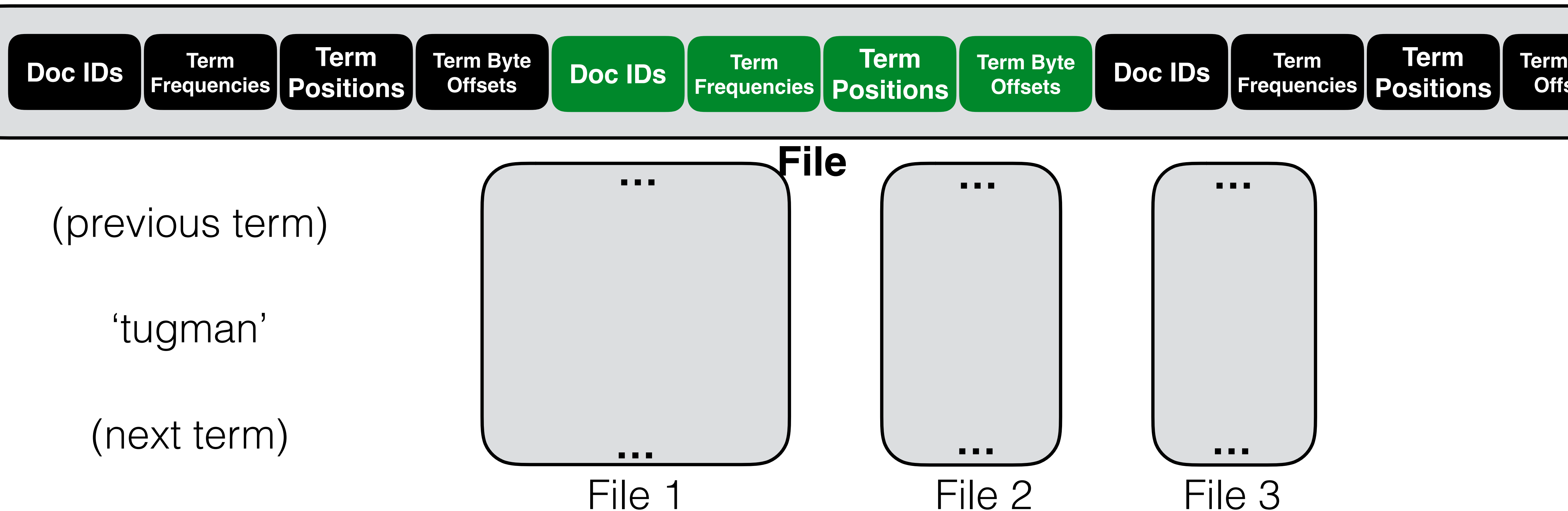


File 2



File 3

# WiSER groups data by term



# How much I/O does WiSER need for ‘tugman’?

I/O Count: 0

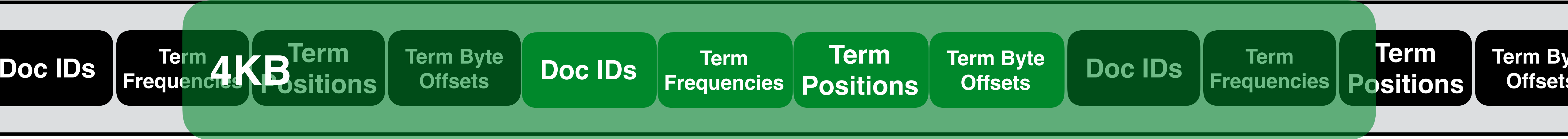
Reduce I/O by up to 3x

CPU

(previous term)

‘tugman’

(next term)



File

## **Techniques**

1. Cross-stage Data Grouping
- 2. Two-way Cost-aware Bloom Filters**
3. Adaptive Prefetching
4. Trade Disk Space for I/O



# Background: phrase queries offer high accuracy

Search phrase “**united** **AND** **states**”

The **United States** president...

The spokesman of **united** airlines **states** that...

# Background: phrase queries offer high accuracy

But demand a lot more data

**Term**  
united  
states

**Position**  
2,X,X,X,X,X,X...  
3,X,X,X,X,X,X...

Our **United States**  
president...

**Position**  
4,X,X,X,X,X,X...  
6,X,X,X,X,X,X...

The spokesman of  
**united** airlines **states**  
that...

**Tried regular Bloom filters, it doesn't work...**

# WiSER adds two Bloom filters to enable “two-way” filtering

Term	bloom filter <b>before</b>	bloom filter <b>after</b>	Position
united	our	states	2
states	united	president	3
united	of	airline	4
states	airlines	that	6

Our **United States**  
president...

The spokesman of  
**united** airlines **states**  
that...

Check if “united states” is a phrase in a document by:

1. Checking if “states” is in “united”.after
2. Or, checking if “united” is in “states”.before.

**Why adding two filters?**

**WiSER dynamically chooses the  
smaller filter for filtering (two-way)**

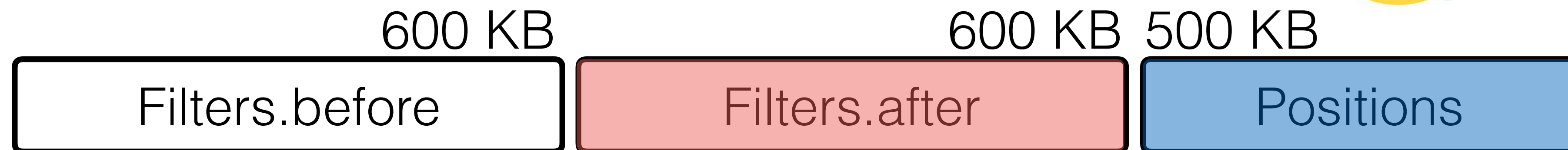
# WiSER dynamically chooses the smaller filter for filtering (two-way)

Check if 'states' is after 'united'

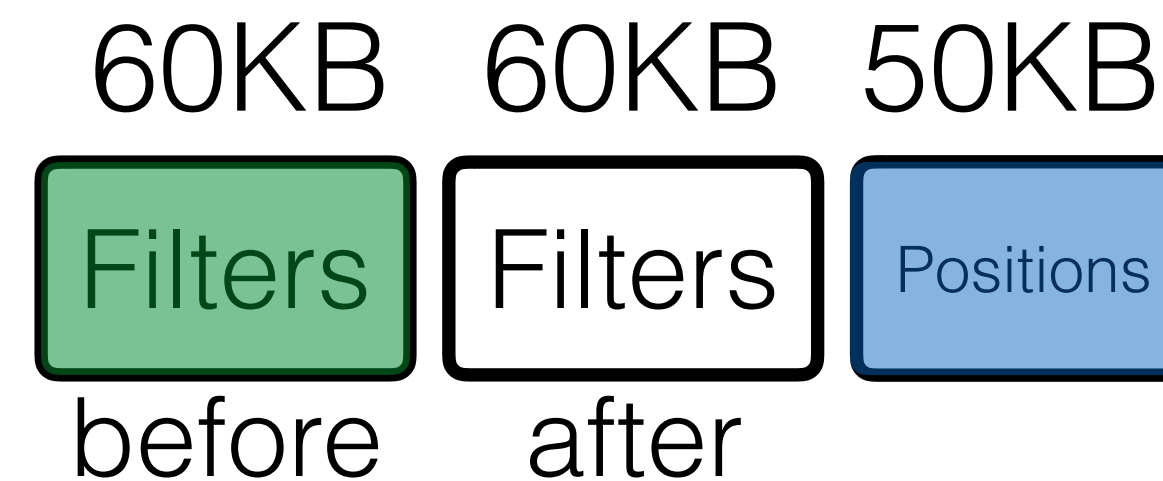
$$600 \text{ KB} > 500 \text{ KB} + 50 \text{ KB}$$



**united**



**states**



Check if 'united' is before 'states'

$$60 \text{ KB} < 500 \text{ KB} + 50 \text{ KB}$$



**WiSER only uses filters when  
it can reduce I/O (cost-aware)**



# WiSER only uses filters when it can reduce I/O (cost-aware)

Check if 'states' is after 'united'

$$600 \text{ KB} > 200 \text{ KB} + 200 \text{ KB}$$

Not worth it

**united**



**states**



Check if 'united' is before 'states'

$$600 \text{ KB} > 200 \text{ KB} + 200 \text{ KB}$$

Not worth it

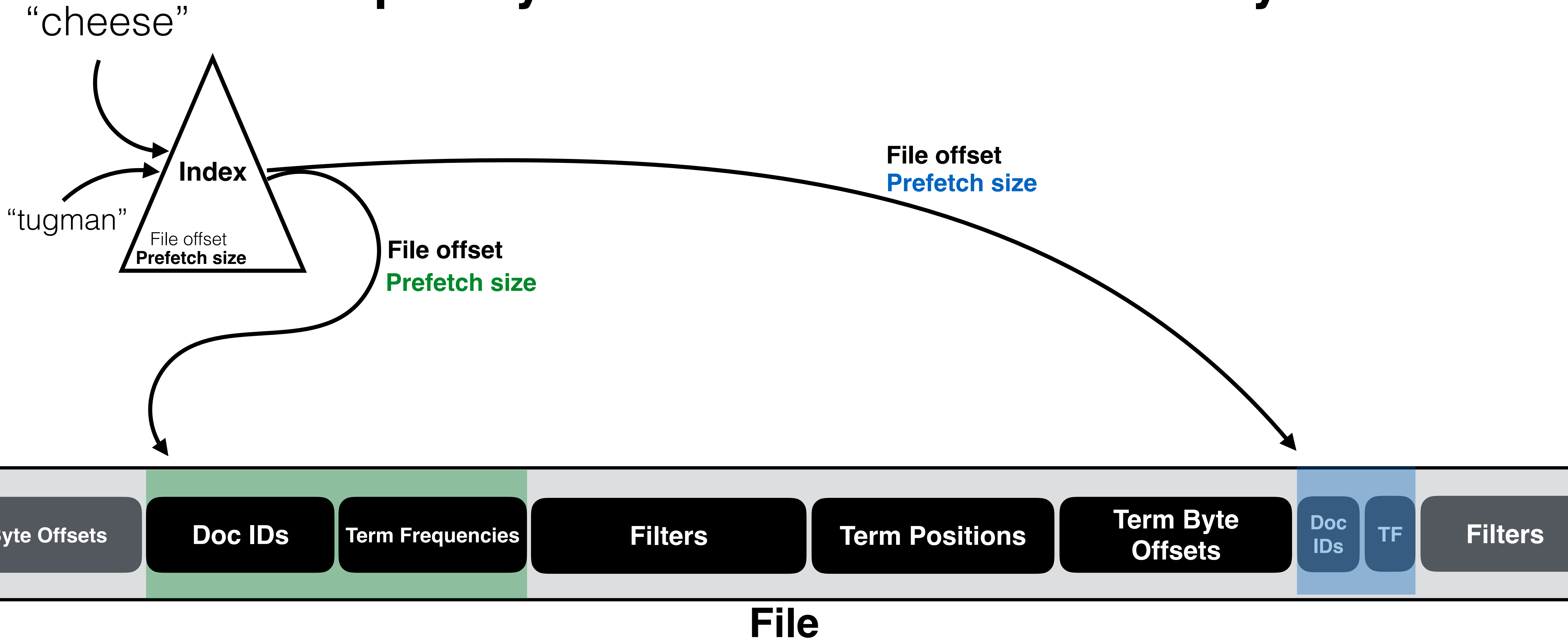
## **Techniques**

1. Cross-stage Data Grouping
2. Two-way Cost-aware Bloom Filters
- 3. Adaptive Prefetching**
4. Trade Disk Space for I/O

**Elasticsearch relies on the OS to prefetching...**

**WiSER adaptively prefetches  
frequently-used data to hide I/O latency**

# WiSER adaptively prefetches frequently-used data to hide I/O latency

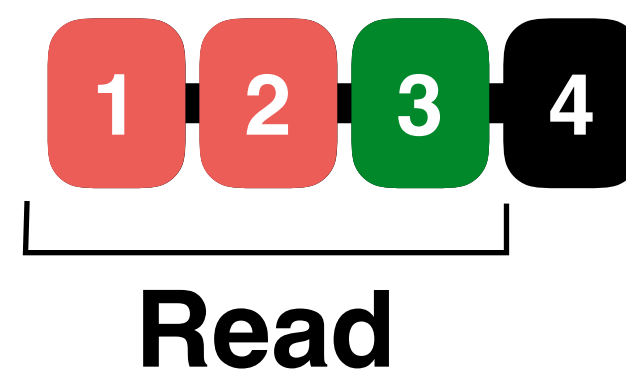
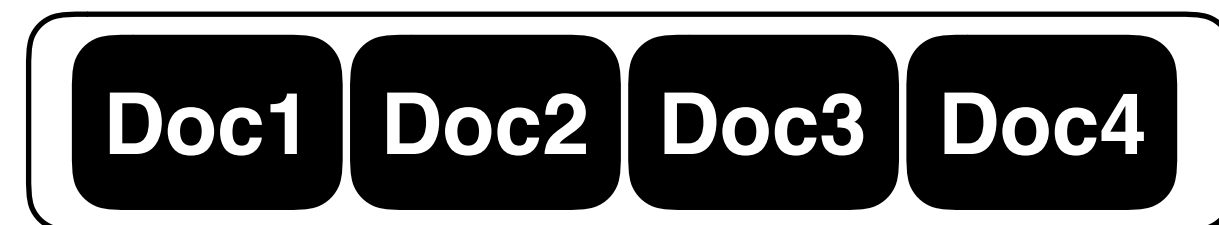


## **Techniques**

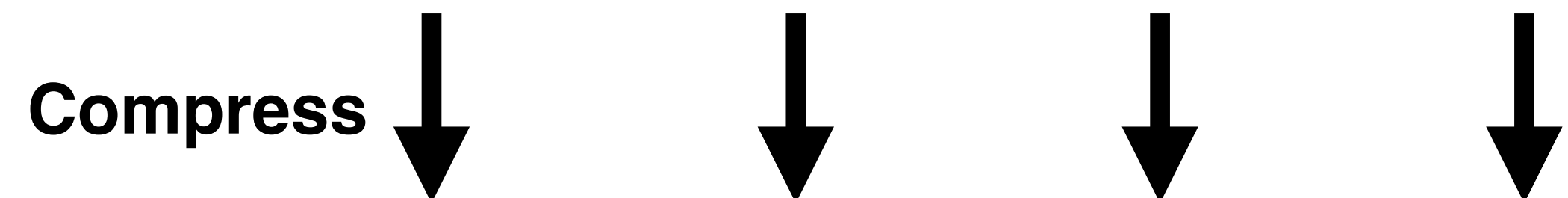
1. Cross-stage Data Grouping
2. Two-way Cost-aware Bloom Filters
3. Adaptive Prefetching
4. **Trade Disk Space for I/O**

# WiSER compresses documents individually to reduce read amplification

## Elasticsearch



## WiSER



Slightly larger  
(the tradeoff)

Read



## **Techniques**

- 1. Cross-stage Data Grouping**
- 2. Two-way Cost-aware Bloom Filters**
- 3. Adaptive Prefetching**
- 4. Trade Disk Space for I/O**

# Outline

Overview

Techniques

**Evaluation**

Final Thoughts

# Evaluate with *WSBench*

## **Dataset: Wikipedia**

- 6 million documents, 6 million terms, 18GB

## **Queries**

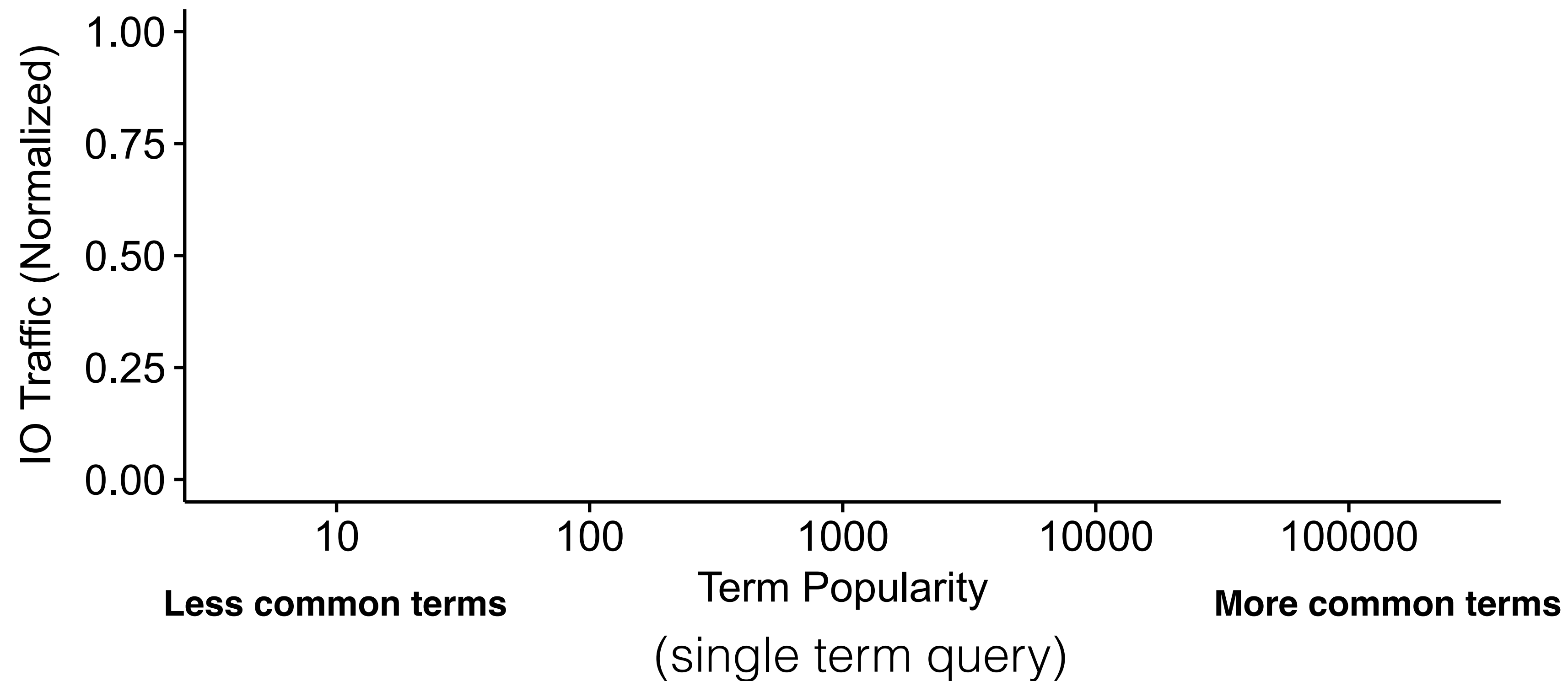
- single term queries, “and” queries, “phrase” queries, real queries
- vary term popularities in wikipedia

## **Machine**

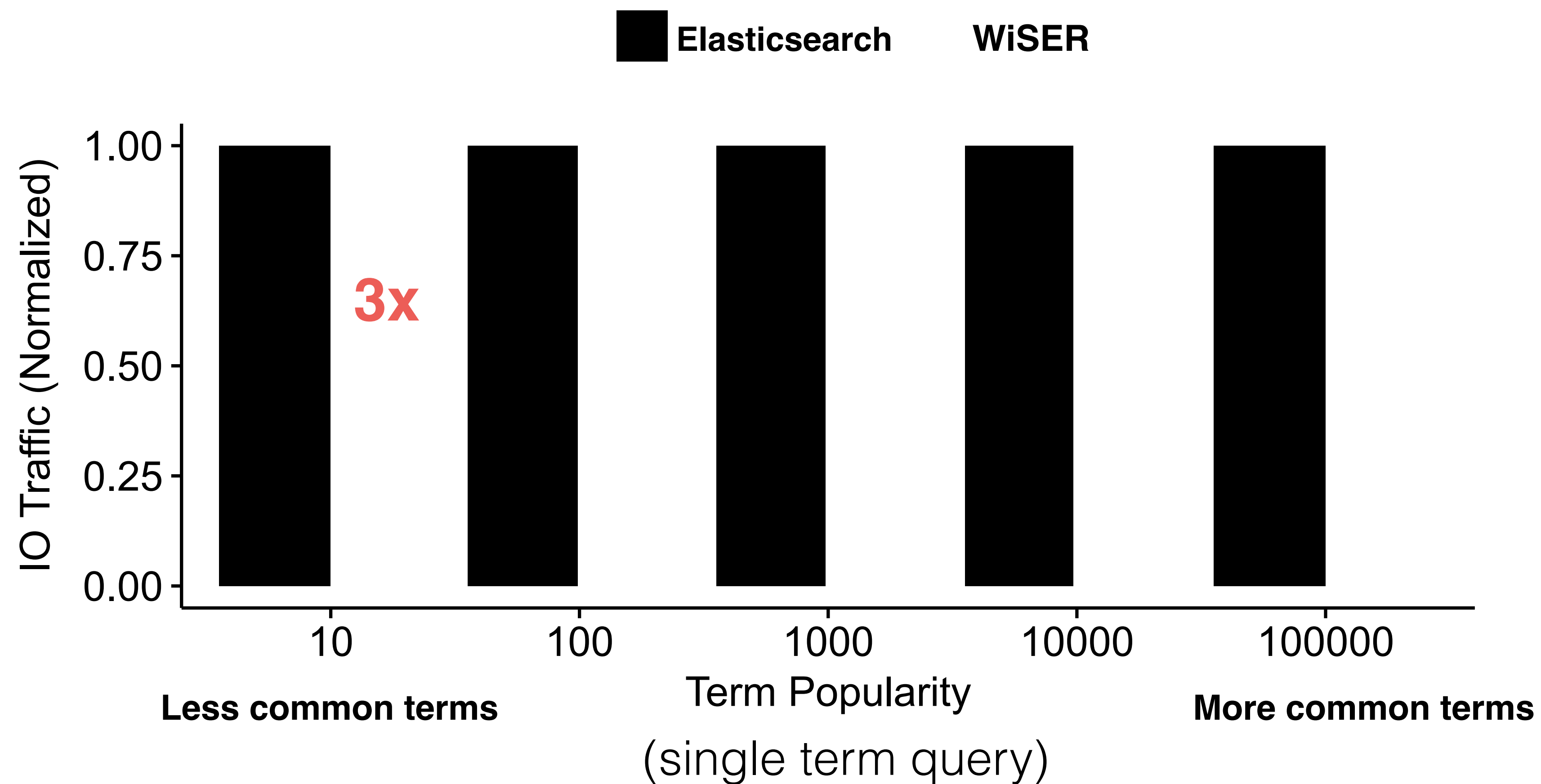
- 16 cores
- NVMe SSD: peak read 2 GB/s, 200,000 IOPS
- Linux container with 512-MB RAM
  - how well each search engine can scale up to large working sets that do not fit in main memory?

**How much read traffic can “grouping by term” (technique 1) reduce?**

# How much read traffic can “grouping by term” (technique 1) reduce?

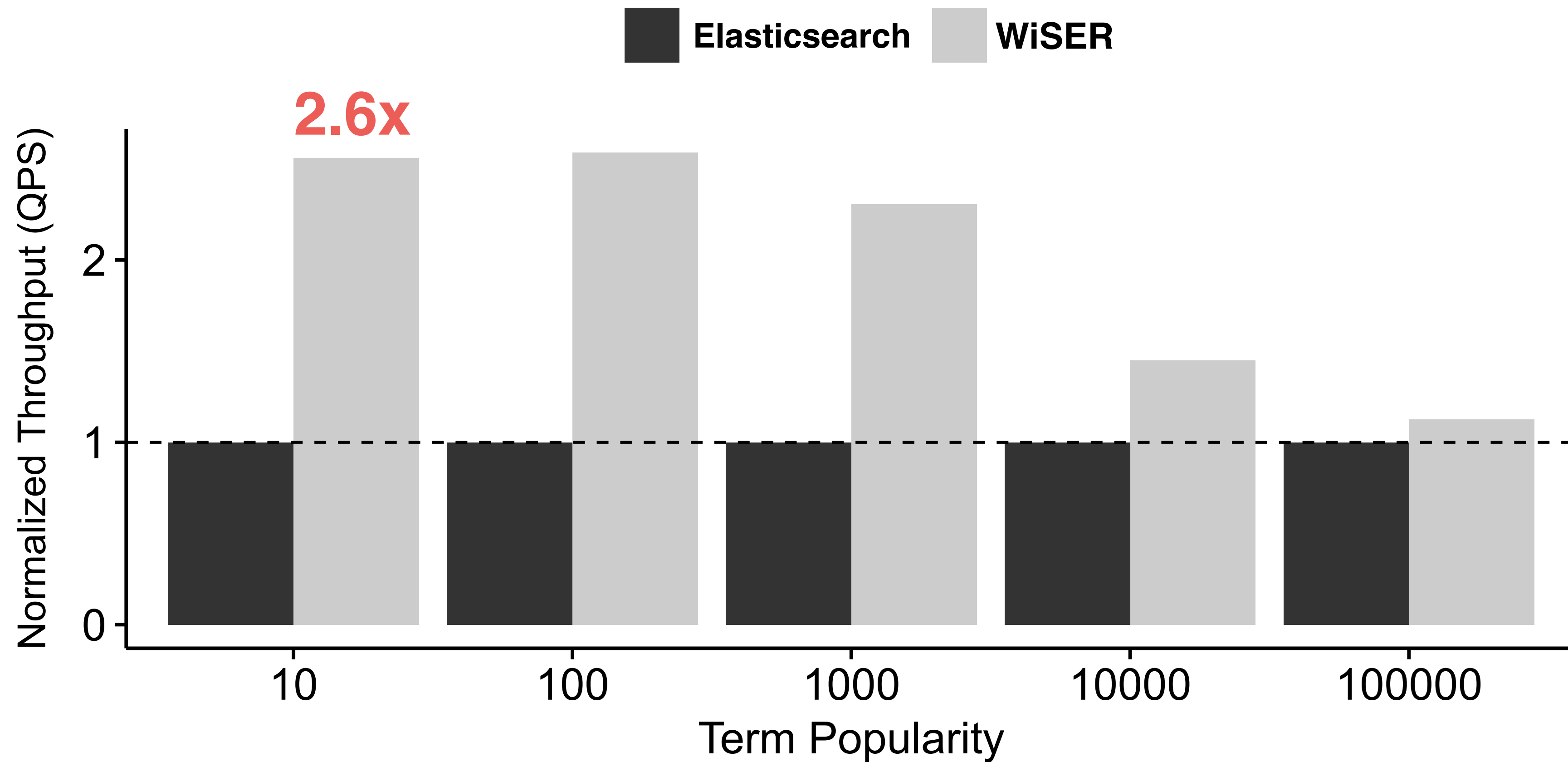


# How much read traffic can “grouping by term” (technique 1) reduce?



**How much query throughput can  
“grouping by term”(technique 1) increase?**

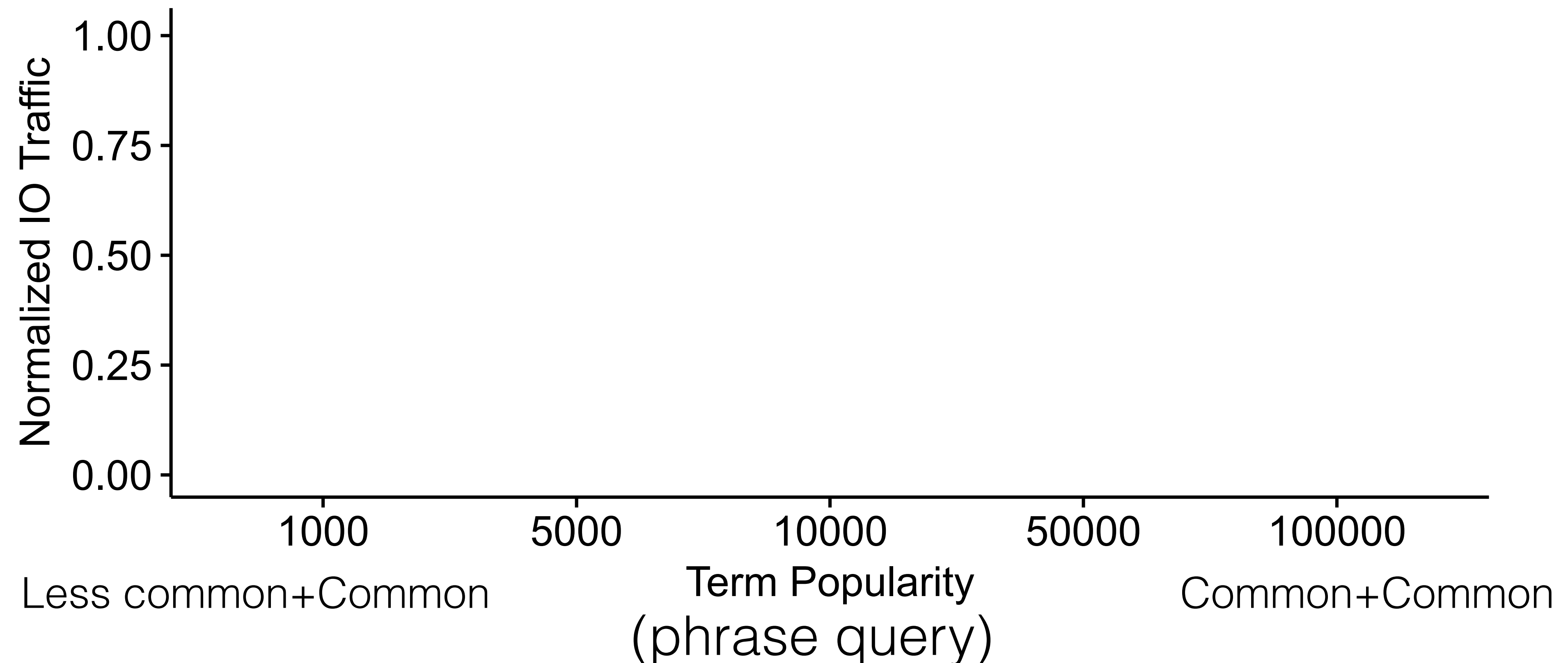
# How much query throughput can “grouping by term”(technique 1) increase?



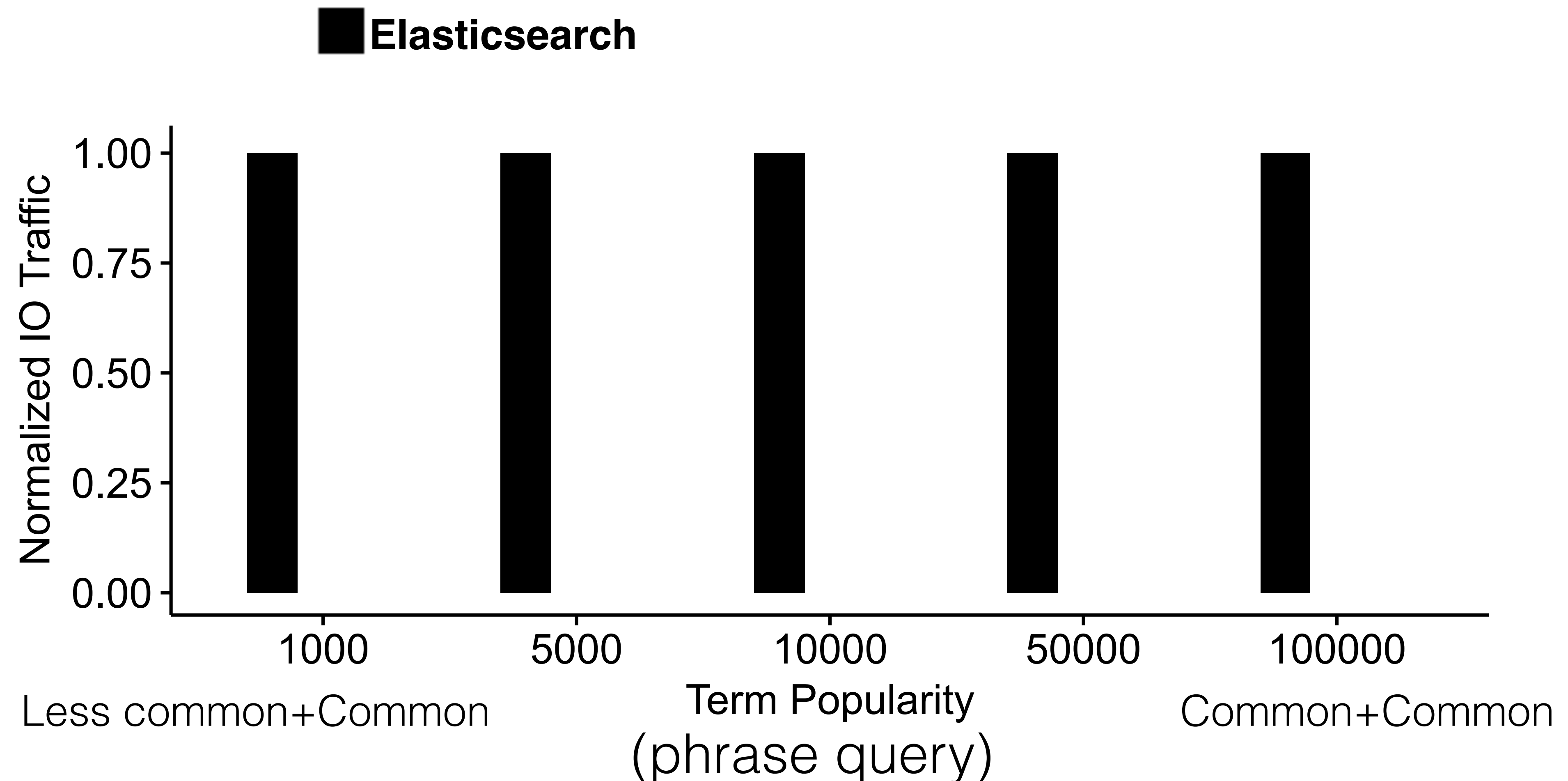


**How much read traffic can two-way cost-aware Bloom filters (technique 2) reduce?**

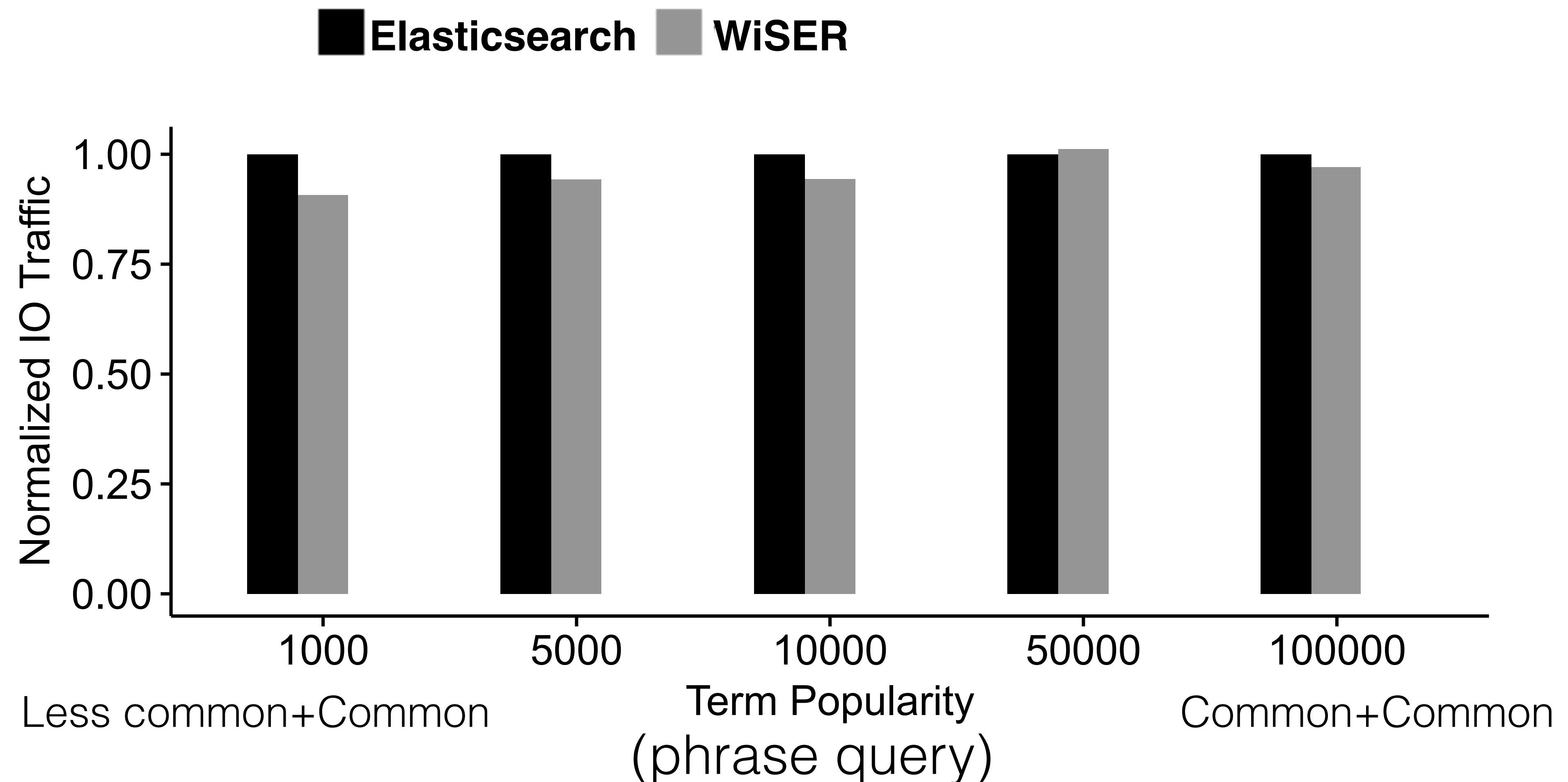
# How much read traffic can two-way cost-aware Bloom filters (technique 2) reduce?



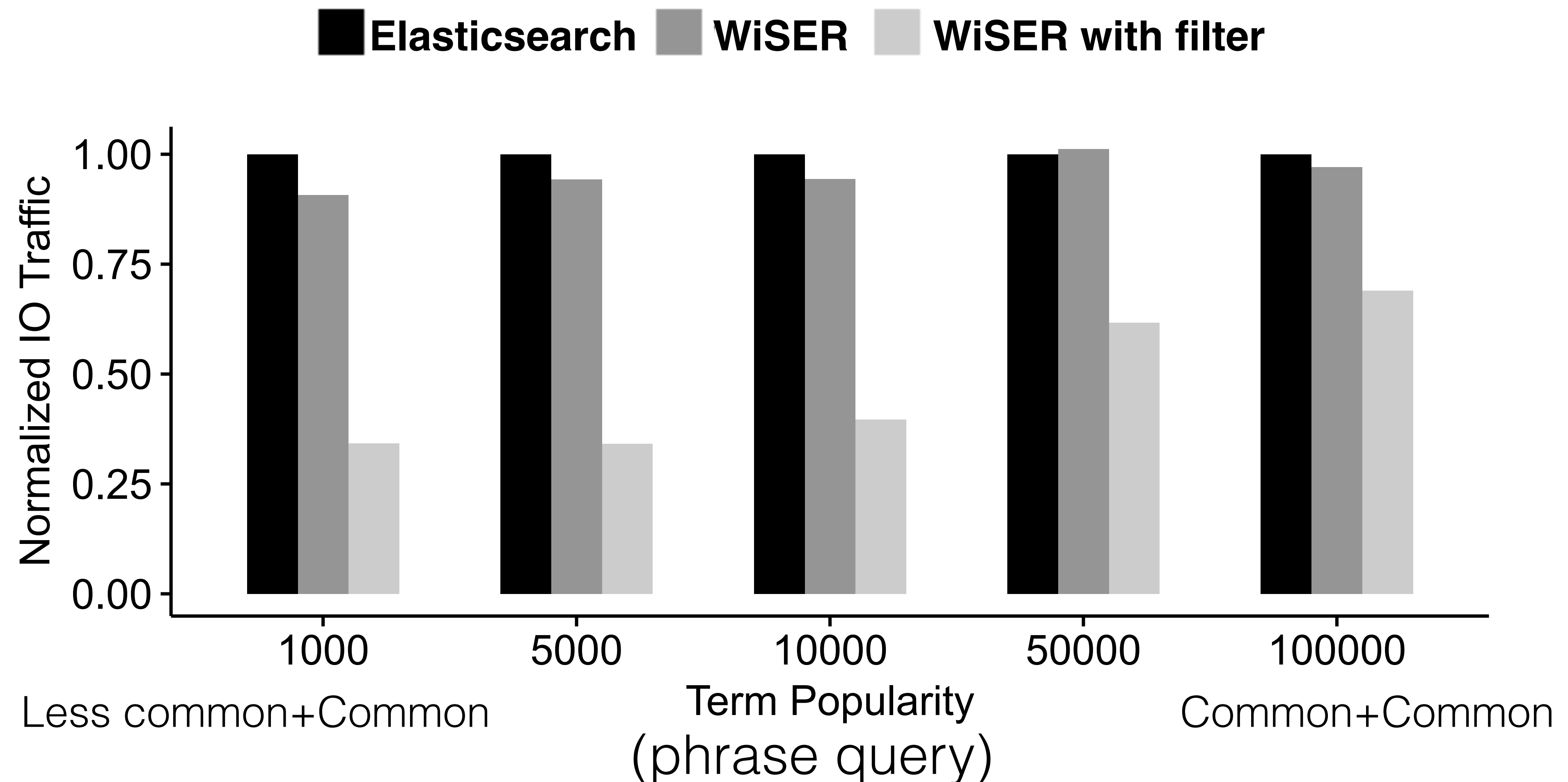
# How much read traffic can two-way cost-aware Bloom filters (technique 2) reduce?



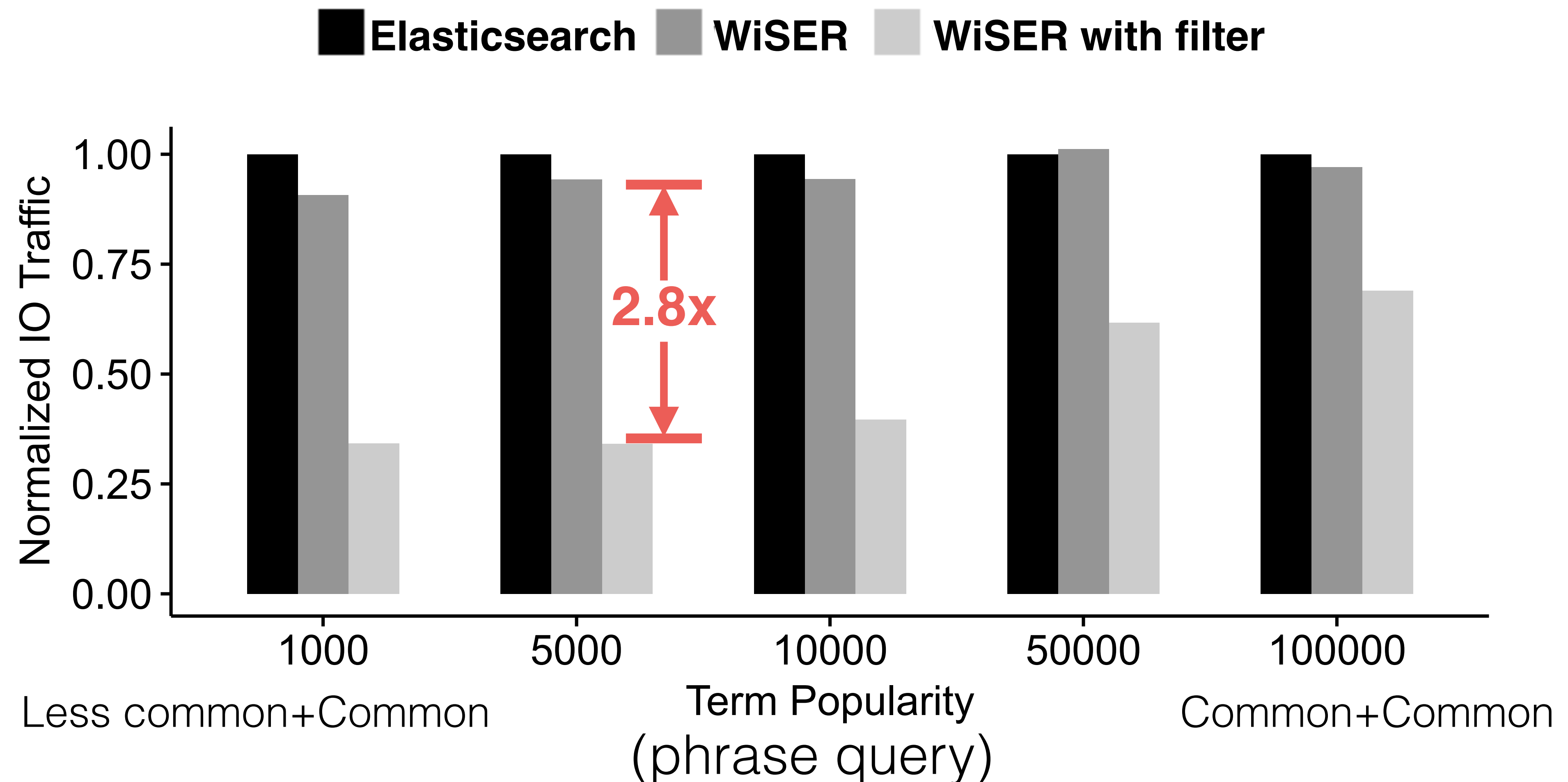
# How much read traffic can two-way cost-aware Bloom filters (technique 2) reduce?



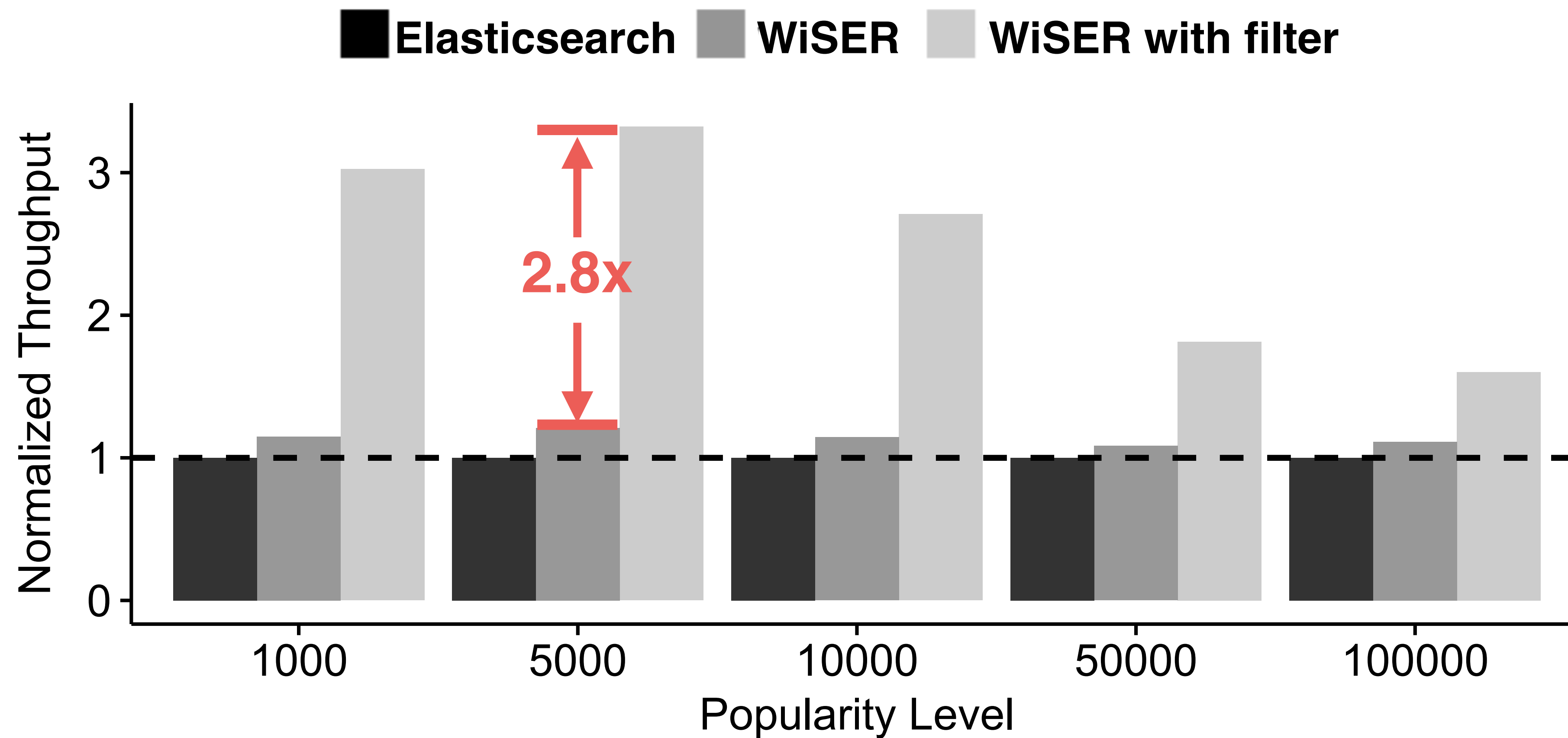
# How much read traffic can two-way cost-aware Bloom filters (technique 2) reduce?



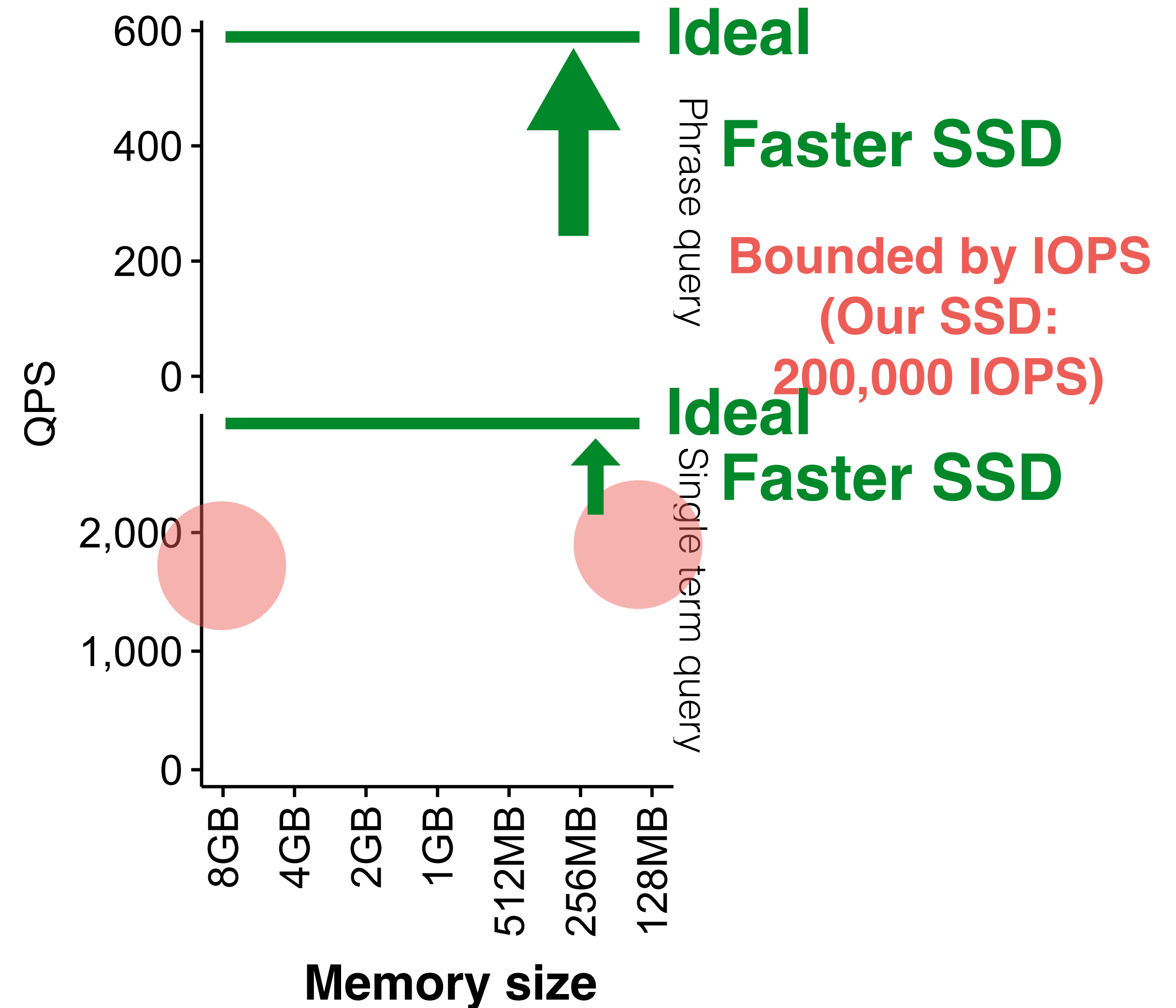
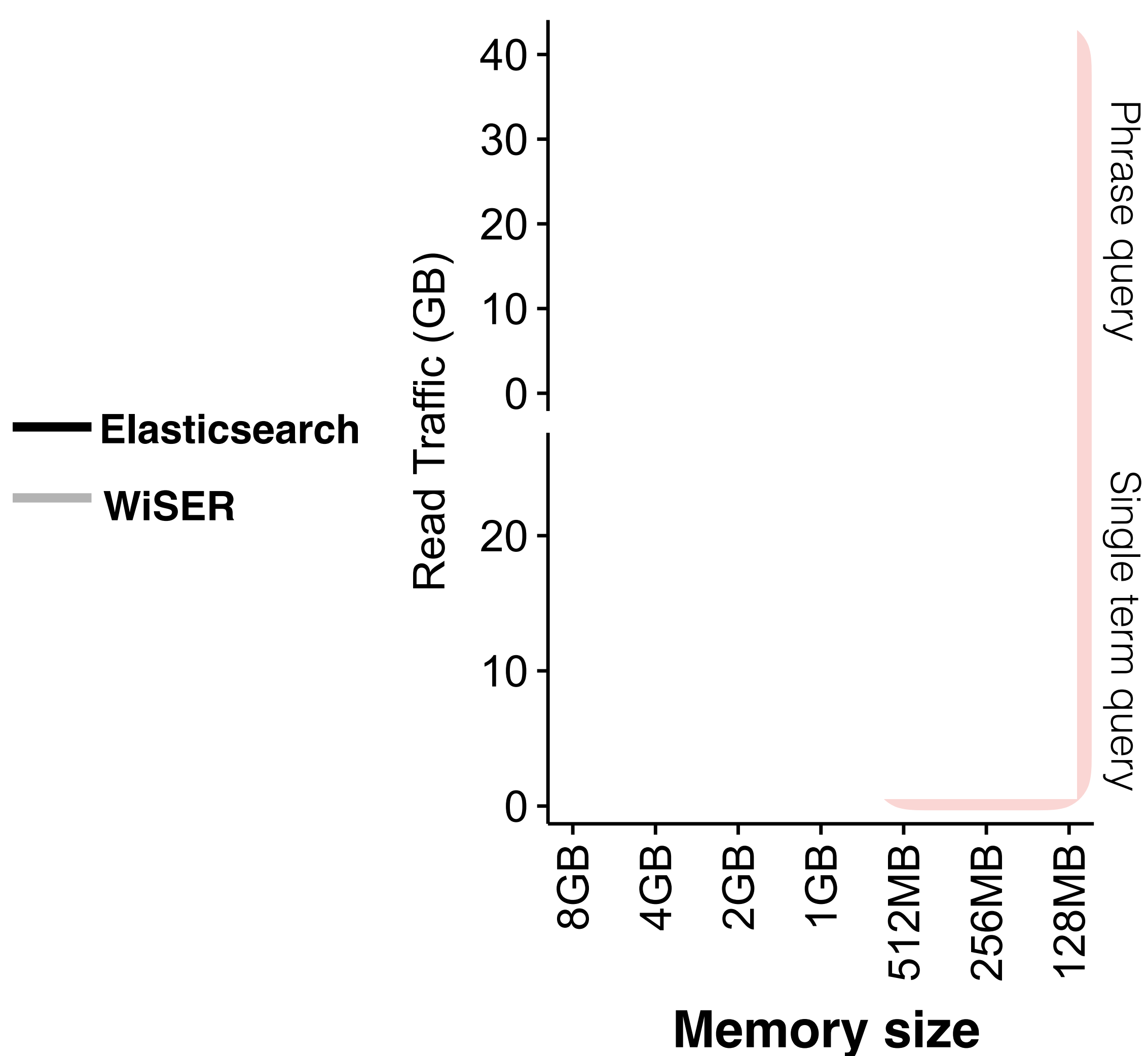
# How much read traffic can two-way cost-aware Bloom filters (technique 2) reduce?



# How much query throughput can two-way cost-aware filtering (technique 2) increase?



# How does small-memory performance compare with larger-memory performance?





**More results are in the paper...**

# Final Thoughts

**Does your application really need large cache/RAM?**

**Will your application works just fine if it reads data as needed from fast SSDs?**

**Q & A**  
**Ask as Needed**