

# Zwift : A Programming Framework for High Performance Text Analytics on Compressed Data

Feng Zhang †◊, Jidong Zhai ◊, Xipeng Shen #, Onur Mutlu ★, Wenguang Chen ◊

†Renmin University of China

◊Tsinghua University

#North Carolina State University

★ETH Zürich



**ETH** Zürich

# Motivation

- Data management evolution in history.



Every day, 2.5 quintillion bytes of data created – 90% data in the world today has been created in the last two years alone[1].

- QUESTION:

**What if the data are too big  
that exceed storage capacity?**

[1] What is Big Data?

<https://www-01.ibm.com/software/data/bigdata/what-is-big-data.html>

# Motivation

- Challenge 1:
  - SPACE: Large Space Requirement
- Challenge 2:
  - TIME: Long Processing Time
- Observation:
  - Using Hash Table to check redundant content

REUSE



# Motivation

- Our Idea:
  - Text analytics directly on compressed data (TADOC)

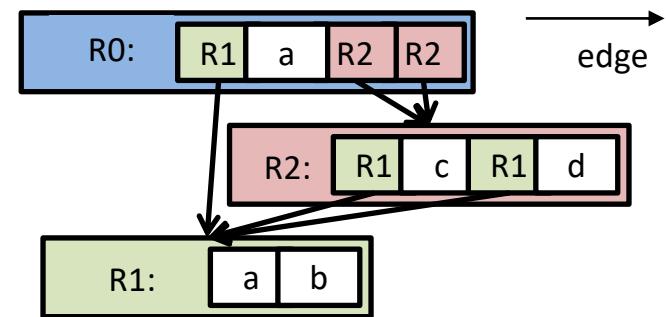
Input:

```
a b a a b c a b d a b c  
a b d
```

Rules:

```
R0 → R1 a R2 R2  
R1 → a b  
R2 → R1 c R1 d
```

(a) Original data

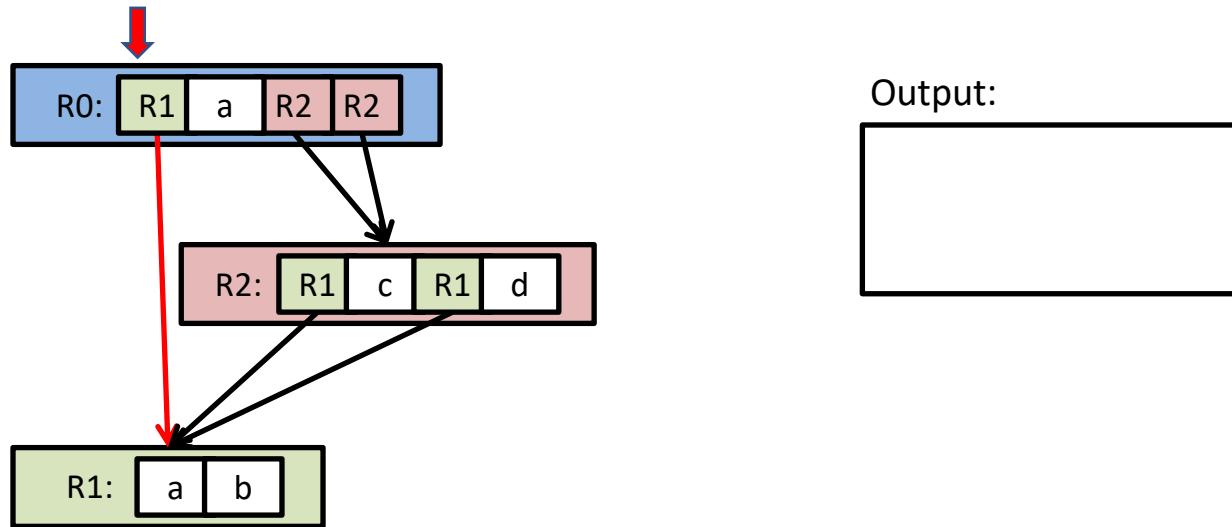


(b) Grammar-based compression

(c) DAG Representation

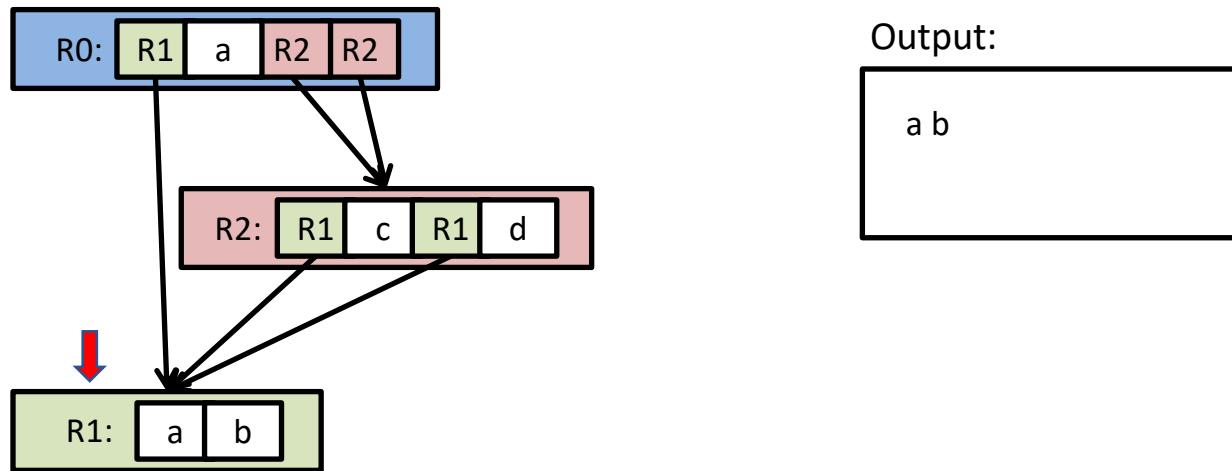
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



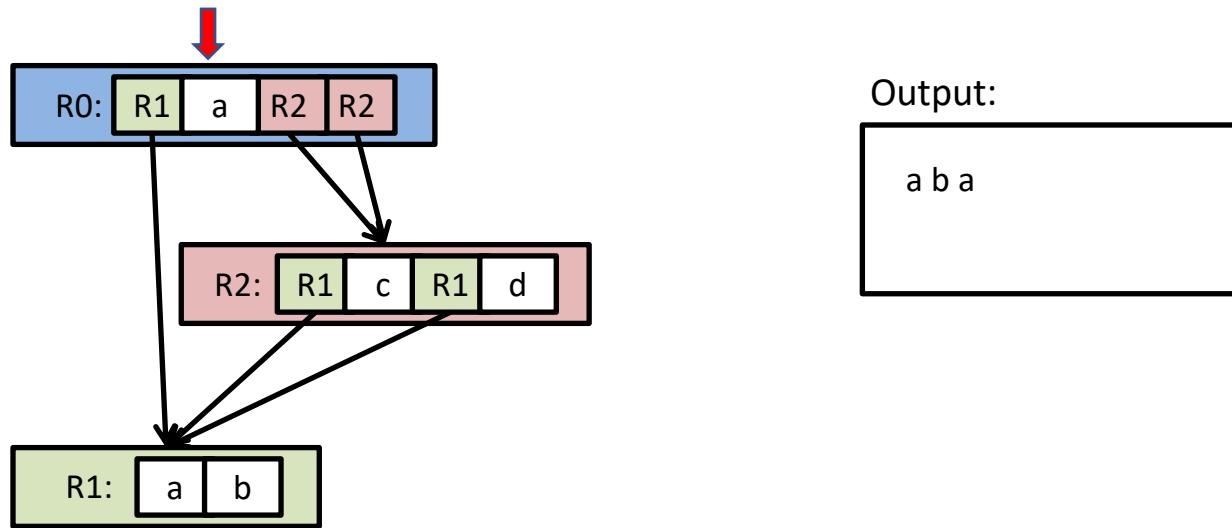
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



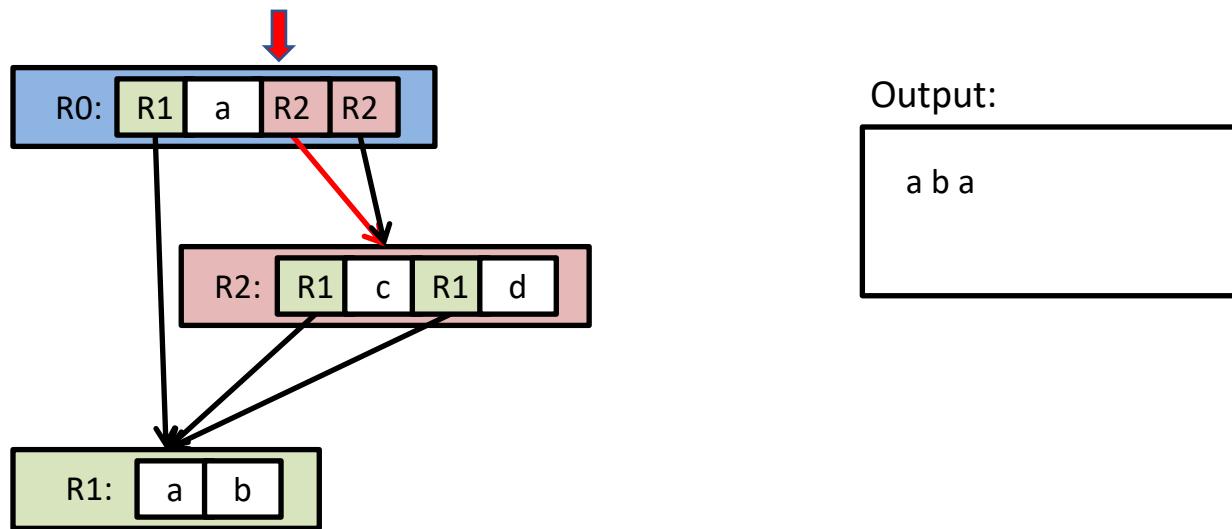
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



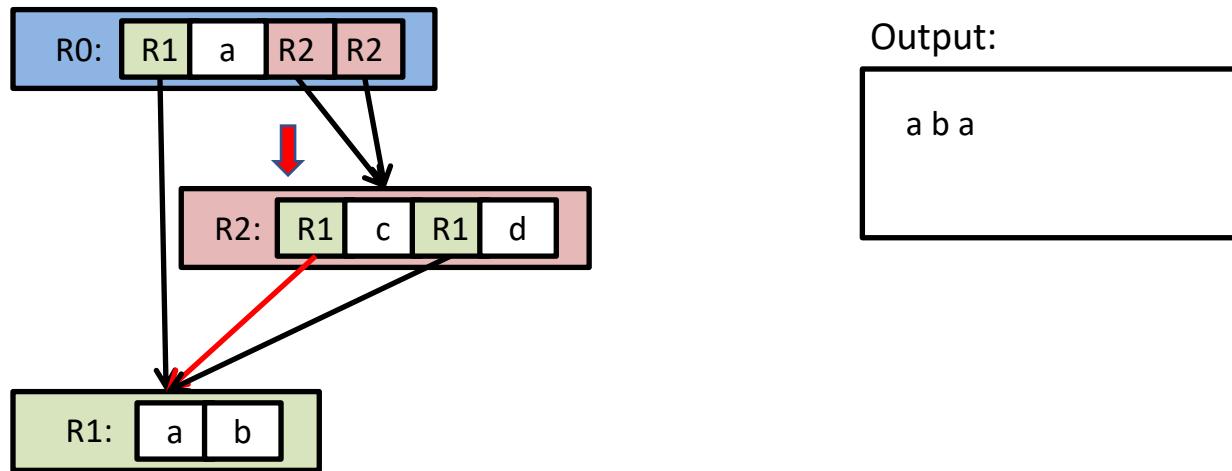
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



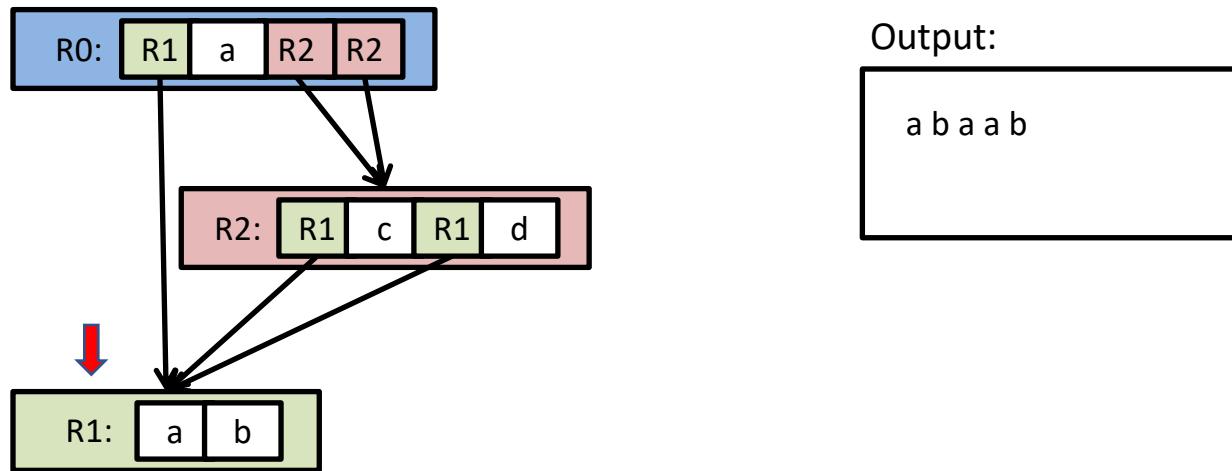
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



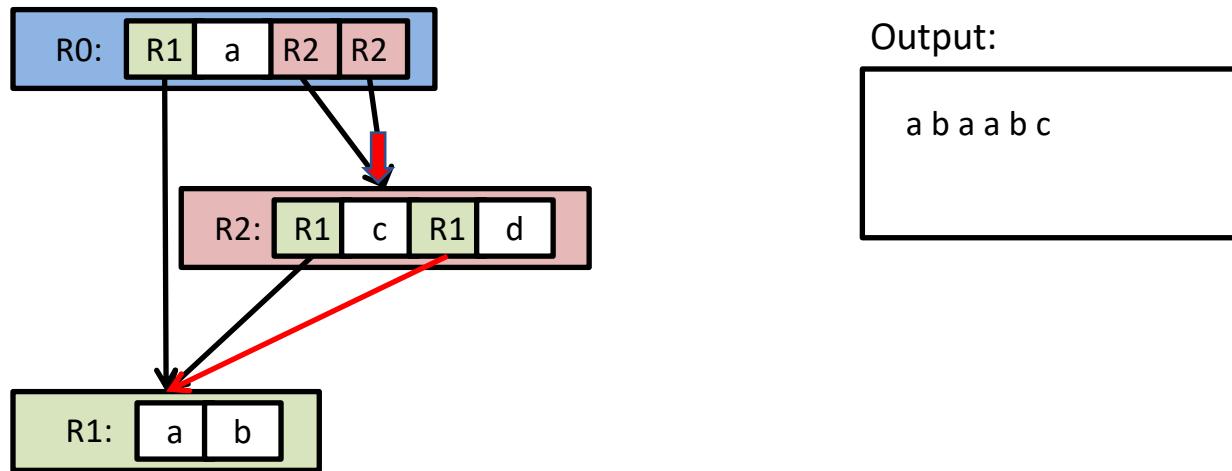
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



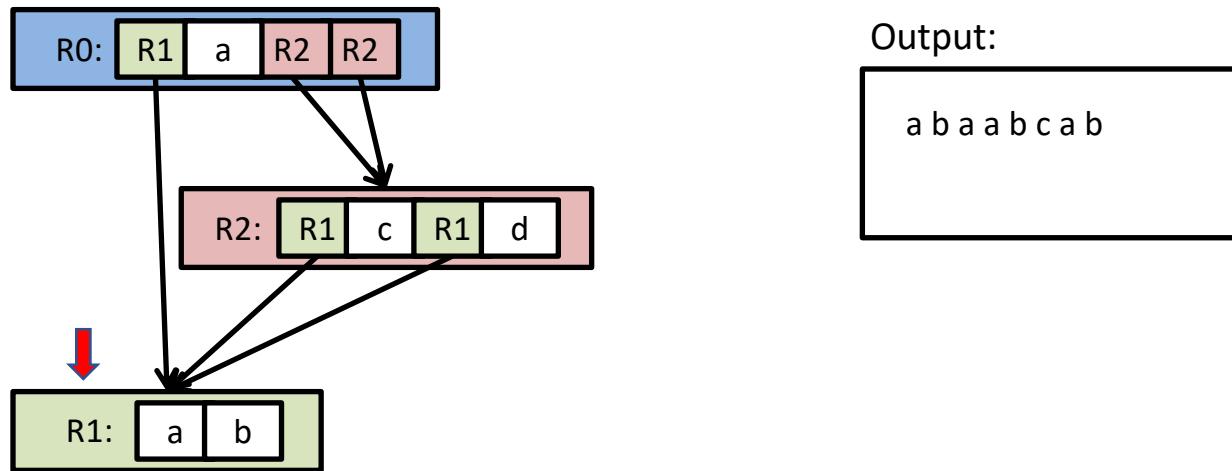
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



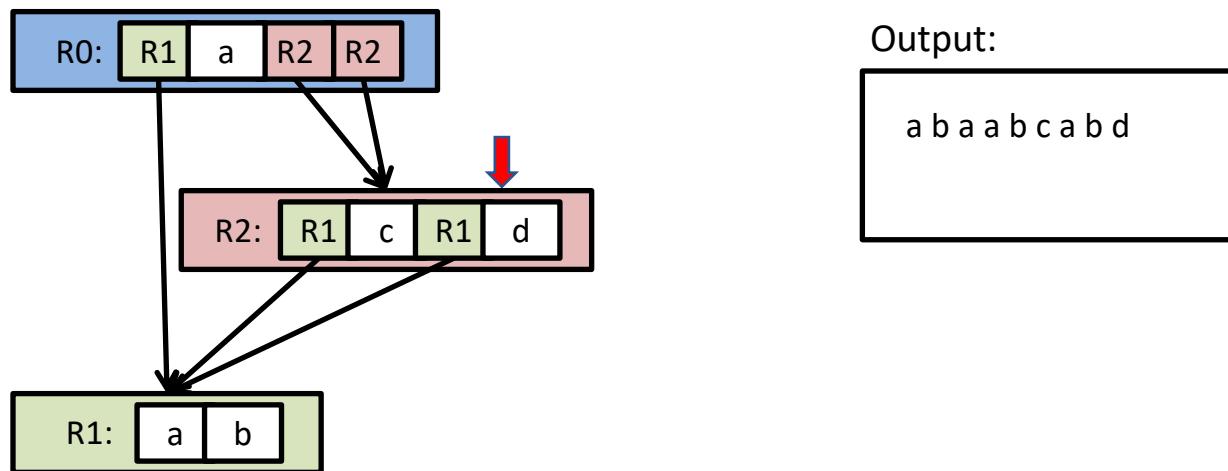
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



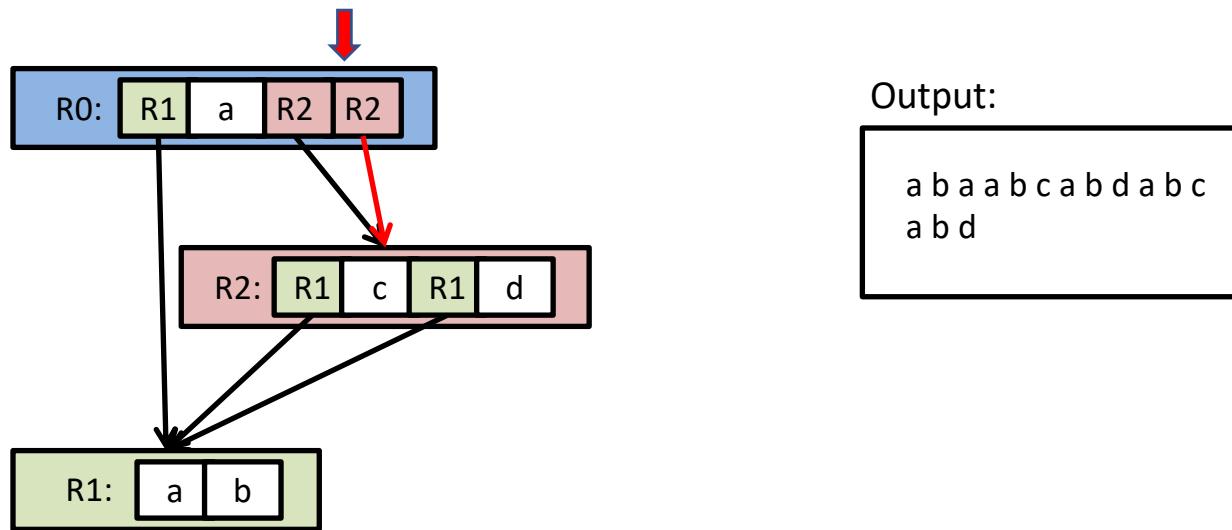
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



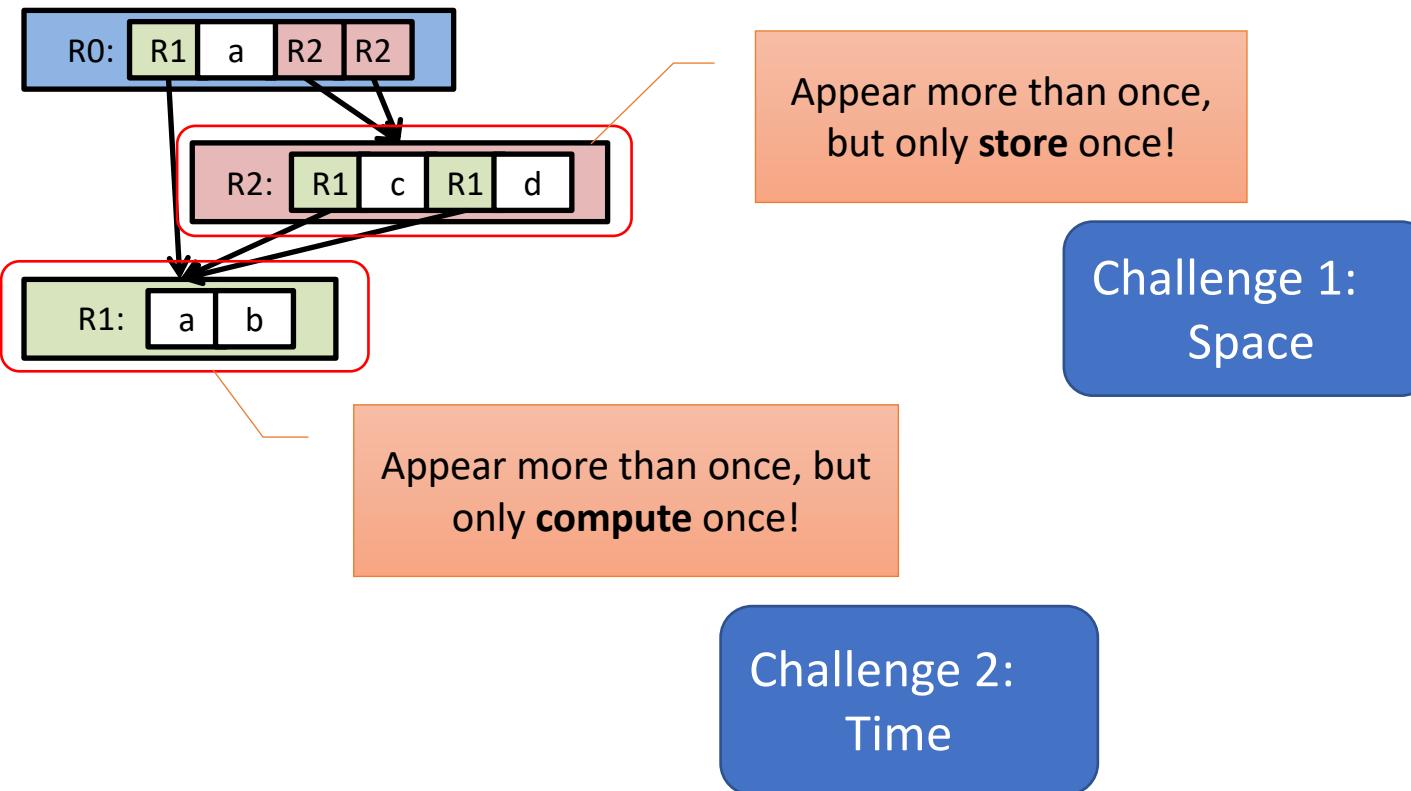
# Motivation

- How text analytics directly on compressed data (TADOC) recovers data?



# Motivation

- Text analytics directly on compressed data (TADOC), why it is good?
- In format.

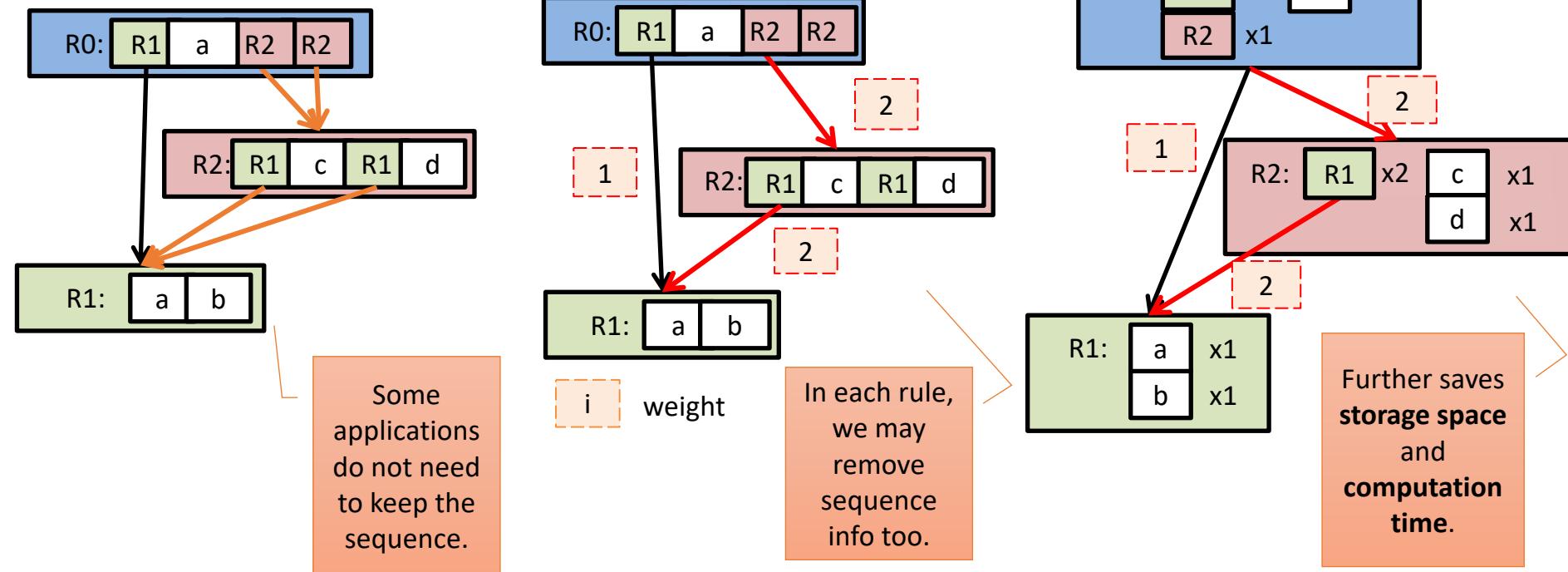


What's more?

Even BETTER!

# Motivation

- Text analytics directly on compressed data (TADOC), why it is good?
- In Process:

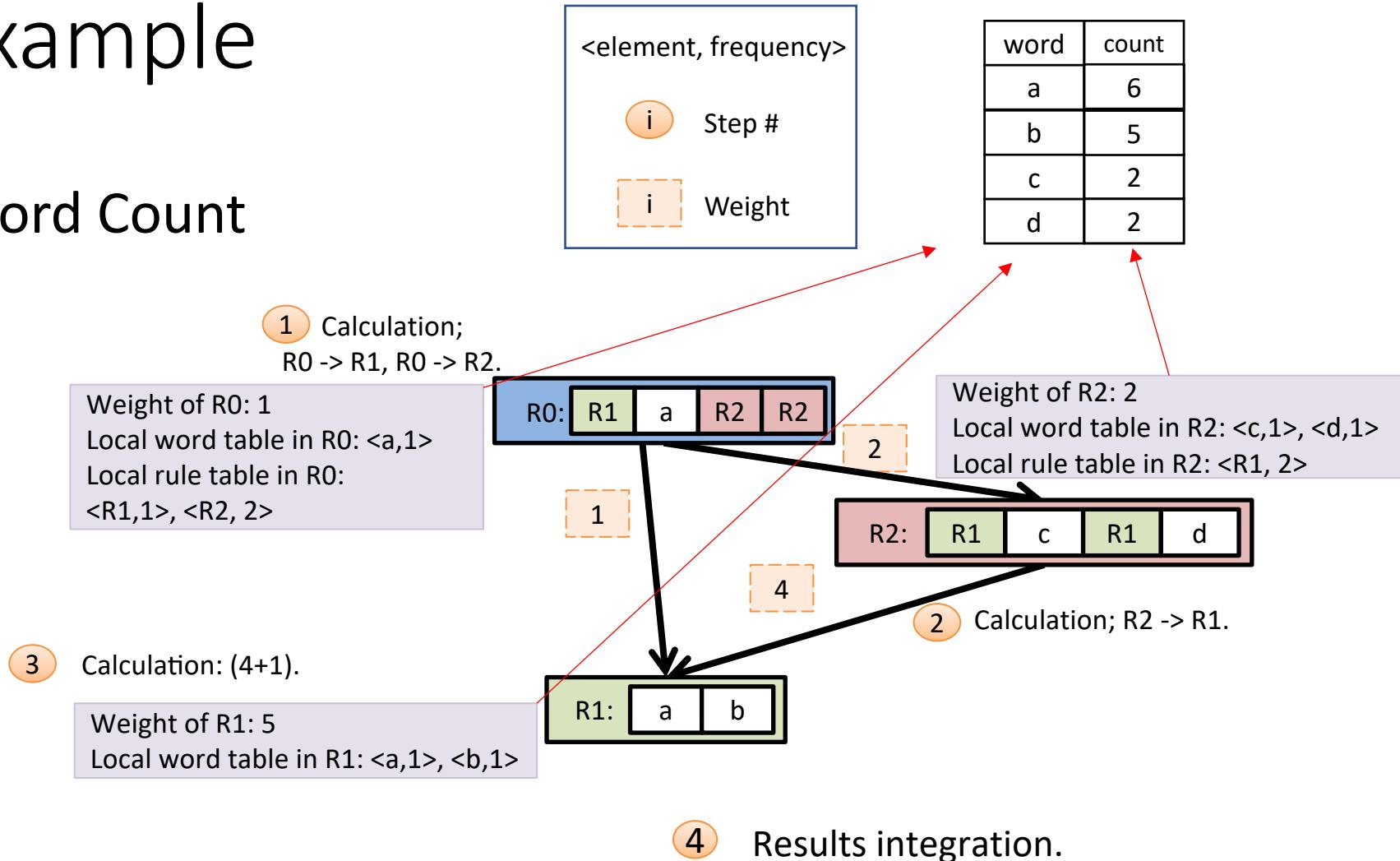


# Outline

- Introduction
    - Motivation & Example
    - Programming Challenges
    - Solution
  - Swift Implementation
    - Language Overview
    - Language Features
  - Results
    - Benchmarks
    - Performance
    - Storage Savings & Productivity
  - Conclusion
- 
- The diagram illustrates the flow of information in the outline. A blue line connects the 'Introduction' section to a blue box labeled 'We are here!'. From the 'Motivation & Example' item, a green line points to a green box labeled 'Use word count for illustration'. From the 'Programming Challenges' item, a yellow line points to an orange box labeled 'Not easy to generalize this technique'. A horizontal line connects the 'Solution' item to a white box labeled 'A conceptual framework and a high level solution'. An orange line connects the 'Swift Implementation' section to a white box labeled 'How we implement our idea'. A blue line connects the 'Results' section to a white box labeled 'Exhibition of the Swift from several perspectives'.

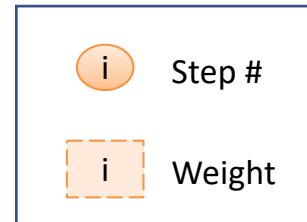
# Example

- Word Count



# Example

- Word Count

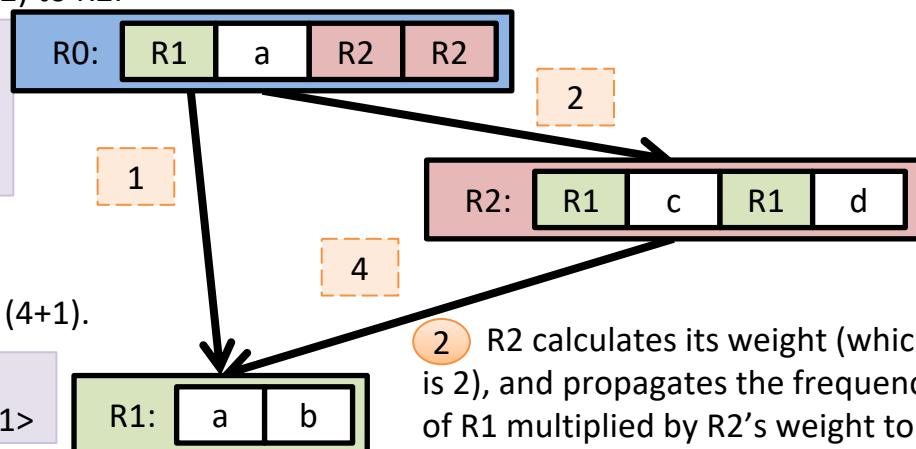


- 1 R0 propagates the frequency of R1 (which is 1) to R1, and the frequency of R2 (which is 2) to R2.

Weight of R0: 1

Local word table in R0: <a,1>

Local rule table in R0: <R1,1>, <R2, 2>



- 3 R1 calculates its weight, which is 5 (4+1).

Weight of R1: 5

Local word table in R1: <a,1>, <b,1>

- 4 We integrate the local word table in each node multiplied by its weight as the final result.

2 R2 calculates its weight (which is 2), and propagates the frequency of R1 multiplied by R2's weight to R1.

Weight of R2: 2

Local word table in R2: <c,1>, <d,1>

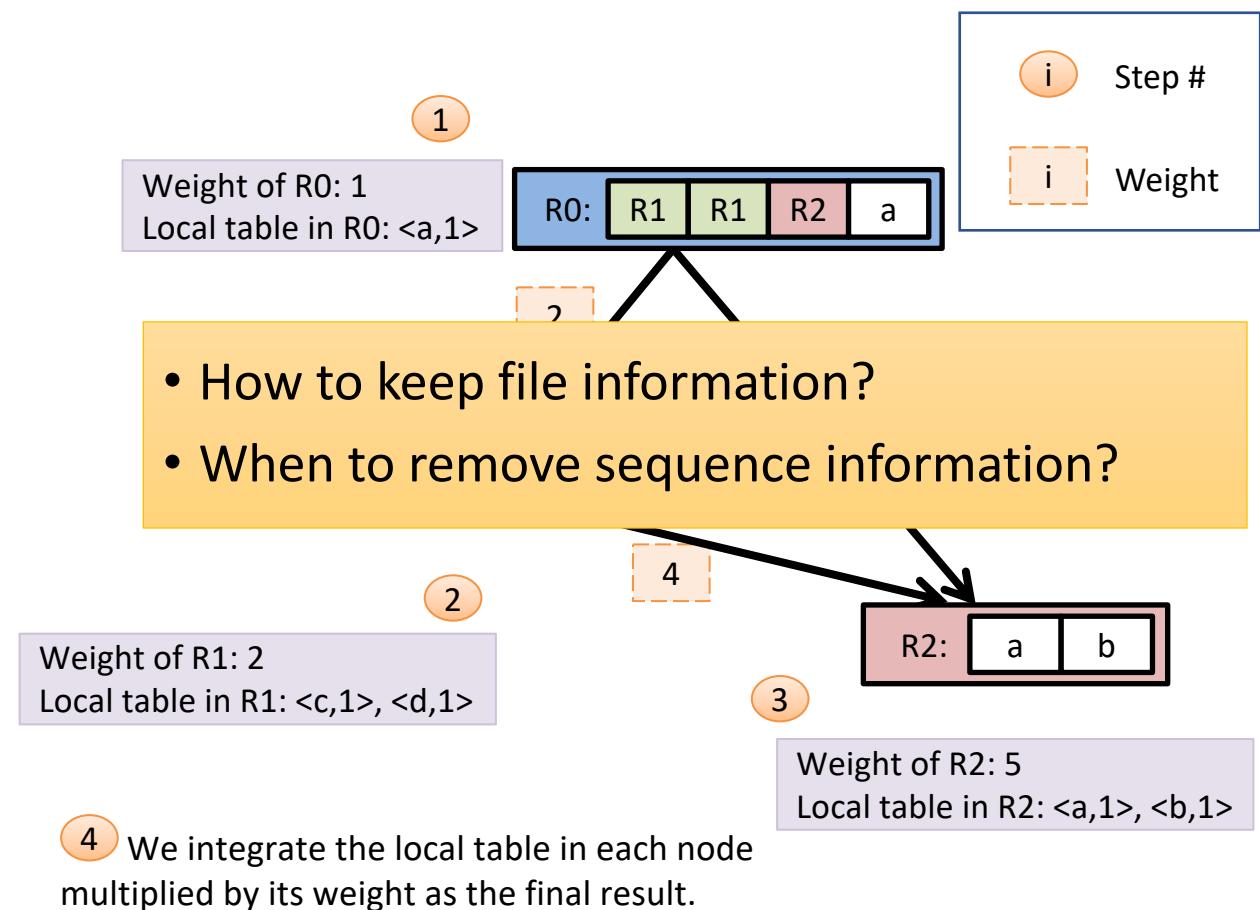
Local rule table in R2: <R1, 2>

# Programming Challenges

Problem dimension

Implementation dimension

Dataset dimension

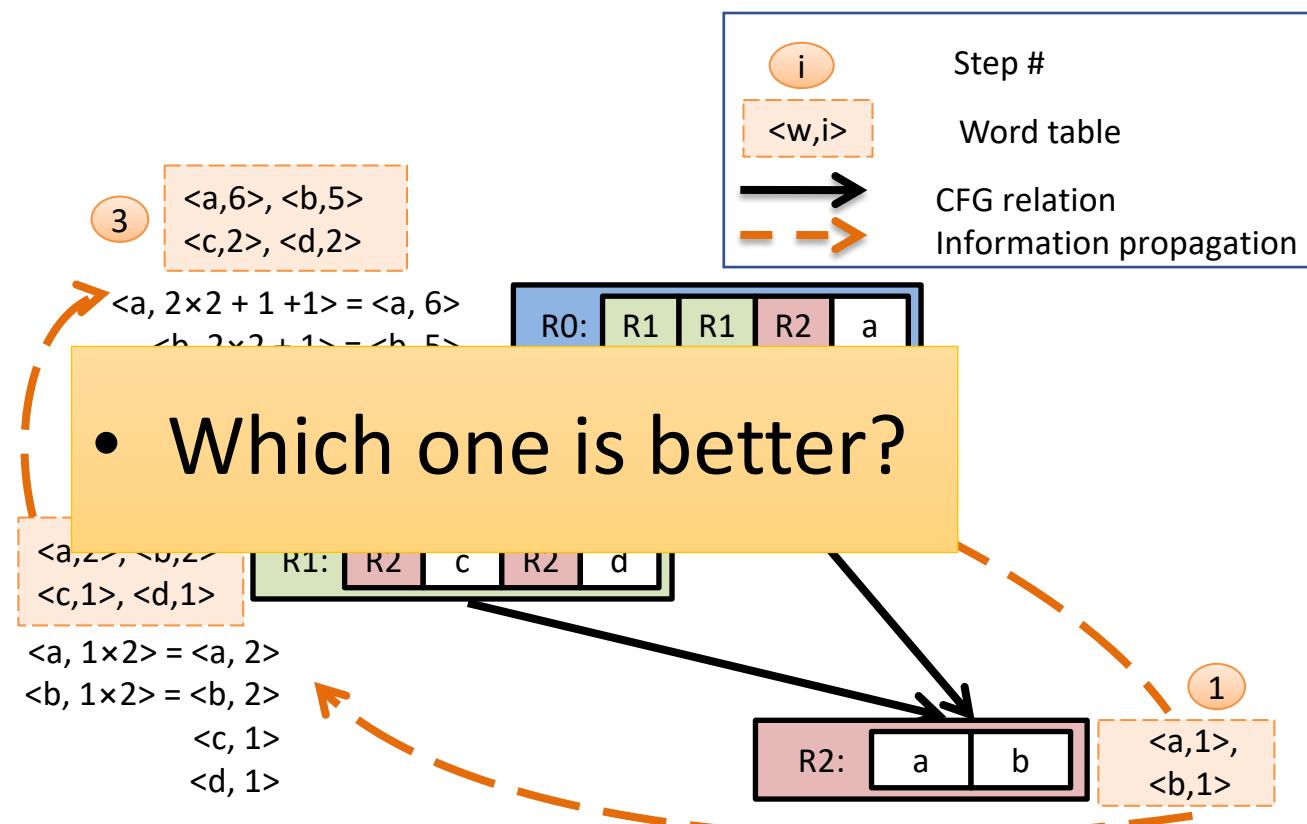


# Programming Challenges

Problem dimension

Implementation dimension

Dataset dimension

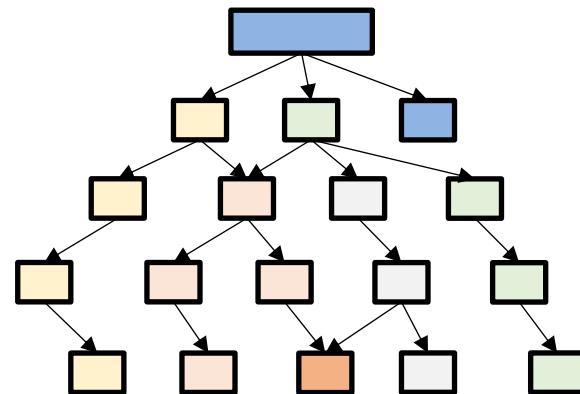


# Programming Challenges

Problem  
dimension

Implementation  
dimension

Dataset  
dimension

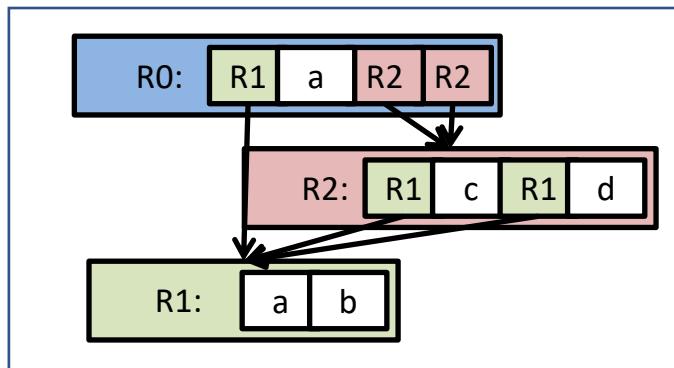


Traversal Order  
↓ or ↑ or ?

The best traversal order may  
depend on input.

# Solution

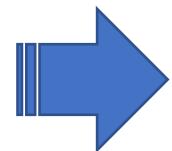
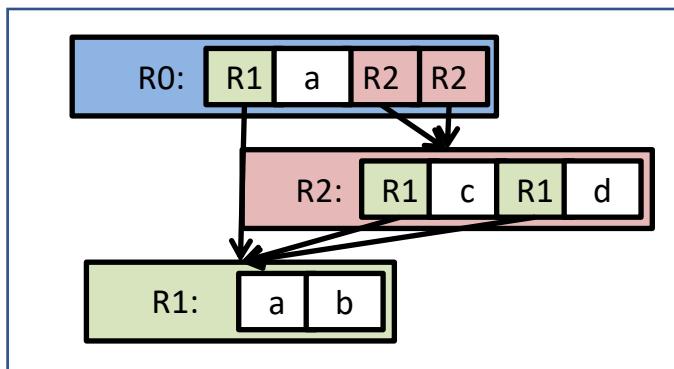
- A conceptual framework, a six-element tuple  $(G, V, F, V, D, \Lambda)$ 
  - A graph  $G$ , DAG representation



- G, how to represent it.

# Solution

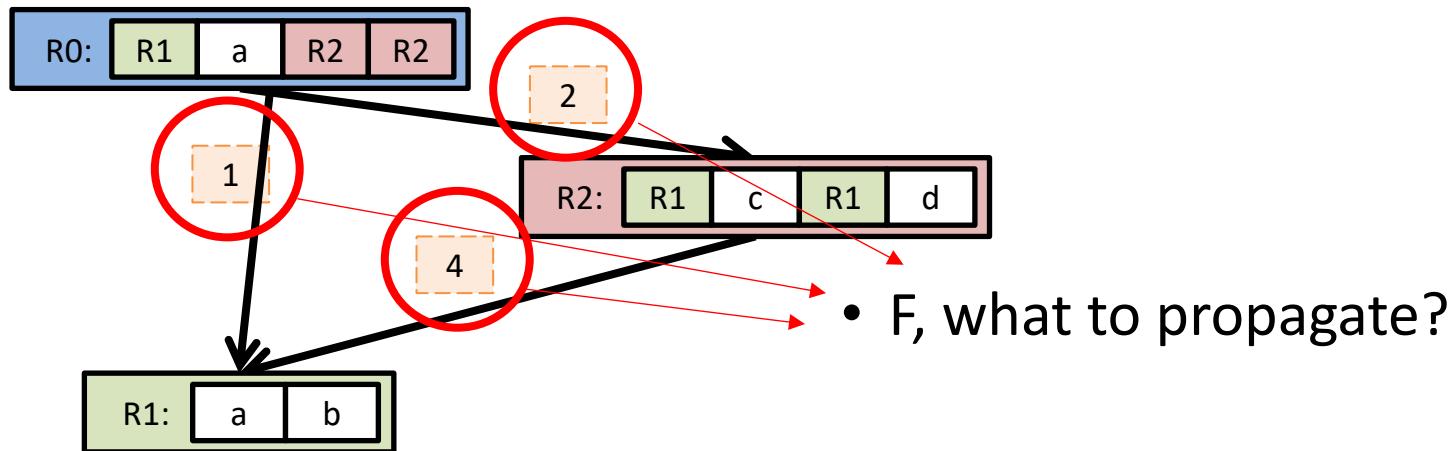
- A conceptual framework, a six-element tuple  $(G, V, F, V, D, \Lambda)$ 
  - A domain of values  $V$ , possible values as the outputs



- $V$ , what do we want to get from  $G$ ?

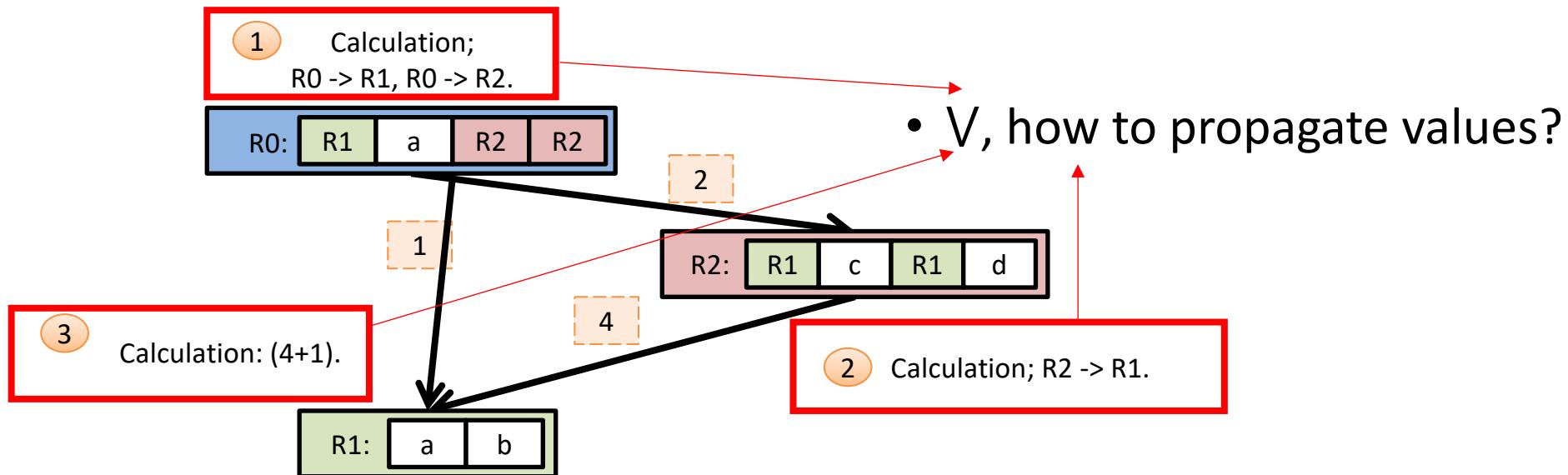
# Solution

- A conceptual framework, a six-element tuple  $(G, V, F, \nabla, D, \Lambda)$ 
  - A domain of values  $F$ , possible values to be propagated



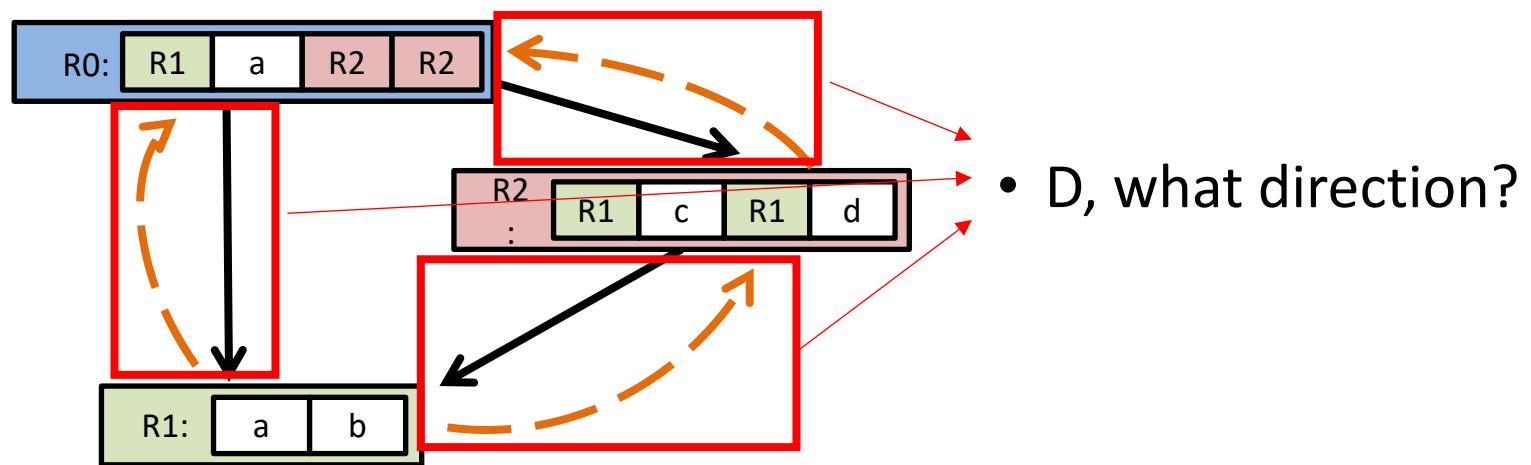
# Solution

- A conceptual framework, a six-element tuple  $(G, V, F, V, D, \Lambda)$ 
  - A meet operator  $V$ , how to transform values



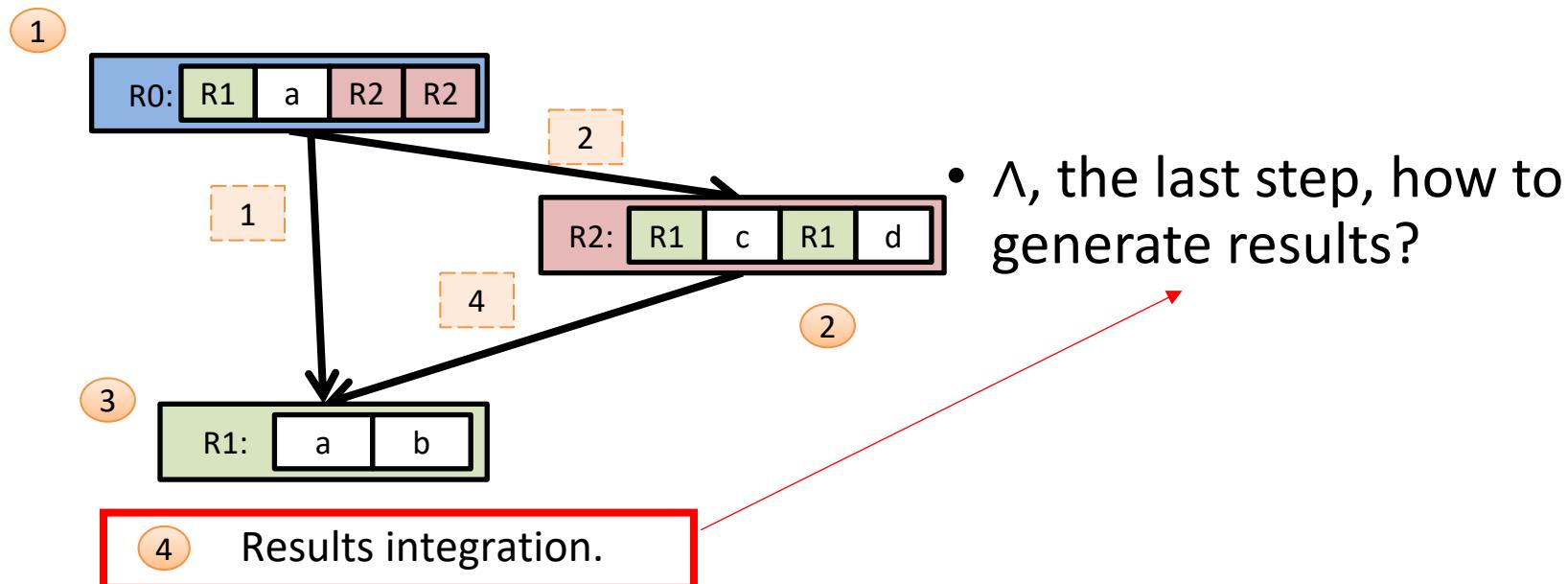
# Solution

- A conceptual framework, a six-element tuple  $(G, V, F, V, D, \Lambda)$ 
  - A direction of the value flow  $D$



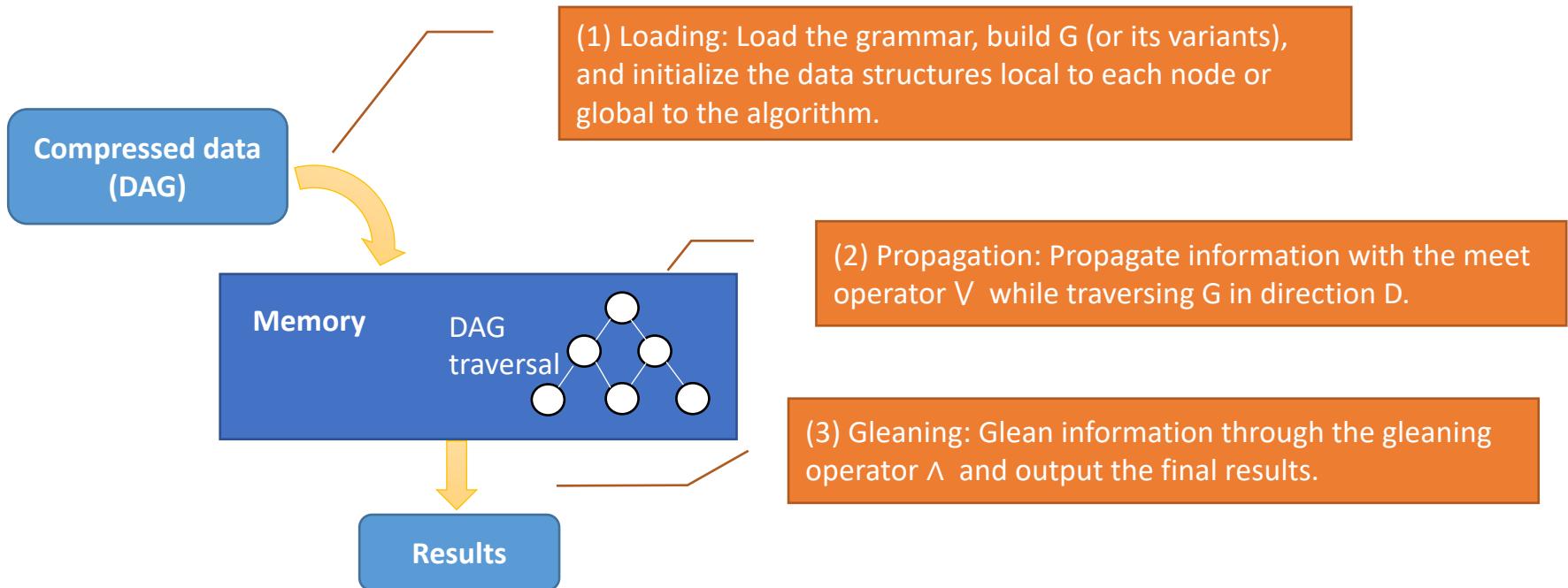
# Solution

- A conceptual framework, a six-element tuple  $(G, V, F, V, D, \Lambda)$ 
  - A gleaning operator  $\Lambda$ , final stage of the analytics



# Solution

- High-level algorithm for TADOC:



# Contribution

- How to generalize the idea of TADOC
- Swift, the first programming framework for TADOC
  - The Swift Language
  - A compiler and runtime
  - A utility library

# Outline

- Introduction
  - Motivation & Example
  - Programming Challenges
  - Solution
- **Zwift Implementation**
  - **Language Overview**
  - **Language Features**
- Results
  - Benchmarks
  - Performance
  - Storage Savings & Productivity
- Conclusion

# The Zwift Language

Zwift DSL template.

- The Conceptual Framework  
 $(G, V, F, \vee, D, \Lambda)$

- A graph G.
- A domain of values V.
- A domain of values F.
- A meet operator  $\vee$ .
- A direction of the value flow D.
- A gleaning operator  $\Lambda$ .

```
01 ELEMENT = LETTER/WORD/SENTENCE
02 USING_FILE = true/false
03 NodeStructure = {
04 //data structures in each node
05 }
06 Init = {
07 //initializations for each node in ZwiftDAG
08 }
09 Direction = FORWARD/BACKWARD/DEPTHFIRST
10 Action = {
11 //actions taken at each node during a traversal
12 }
13 Result = {
14 // result structure
15 }
16 FinalStage = {
17 // final operations to produce the results
18 }
```

# Language Features

- Edge Merging
  - Edge merging merges multiple edges between two nodes into one, and the weight of the edge may be used to indicate the number of original edges in the Swift code.
- Node Coarsening
  - Node coarsening reduces the size of the graph, and at the same time, reduces the number of substrings spanning across nodes.
- Two-Level Bit Vector
  - Level Two contains a number of N -bit vectors (where N is a configurable parameter). Level One contains a pointer array and a level-1 bit vector.

# Language Features

- Coarse-Grained Parallelism
  - The parallelism is at the data level. The coarse-grained method also simplifies the conversion from sequential code to parallel and distributed versions of the code.
- Version Selection
  - Swift allows programmers to easily run the different versions on some representative inputs.
- Double Compression
  - Double compression performs further compression (e.g., using Gzip) of the compressed results from Sequitur.

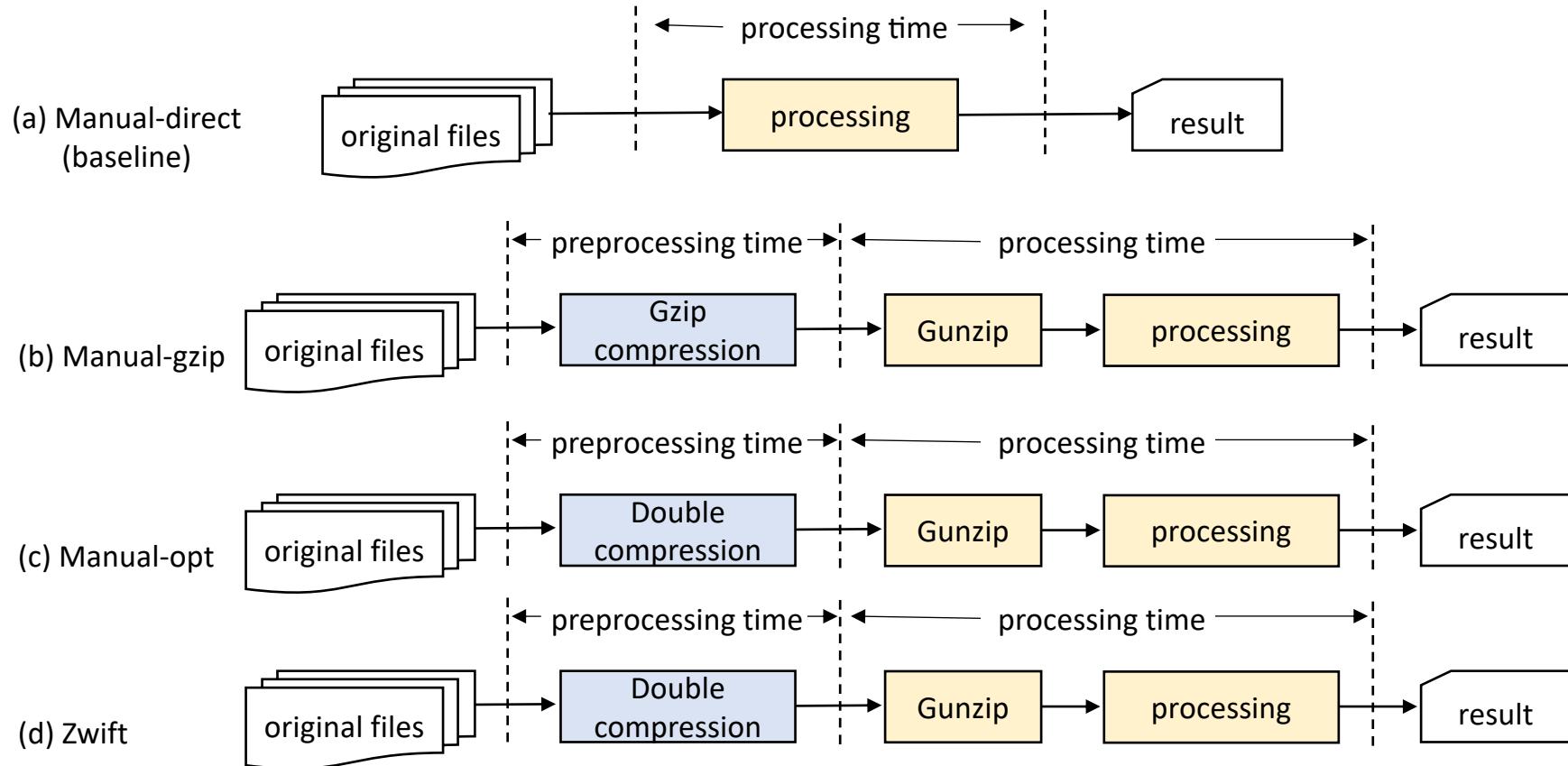
# Outline

- Introduction
  - Motivation & Example
  - Programming Challenges
  - Solution
- Swift Implementation
  - Language Overview
  - Language Features
- **Results**
  - **Benchmarks**
  - **Performance**
  - **Storage Savings & Productivity**
- Conclusion

# Benchmarks

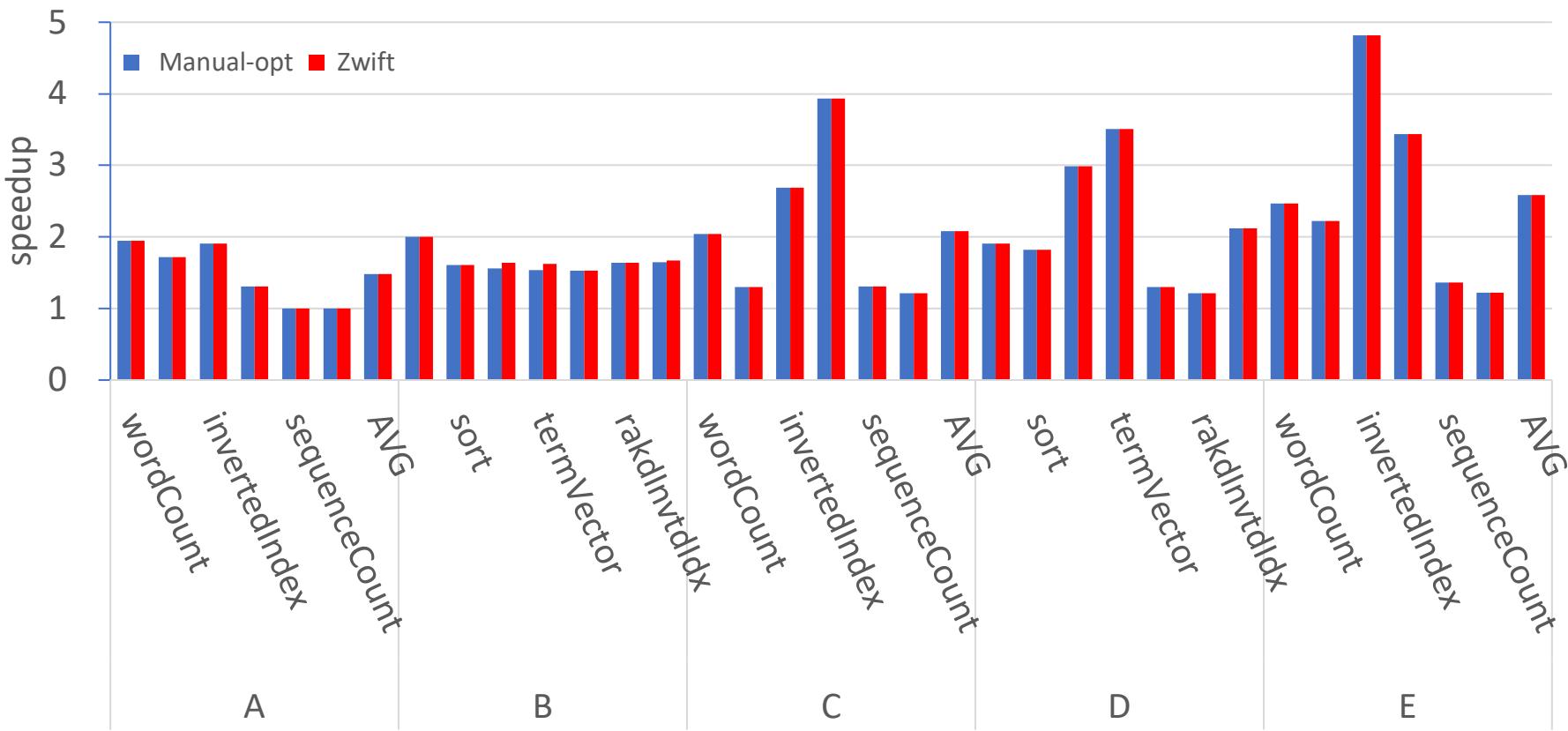
- Six benchmarks
  - Word Count, Inverted Index, Sequence Count, Ranked Inverted Index, Sort, Term Vector
- Five datasets
  - 580 MB ~ 300 GB
- Two platforms
  - Single node
  - Spark cluster (10 nodes on Amazon EC2)

# Benchmarks



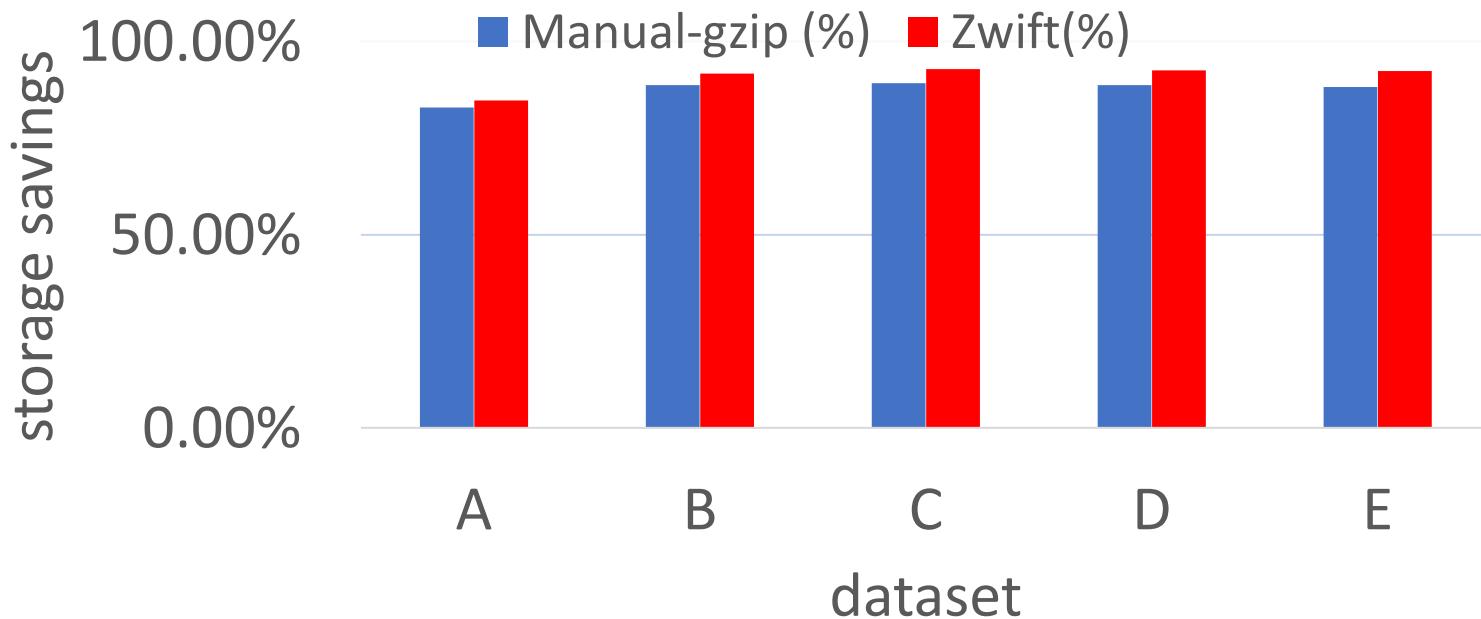
# Performance

- Execution time benefits of speedup over *manual-direct*.
  - Zwift yields **2X** speedup, on average, over *manual-direct*.
  - Zwift and *manual-opt* show similar performance, which indicates that Zwift successfully unleashes most power of TADOC.



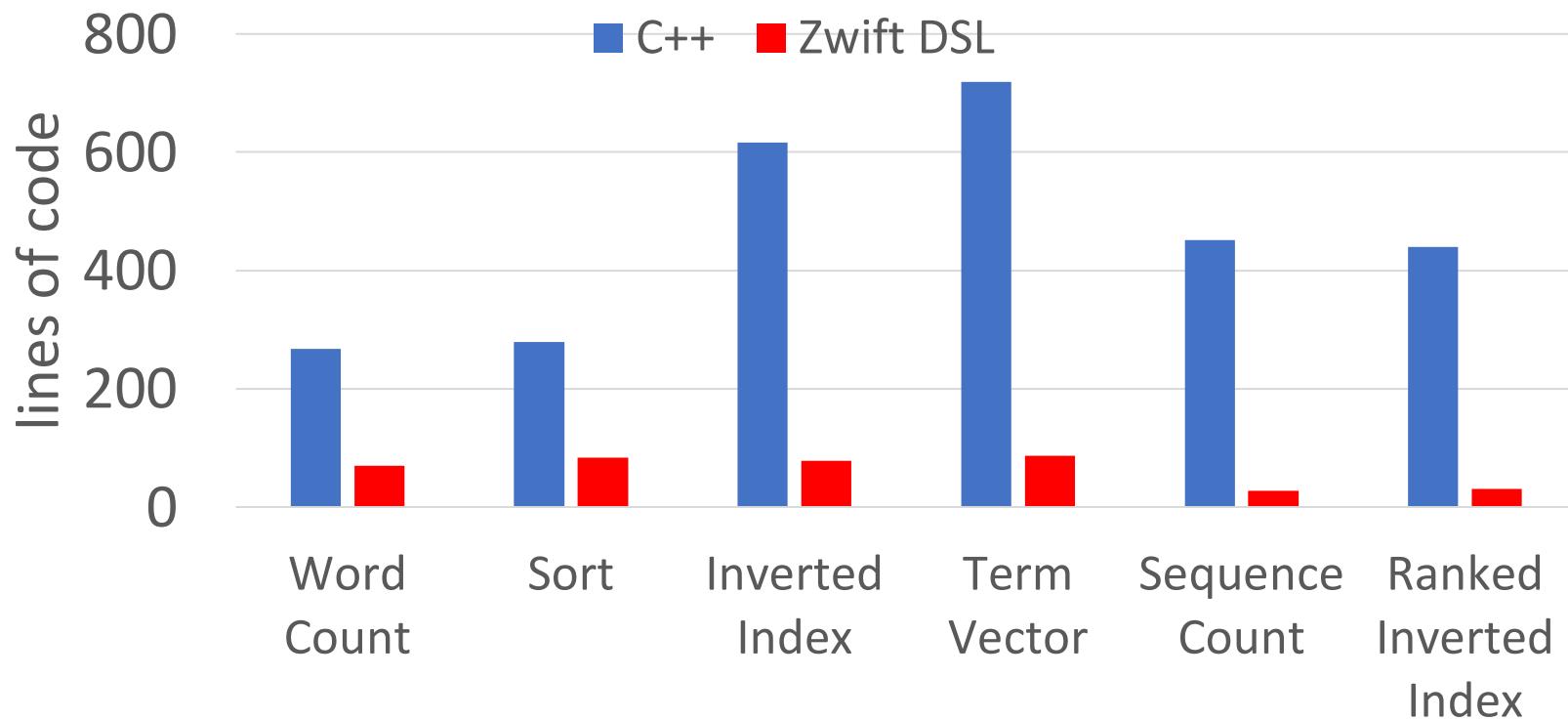
# Storage Savings

Zwift reduces storage usage by **90.8%**, even more than *manual-gzip* does.



# Productivity

On average, the Swift code is **84.3%** shorter than the equivalent C++ code.



# Conclusion

- Swift is an effective framework for TADOC.
  - It reduces storage usage by **90.8%**
  - It reduces execution time by **41.0%**
- Swift significantly improves programming productivity.

# Thanks!

- Any questions?