

# COMP9313: Big Data Management

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

**Lecturer: Xin Cao**

Course web site: <http://www.cse.unsw.edu.au/~cs9313/>

Assignment Project Exam Help

**Set Si** <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Set-Similarity Join

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Finding pairs of records with a **similarity** on their join attributes > t

# Application: Record linkage

Table R

Star
Keanu Reeves
Samuel Jackson
Schwarzenegger
...

Table S

Star
Reeves
Jackson
negger
...



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Two-step Solution

Table R

Star
...

**Step 1:**  
**Similarity Join**



Table S

Star
...

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

**Step 2: Verification**

# Why Hadoop?

- Large amounts of data
- Data or processing does not fit in one machine

- Assumptions:

- Self join  $R \bowtie R$
  - Two similar

<https://eduassistpro.github.io/>

- Efficient Parallel Set-Similarity Join p (SIGMOD'10)

Add WeChat edu\_assist\_pro

## A naïve solution

- Map:  $\langle 23, (a,b,c) \rangle \rightarrow (a, 23), (b, 23), (c, 23)$
- Reduce:  $(a,23), (a,29), (a,50), \dots \rightarrow$  Verify each pair  $(23, 29), (23, 50), (29, 50) \dots$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

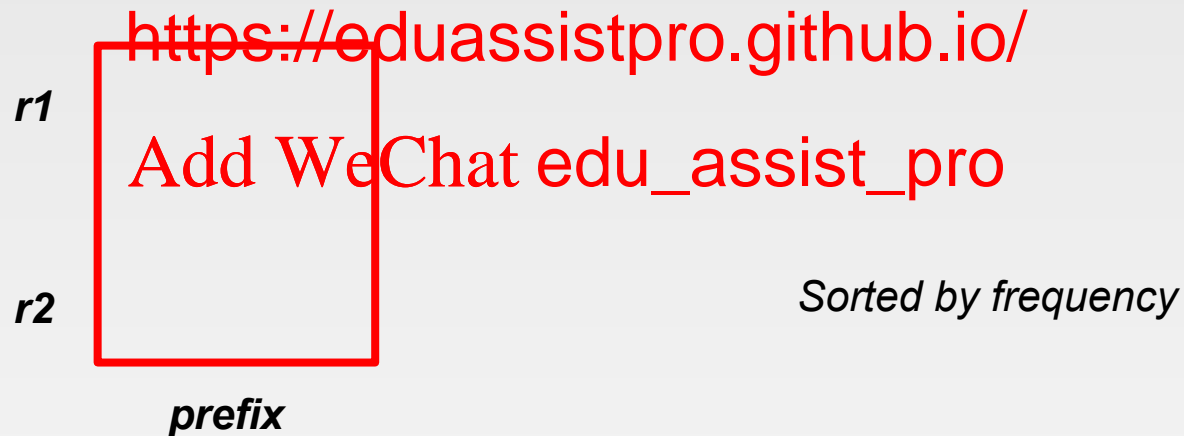
- Too much data to transfer 😞
- Too many pairs to verify 😞.

# Solving frequency skew: prefix filtering

- Sort tokens by frequency (ascending)



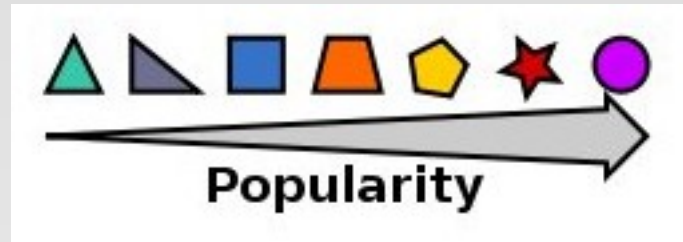
- **Prefix** of a set: least frequent tokens



- Prefixes of similar sets should share tokens



# Prefix filtering: example



Record 1 Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Record 2 Add WeChat edu\_assist\_pro

- Each set has 5 tokens
- “Similar”: they share at least 4 tokens
- Prefix length: 2

# Hadoop Solution: Overview

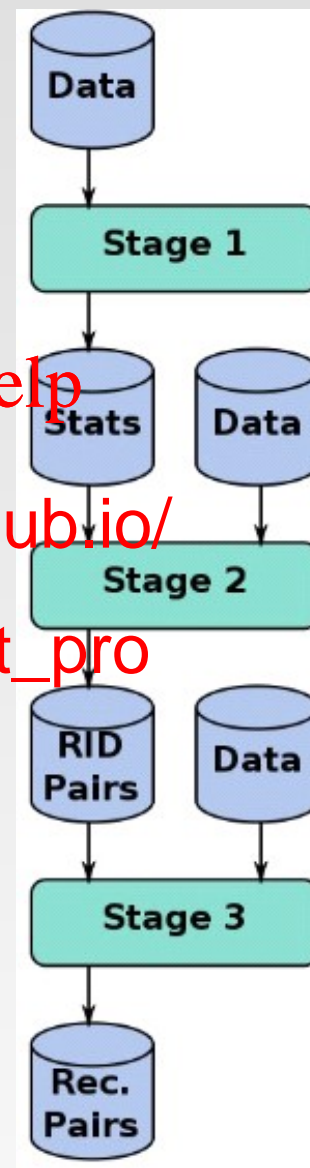
- Stage 1: Order tokens by frequency  
(Already done in the given example data)

Assignment Project Exam Help

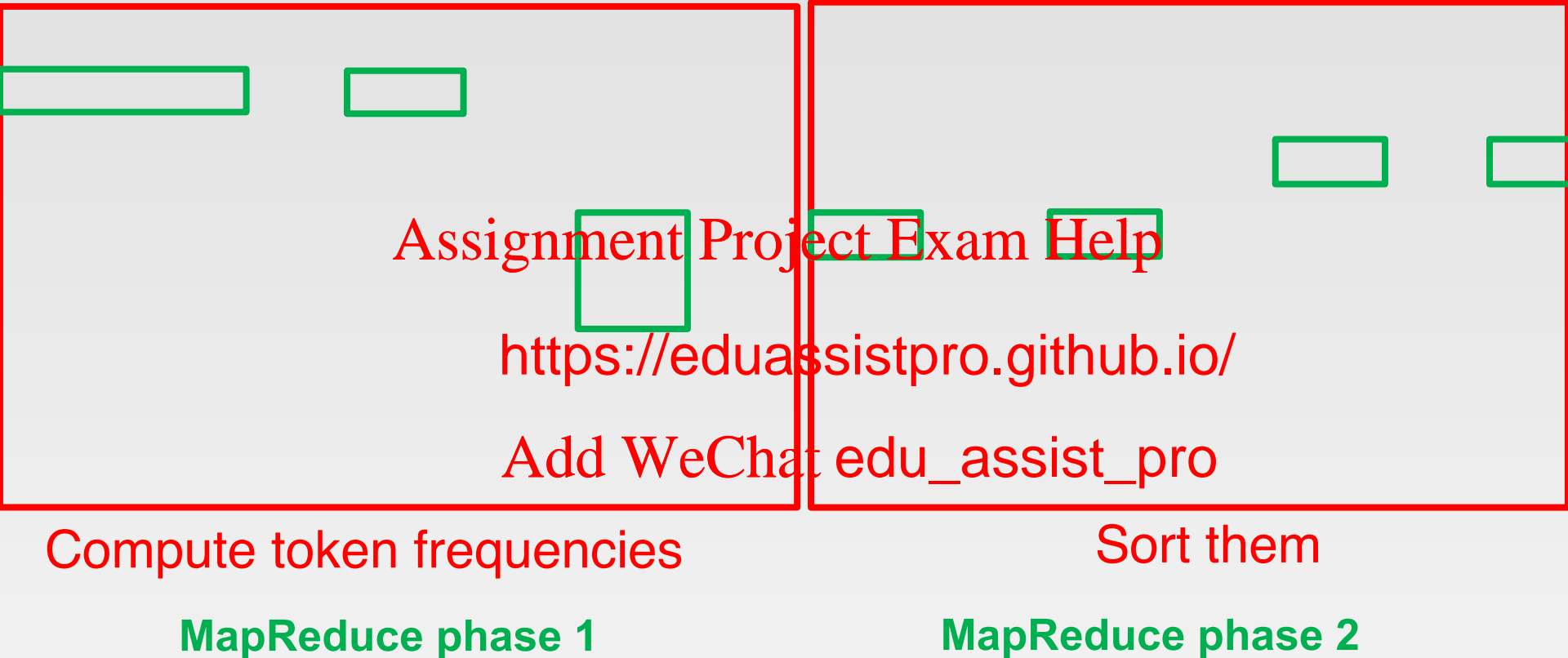
- Stage 2: Finding <https://eduassistpro.github.io/>  
(verification)

Add WeChat edu\_assist\_pro

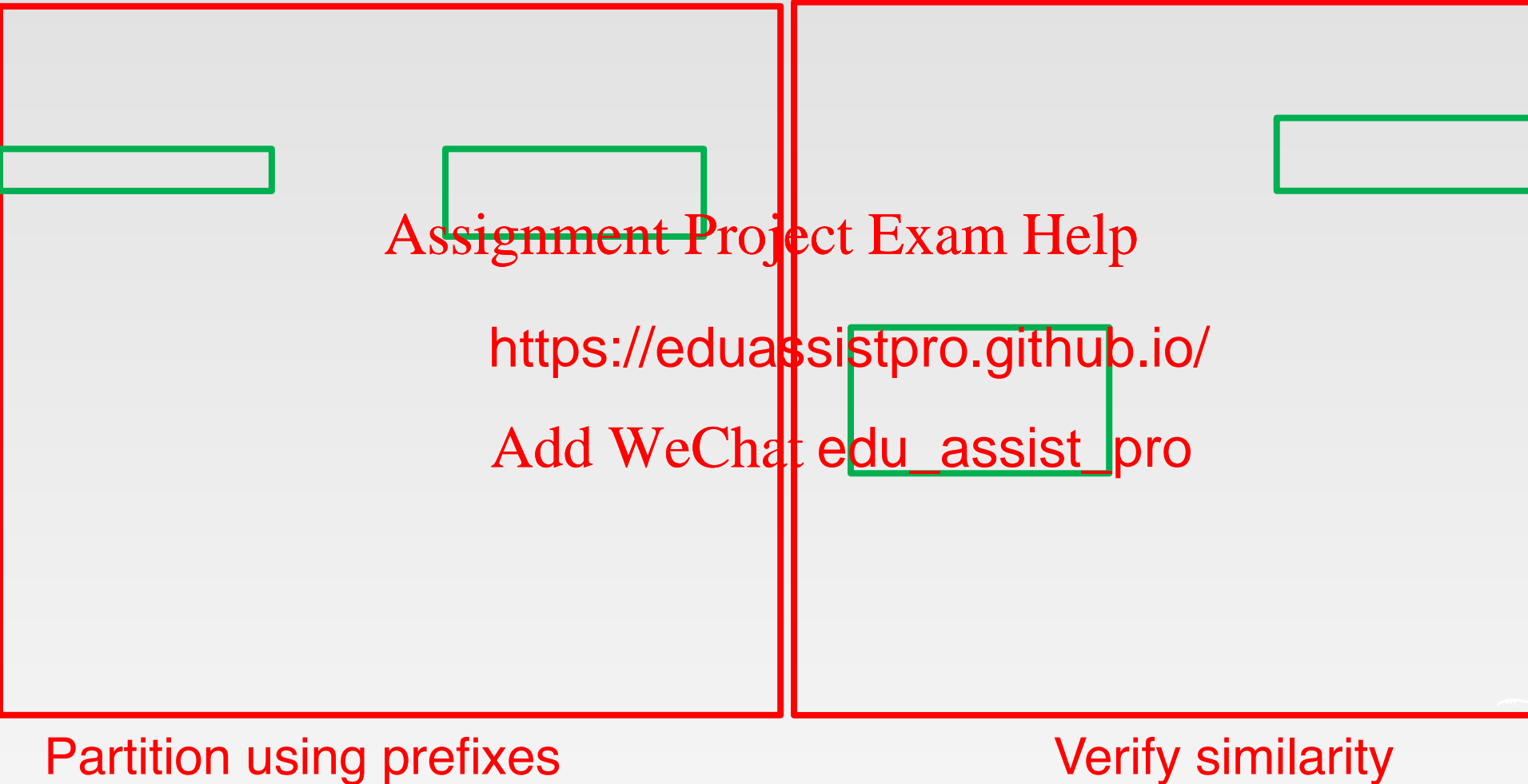
- Stage 3: remove duplicates



# Stage 1: Sort tokens by frequency



## Stage 2: Find “similar” id pairs



# Compute the Length of Shared Tokens

- Jaccard Similarity:  $\text{sim}(r, s) = |r \cap s| / |r \cup s|$
- If  $\text{sim}(r, s) \geq \tau$ ,  $l = |r \cap s| \geq |r \cup s| * \tau \geq \max(|r|, |s|) * \tau$
- Given a record  $r$ , you can compute the prefix length as  $p = |r| - l + 1$
- $r$  and  $s$  is a candidate if and only if  $r$  and  $s$  share at least one token in the first  $(|r| - l + 1)$  tokens
- Given a record  $r = (A, B, C, D)$  and  $p = 2$ , the mapper emits  $(A, r)$  and  $(B, r)$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Stage 3: Remove Duplicates

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# More Optimization Strategies

- **It is your job!!!**

- **The faster the better**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro