

DSMP Project

Chosen Dataset: Reduced_57k.txt

List of Real Users
Thuncher
akimboart
teachmesocial
TheLeeWilliams
EliasTheodorou
wichitaksgirl
HRBlockCanada
CBCcathyaalex
MykelJEstes
HemingwayCats
leftdog
WeighLossDrinks
ruthvenb
yumicute
jamestakeo
caitlinrush
MarshallJulius
USBCanada
michael_cruises
lucyheaps24
RyersonU
TheCatTweeting
weathernetwork
TorontoStar

For each algorithm, we ran the simulator on the dataset 3 times for each node. For the calculation, we took the average of the 3 runs to get a more accurate value.

Algorithm 1: K-Means

Node 1: $(545283 \text{ ms} + 563861 \text{ ms} + 568141 \text{ ms}) / 3 = 559095 \text{ ms}$

Node 2: $(229088 \text{ ms} + 264556 \text{ ms} + 285418 \text{ ms}) / 3 = 259687.3 \text{ ms} \sim 259687 \text{ ms}$

Node 4: $(203888 \text{ ms} + 231236 \text{ ms} + 259480 \text{ ms}) / 3 = 231534.7 \text{ ms} \sim 231535 \text{ ms}$

Node 8: $(212028 \text{ ms} + 221504 \text{ ms} + 181487 \text{ ms}) / 3 = 205006.3 \text{ ms} \sim 205006 \text{ ms}$

Algorithm 2: Similarity

Node 1: $(79382 \text{ ms} + 81263 \text{ ms} + 80495 \text{ ms}) / 3 = 80380 \text{ ms}$

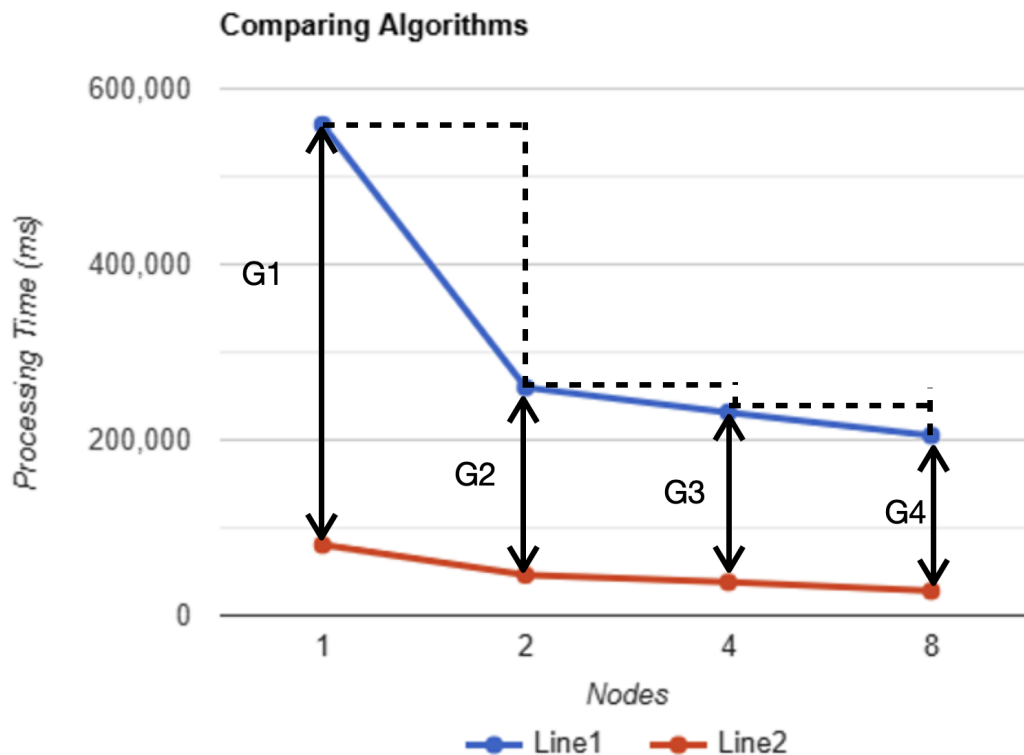
Node 2: $(45063 \text{ ms} + 44852 \text{ ms} + 48431 \text{ ms}) / 3 = 46115.3 \text{ ms} \sim 46115 \text{ ms}$

Node 4: $(36894 \text{ ms} + 37789 \text{ ms} + 38982 \text{ ms}) / 3 = 37888.3 \text{ ms} \sim 37888 \text{ ms}$

Node 8: $(26788 \text{ ms} + 29793 \text{ ms} + 26985 \text{ ms}) / 3 = 27855.3 \text{ ms} \sim 27855 \text{ ms}$

Questions

a. Plot the timing vs. nodes for each algorithm (you will plot 2 curves).



b. How do you determine the processing time gain for any number of nodes used in any algorithm? Write this in a formula.

The processing time gain can be calculated as the percentage decrease in processing time when using multiple nodes compared to a single node.

The formula is:

$$\text{Gain}(\%) = ((T_1 - T_n) / T_1) \times 100$$

Where:

- T_1 is the processing time with 1 node.
- T_n is the processing time with n nodes.

For example:

$T_1 = 559095$ ms (K-Means with 1 node)

$T_2 = 259687$ ms (K-Means with 2 nodes)

$$\begin{aligned} G\% &= ((559095 - 259687) / 559095) \times 100\% \\ &= (299408 / 559095) \times 100\% \\ &= (0.5335522) \times 100\% \\ &= 53.55\% \end{aligned}$$

- So, the processing time gain for K-Means at 2 nodes is **53.55%**

Comparison Table of Processing Time Gains (%)

Nodes	K-Means Gain (%)	Similarity Gain (%)
2 Nodes	53.55%	42.63%
4 Nodes	58.59%	52.86%
8 Nodes	63.33%	65.35%

- This table shows that the **Similarity algorithm** has better scaling at 8 nodes compared to K-Means.

c. Which algorithm yields the best timing?

Based on the results,

- K-Means Algorithm: The processing time with 8 nodes is 205006 ms.
- Similarity Algorithm: The processing time with 8 nodes is 27855 ms.

Since the Similarity Algorithm has a significantly lower processing time compared to K-Means, it is the better algorithm in terms of timing.

Parallelization Gain (G1, G2, G3, G4)

K-Means Algorithm:

- G1 = 299408 ms (Reduction from 1 to 2 nodes)
- G2 = 28152 ms (Reduction from 2 to 4 nodes)
- G3 = 26529 ms (Reduction from 4 to 8 nodes)
- G4 = 205006 ms (Final execution time with 8 nodes)

Similarity Algorithm:

- G1 = 34265 ms (Reduction from 1 to 2 nodes)
- G2 = 8227 ms (Reduction from 2 to 4 nodes)
- G3 = 10033 ms (Reduction from 4 to 8 nodes)
- G4 = 27855 ms (Final execution time with 8 nodes)

Since Similarity has a lower final execution time and better scaling, it is the better algorithm in terms of performance.