

# An Improved Taxonomy Re-structuring Using a Modified K-means Clustering for Efficient Large-scale Text Classification

Abubakar Ado<sup>1</sup>, Abdurra'uf Garba Sharifai<sup>2</sup>, Mansir Abubakar<sup>3</sup>, Bashir Salisu<sup>4</sup>, Usman Mahmud<sup>5</sup>, Abdulkadir Abubakar Bichi<sup>6</sup>

<sup>1,2,4,5,6</sup>Department of Computer Science, Northwest University Kano

<sup>4</sup>Department of Computer Science, Al-Qalam University, Katsina



## Introduction

Textual classification to a hierarchical taxonomy of classes is a common and known problem associated with Large-Scale Text classification (LSTC). When there are many labels, hierarchical restructuring of classes has been recognized as a natural and effective way to organize similar classes and it has been well studied in the past two decades.

## Problem Statement

When there are many classes with increase number of features, the exiting hierarchy restructuring methods tend to produce many nodes with similar granularities. This results in mis-classification, computationally expensive, not scalable for many classification models.

## mKmeans Algorithm



## Proposed Method

### Algorithm: Res-mKMeans Algorithm

**Input:** Original hierarchy  $H$ , dataset  $X$ , class labels  $y$ , number of clusters  $K$  for restructuring

**Output:** Modified hierarchy  $H_{new}$

$H_{new} \leftarrow H$

for each parent node  $p$  in  $H_{new}$  do

$C_p \leftarrow$  children classes of  $p$

$X_p \leftarrow$  samples belonging to classes in  $C_p$

if  $|C_p| > K$  then

$(cluster\_ids, centroids) \leftarrow$  mKMeans( $X_p, K$ )

$NewGroups \leftarrow \{G_1, G_2, \dots, G_K\}$

for each class  $c \in C_p$  do

assign class  $c$  to  $G_{(cluster\_ids(c))}$

end for

replace children( $p$ ) with  $NewGroups$

end if end for

for each parent  $p$  in  $H_{new}$  do

children\_list  $\leftarrow$  children( $p$ )

for each child  $g$  in children\_list do

if size( $g$ ) == 1 then continue

end if

if purity( $g$ ) < threshold then

$(cid, cent) \leftarrow$  mKMeans(samples( $g$ ),  $K_{refine}$ )

$ReGroups \leftarrow \{Rg_1, \dots, Rg_{K_{refine}}\}$

for each class  $c$  in  $g$  do

assign  $c$  to  $Rg_{(cid(c))}$  end for

replace  $g$  in children\_list with  $ReGroups$

end if end for end for

for each node  $n$  in  $H_{new}$  do

if number\_of\_children( $n$ ) < min\_children then

$H_{new} \leftarrow$  NDD( $n, H_{new}$ )

else if number\_of\_children( $n$ ) > max\_children then

$(cid, cent) \leftarrow$  mKMeans(samples( $n$ ),  $K_{split}$ )

$SplitGroups \leftarrow$  assign classes using  $cid$

replace children( $n$ ) with  $SplitGroups$

end if end for

return  $H_{new}$

## Results Analysis

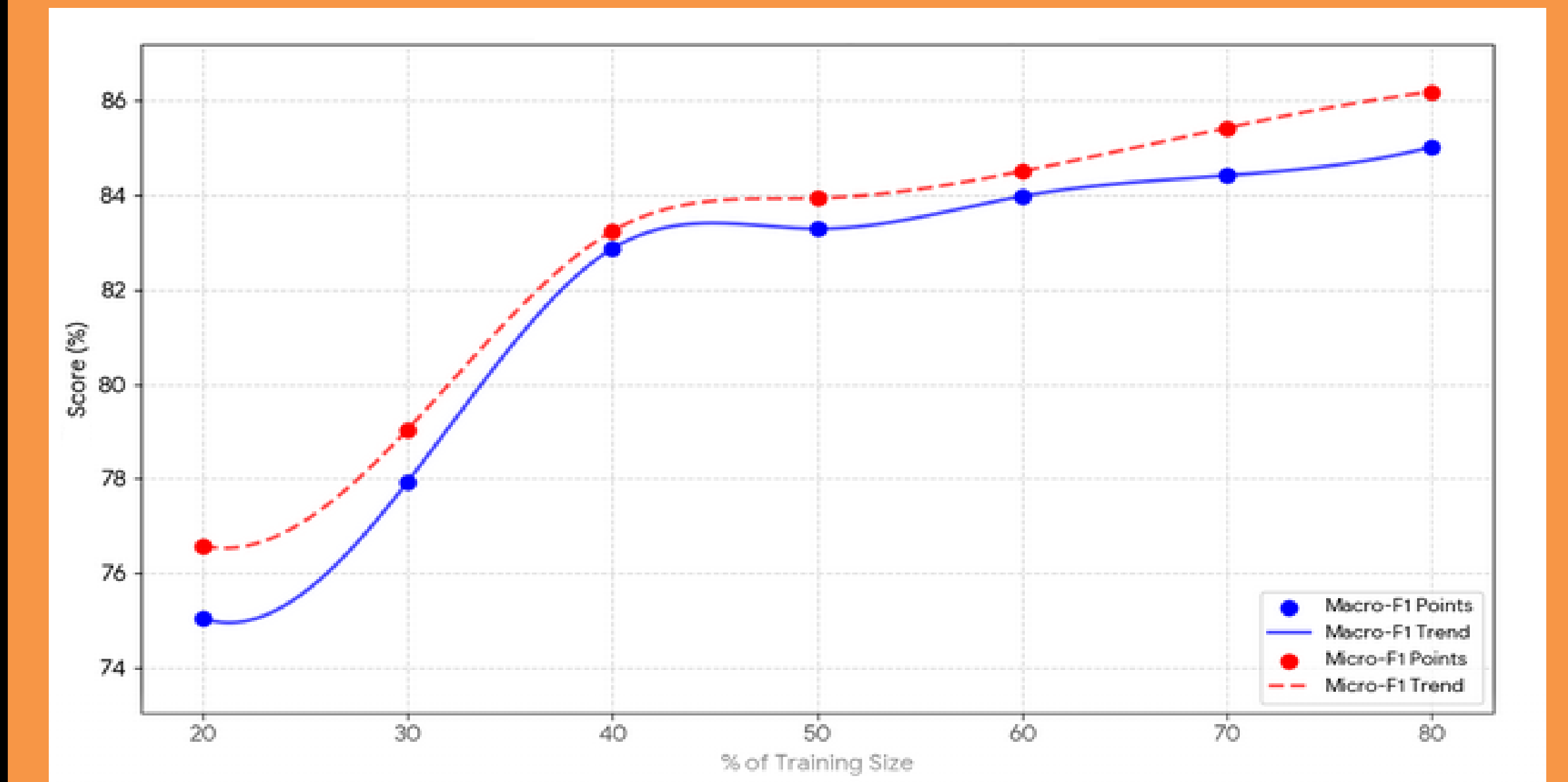


Fig. 2: Macro-f1 and Micro-f1 Performance of the proposed method with varying % of training size

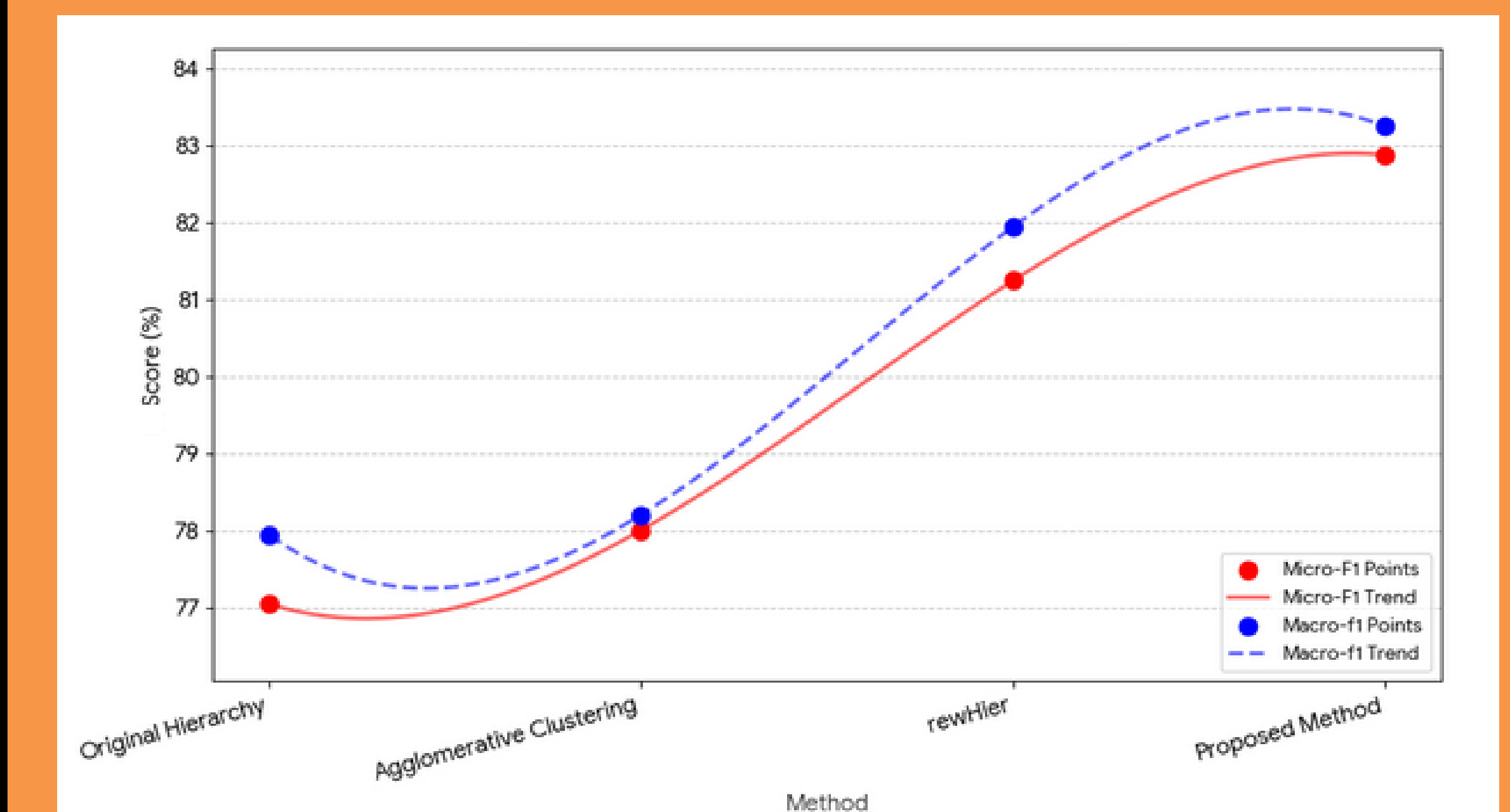


Fig. 2: Micro-f1 and Macro-f1 performance comparison using different hierarchy restructuring approaches on newsgroup dataset

## Conclusion

We propose an approach for restructuring hierarchy that is more suited for HC. In comparison to existing approaches, our method gives better performance that allow HC approaches to significantly scale to LSHC

## Acknowledgments

We would like to thank the TetFund, Nigeria and Northwest University, kano for supporting this research under Institution based research, IBR-2024.