

Enhancing Wikipedia Article “See Also” Section through Article Features and NLP-Generated Semantic Vectors

How can a system be developed to automate the creation and update of the 'See Also' section in Wikipedia articles by utilizing article features and NLP-generated semantic vectors?

By:

Supervised by:

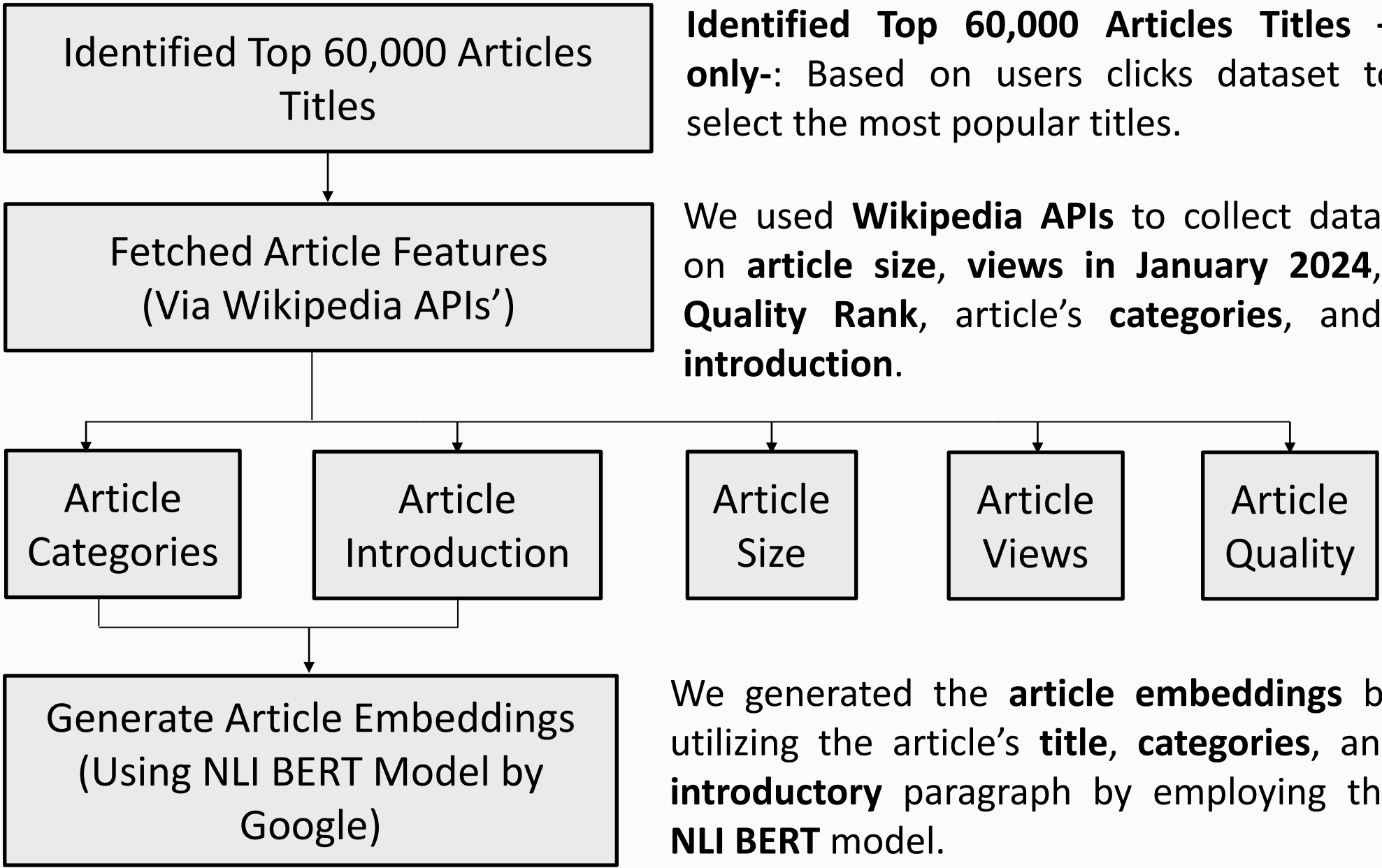
Abstract and Motivation

In today's rapidly expanding information age, Wikipedia stands out as a primary source of free knowledge with over 25 billion monthly visitors*. Understanding the need to keep up with this growth and ensure information is both relevant and high quality, we've started working on creating a system that automates the updating of the "See Also" sections for Wikipedia's top 60,000 articles. Our project employs machine learning and natural language processing techniques to simplify the discovery and linkage of related articles, aiming to enhance navigability and enrich the reader experience on the site. The evaluation phase of our system is set to begin on 10 March and will last for 3 weeks. Early feedback from Wikipedia's admins has been encouraging, showing that our system could effectively save the site’s editors’ efforts in manually updating these sections.

*Source: Wikimedia Statistics

Dataset Preparation Overview

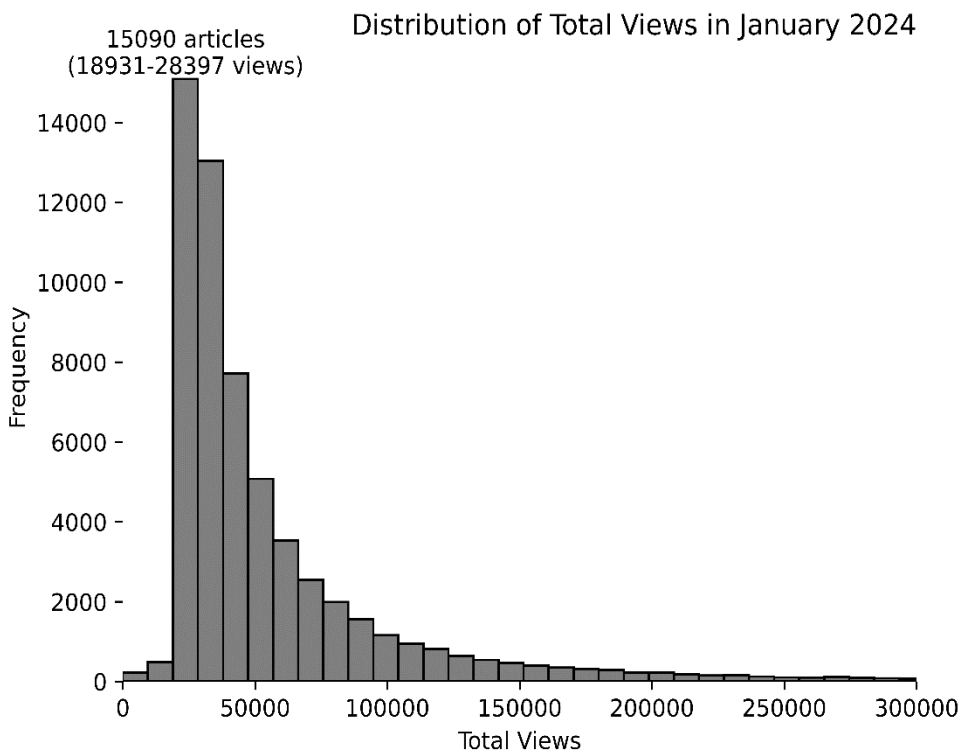
We meticulously compiled a dataset to fuel our automated system, focusing on features that blend popularity with quality. The steps were:



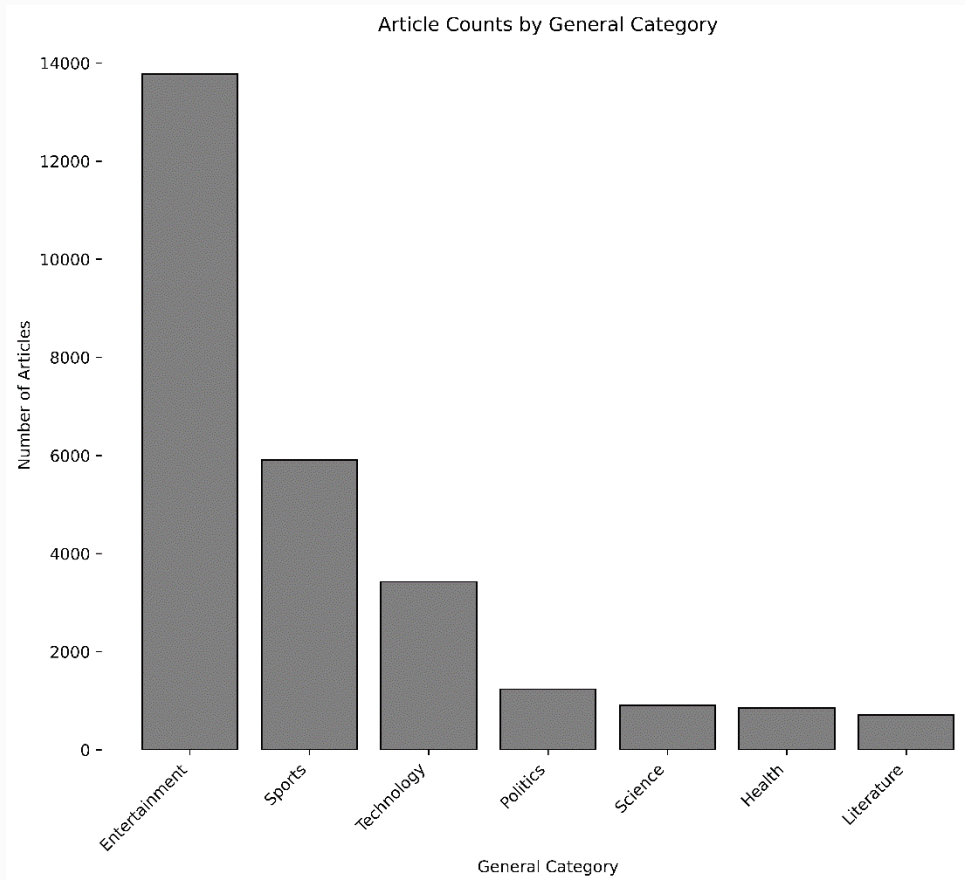
	title	size	total_views	first_paragraph	article_quality	article_categories	bert_0	bert_1	bert_2	bert_3
0	Hyphen-minus	11286	6486	The hyphen-minus symbol - is the form of hyphe...	C	['Punctuation', 'Typographical symbols']	-0.628009	0.816350	-0.294891	0.324033
1	Saltburn_(film)	69335	7450496	Saltburn is a 2023 black comedy psychological ...	B	['2020s American films', '2020s British films'...	-0.447147	0.665674	0.293427	-0.021366

Exploratory Data Analysis

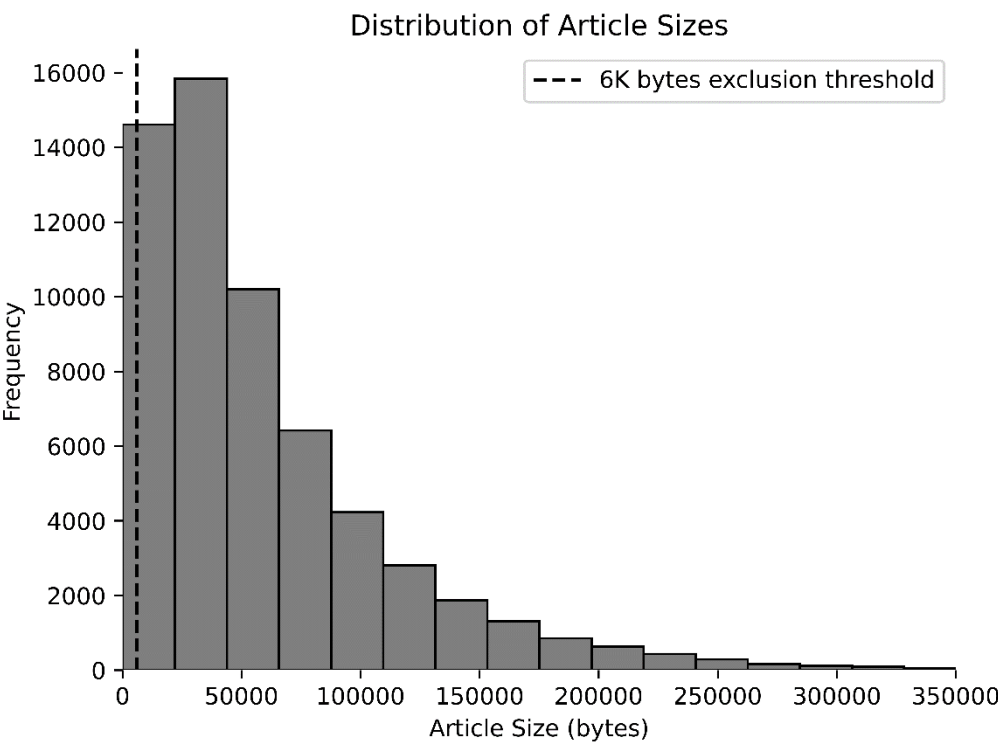
In our Exploratory Data Analysis, we carefully examined the dataset to inform our methodology for automating the 'See Also' section in Wikipedia articles.



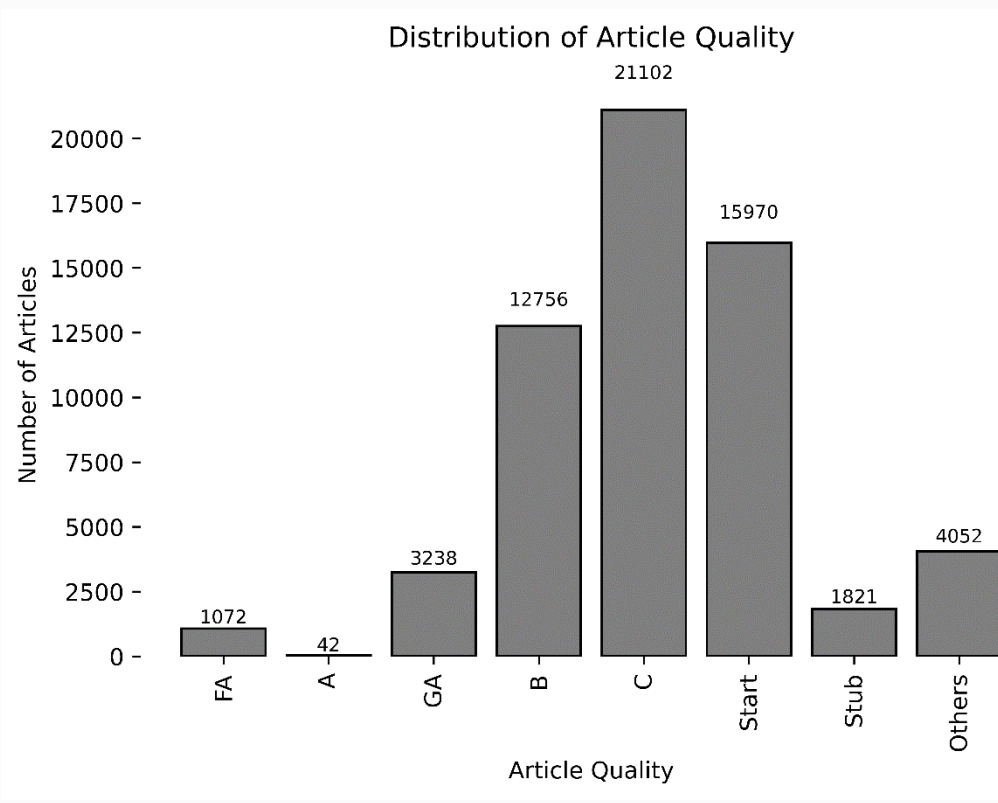
Views Distribution: Assessed to gauge the popularity of articles, confirming our selections resonate with reader interests.



Categories Distribution: Explored to guarantee a diverse selection across various subjects, enhancing the general applicability of our automated system.



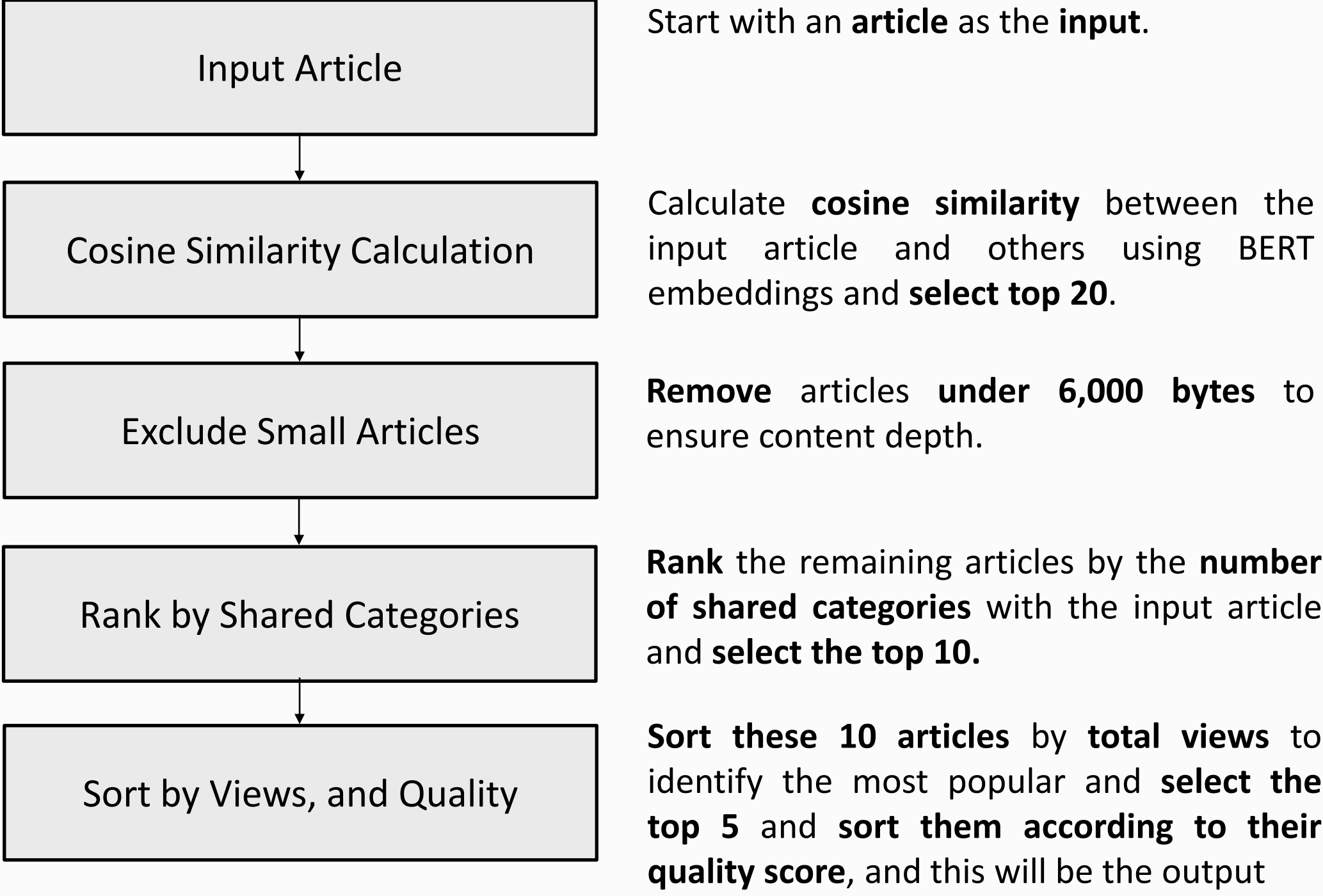
Article Size Distribution: Analyzed to later exclude articles under 6KB, ensuring the recommendations are of substantial content



Quality Distribution: Evaluated with a focus on ascending quality, prioritizing 'Featured Article' (FA) as the highest quality, to ensure recommendations are high quality.

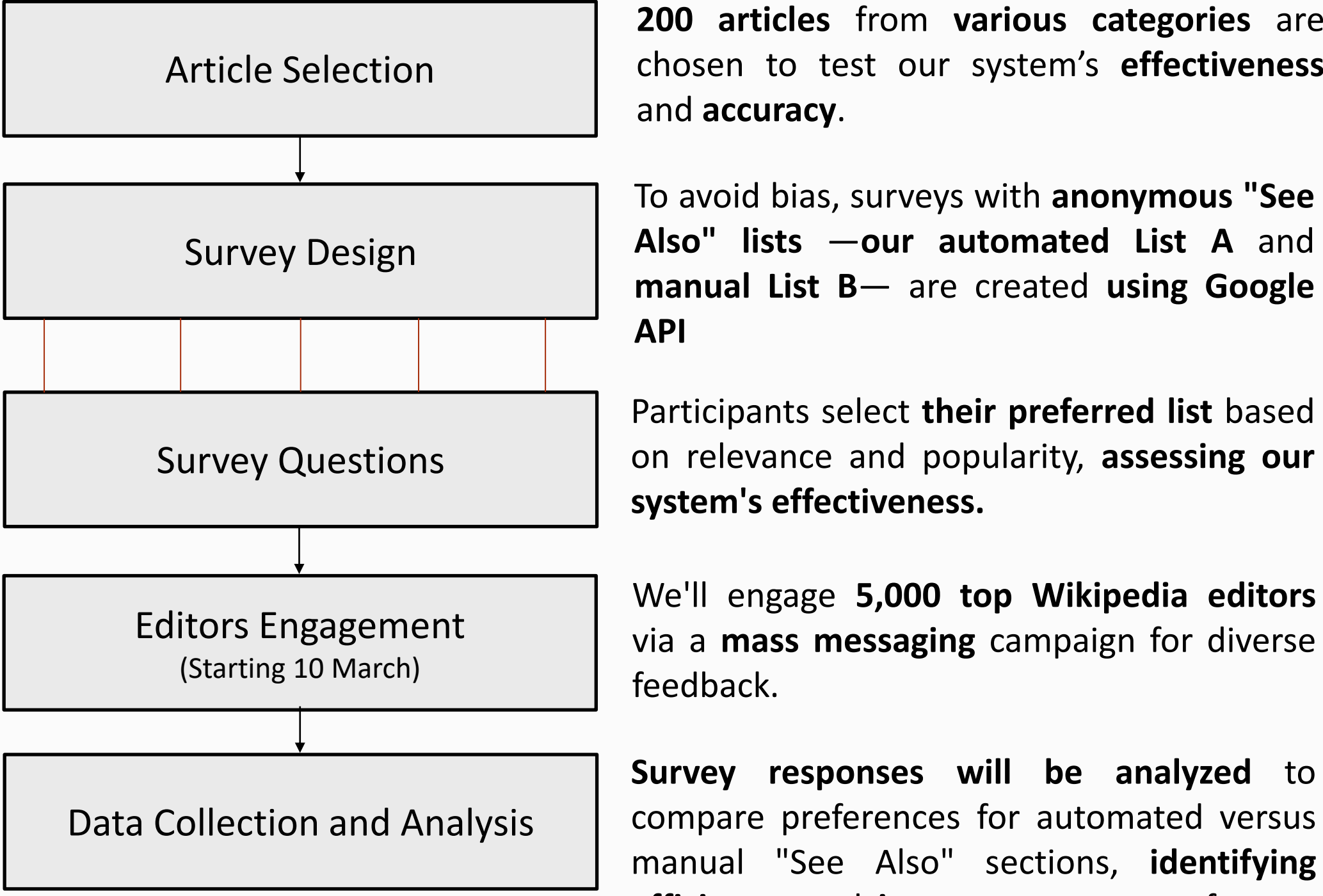
Recommendation System Methodology

Outlined below is our system methodology for automating the Wikipedia’s articles 'See Also' sections, leveraging NLP techniques and other features.



Evaluation, and Results

To assess the effectiveness of our automated "See Also" section generation system and its alignment with Wikipedia's high standards, we have outlined a detailed methodology:



Here are some results from our system's performance compared to current manual lists:

Our System	Current Manual List
	Open Data Science Conference
Branches of science	Scientific Data
Data analysis	Women in Data
Bayesian inference	Python (programming language)
Computer science	R (programming language)
Deep learning	

Data Science Article

Our System	Current Manual List
Niels Bohr	Earth science
Albert Einstein	Neurophysics
Isaac Newton	Psychophysics
	Relationship between mathematics and
James Clerk Maxwell	Science tourism
Theory of relativity	

Physics Article

