

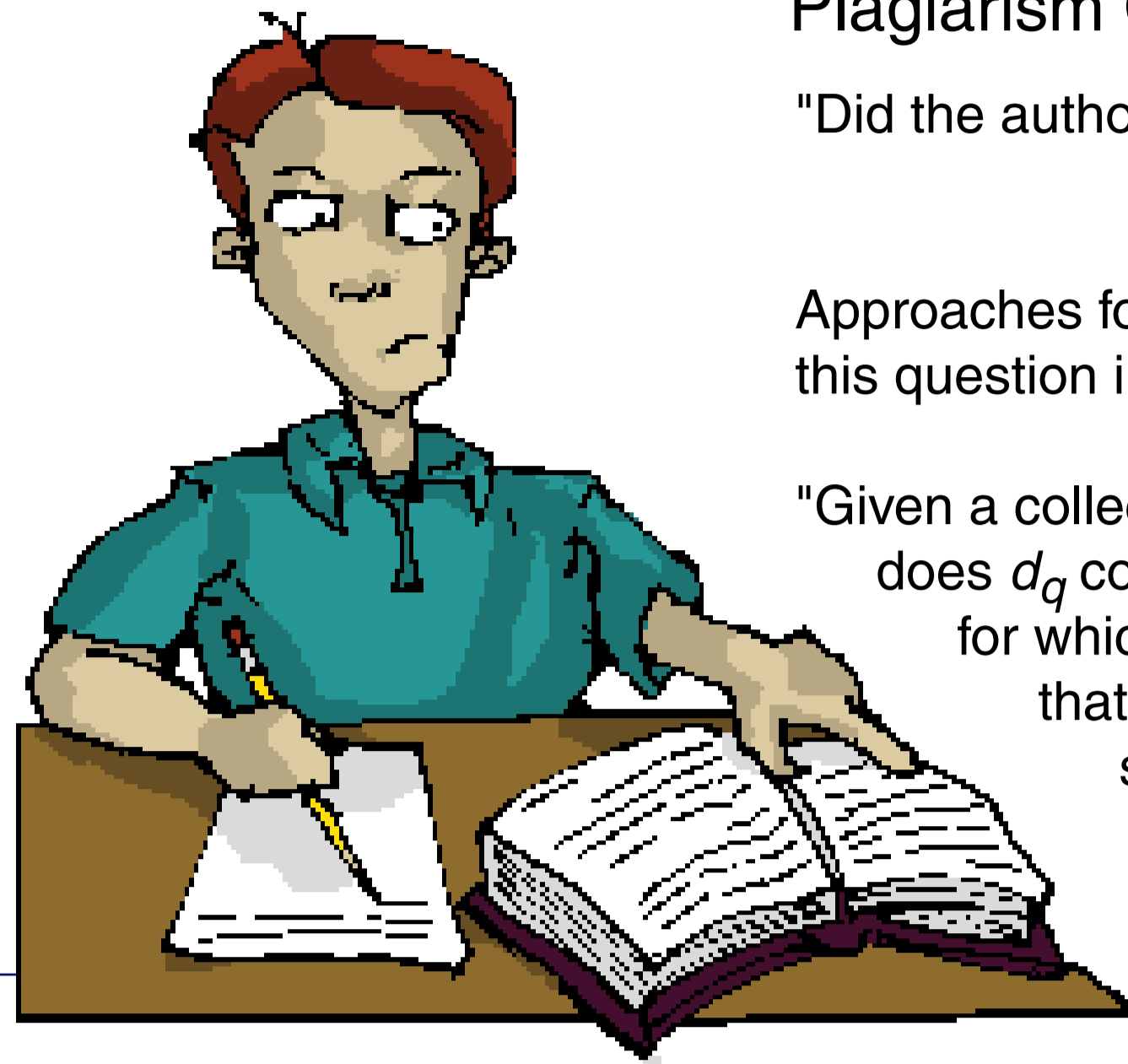
# Strategies for Retrieving Plagiarized Documents

## Contribution:

A tailored index for plagiarism analysis

For the identification of plagiarized passages in large document collections we present retrieval strategies which rely on stochastic sampling and chunk indexes.

Using the entire Wikipedia corpus we compile n-gram indexes and compare them to a new kind of fingerprint index in a plagiarism analysis use case. Our index provides an analysis speed-up by factor 1.5 and is an order of magnitude smaller, while being equivalent in terms of precision and recall.



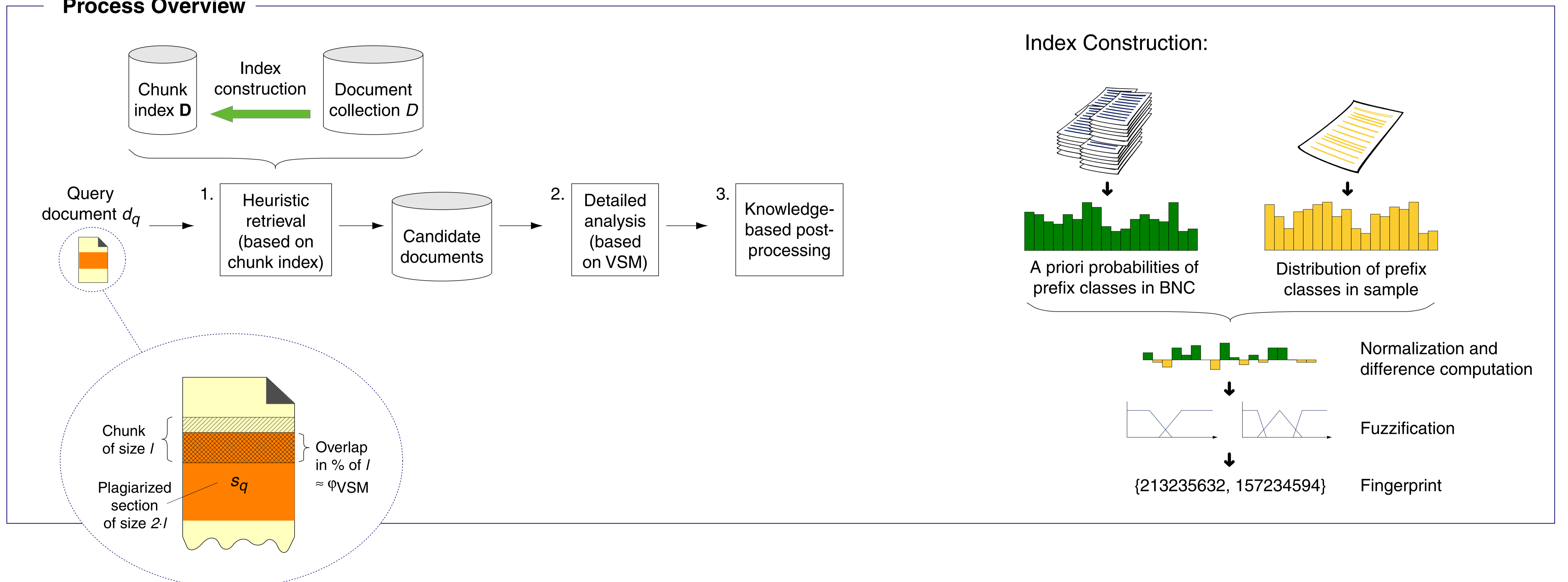
## Plagiarism Challenge:

"Did the author of a document  $d_q$  commit a plagiarism offense?"

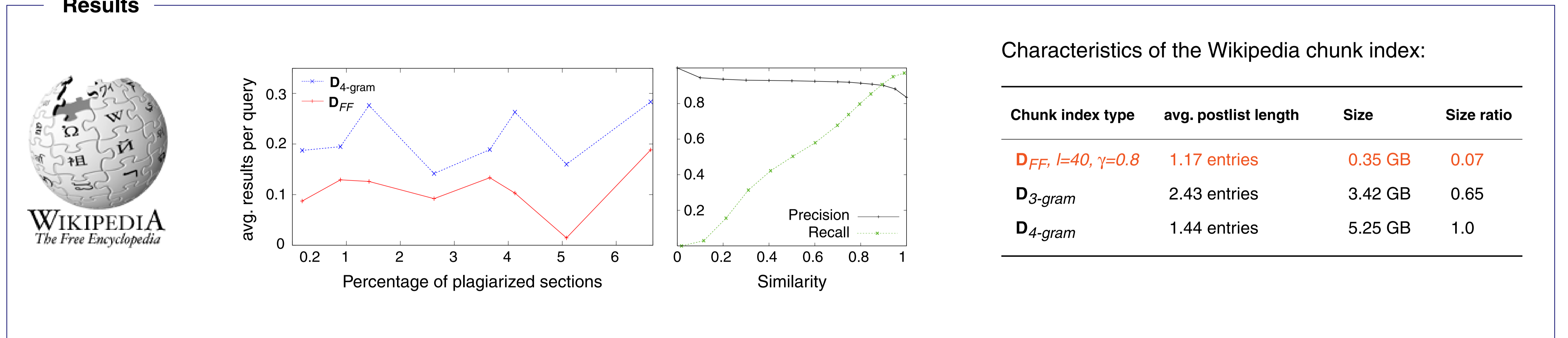
Approaches for computer-based plagiarism analysis break down this question into manageable parts:

"Given a collection  $D$  of documents, does  $d_q$  contain a section  $s_q$  for which one can find a document  $d_x \in D$  that contains a section  $s_x$  such that under some retrieval model  $\mathcal{R}$  the similarity  $\phi_{\mathcal{R}}$  between  $s_q$  and  $s_x$  is close to 1?"

## Process Overview



## Results



## Related Technology [Koppel/Schler 2004]

