

new sensible word count

4000

3000

2000

1000

0

0

1000000

2000000

3000000

4000000

5000000

6000000

7000000

number of tokens in the reference data (scaled to length of smallest language)

