

## Zipf'in laki suomenkieliselle aineistolle

Aineisto: uutisia, kirjoja, aikakauslehtiä yms.  
32 miljoonaa sanaa, 1,6 miljoonaa eri sanamuotoa.

```
1 942024 942024 ja
2 746560 1493120 on
3 282702 848106 ei
4 260709 1042836 että
5 233863 1169315 oli
6 144319 865914 myös
7 139290 975030 ovat
8 137386 1099088 hän
9 128759 1158831 se
10 122759 1227590 mutta
11 115231 1267541 mukaan
12 105534 1266408 kun
13 101741 1322633 kuin
14 99149 1388086 ole
15 91092 1366380 sen
16 82355 1317680 joka
17 78828 1340076 jo
18 73056 1315008 tai
19 68864 1308416 suomen
20 64151 1283020 ollut
21 62161 1305381 niin
22 60909 1339998 vain
23 60273 1386279 vuoden
24 60071 1441704 nyt
25 58951 1473775 sekä
26 57575 1496950 viime
27 55327 1493829 vuonna
28 54845 1535660 jälkeen
29 53778 1559562 jos
30 53117 1593510 olisi
31 52761 1635591 vielä
32 51171 1637472 hänen
33 49381 1629573 noin
```

...

```
1000 3106 3106000 teksti
1001 3103 3106103 ansiosta
1002 3098 3104196 tilaa
1003 3098 3107294 kuuden
1004 3097 3109388 syntyy
1005 3096 3111480 valittu
1006 3096 3114576 sami
1007 3093 3114651 kirjan
1008 3092 3116736 joulukuussa
1009 3088 3115792 hoitaa
1010 3087 3117870 seuraavan
1011 3085 3118935 onnistui
1012 3080 3116960 eivätkä
1013 3079 3119027 suurempi
1014 3076 3119064 edellä
1015 3075 3121125 ensimmäiset
1016 3075 3124200 alun
1017 3070 3122190 tietojen
1018 3065 3120170 hyviä
1019 3063 3121197 hyvät
```

...

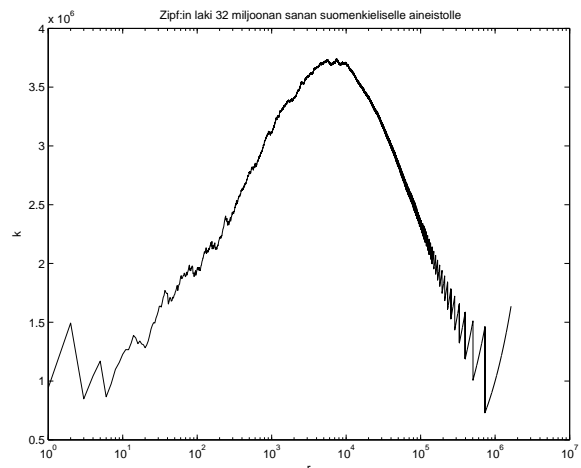
```
10000 370 3700000 sivua
10001 370 3700370 purkamaan
10002 370 3700740 poliisilla
```

```
10003 370 3701110 pienillä
10004 370 3701480 petti
10005 370 3701850 nälkä
10006 370 3702220 muistoja
10007 370 3702590 luokka
10008 370 3702960 lillehammerin
10009 370 3703330 lausunnossaan
10010 370 3703700 kulttuurien
10011 370 3704070 kriittisesti
10012 370 3704440 kiteen
10013 370 3704810 kirjaan
10014 370 3705180 kiinnostava
10015 370 3705550 kielteisesti
10016 370 3705920 jalkapallomaajoukkueen
10017 370 3706290 edellytetään
10018 370 3706660 dos
10019 370 3707030 aikaansa
```

...

```
730884 2 1461768 aakkosittain
730885 2 1461770 aakkosissa
730886 2 1461772 aakkosen
730887 2 1461774 aakkosellista
730888 2 1461776 aakkosellisesta
730889 2 1461778 aak
730890 2 1461780 äähreen
730891 2 1461782 aahos
730892 2 1461784 äähneen
730893 2 1461786 aahmuun
730894 2 1461788 aadamille
730895 2 1461790 aaaah
730896 2 1461792 äää
730897 1 730897 zzzzzzz
730898 1 730898 zzzthzzzjrrrrrrrsttmnkrtrttt
730899 1 730899 zyskowiczistä
730900 1 730900 zyskowiczista
730901 1 730901 zyskowicziltä
730902 1 730902 zyskowiczilta
730903 1 730903 zyskowicziin
730904 1 730904 zyskowiczien
```

...



Kuva 1:  $k$   $r:n$  funktiona