# Frequency Wordlists

**Dictionary Express** 

#### Riddle

100 000 most frq. headwords ~ 1984–2024

#### Riddle

100 000 most frq. headwords ~ 1984–2024

2nd 100 000 ~ 2024-???

#### Riddle

100 000 most frq. headwords ~ 1984–2024

2nd 100 000 ~ 2024— 2064

## Frequency wordlist

verb (25,676 items )

	Word	Lemma	Tag	Frequency	
1	is	be	VBZ	730,929 ***	
2	are	be	VBP	379,918	
3	be	be	VB	344,024 ****	
4	was	be	VBD	278,396	
5	has	have	VHZ	210,972	

#### Frequency types

Absolute Frequency (+ relative F)

Document Frequency

ARF

ALDF

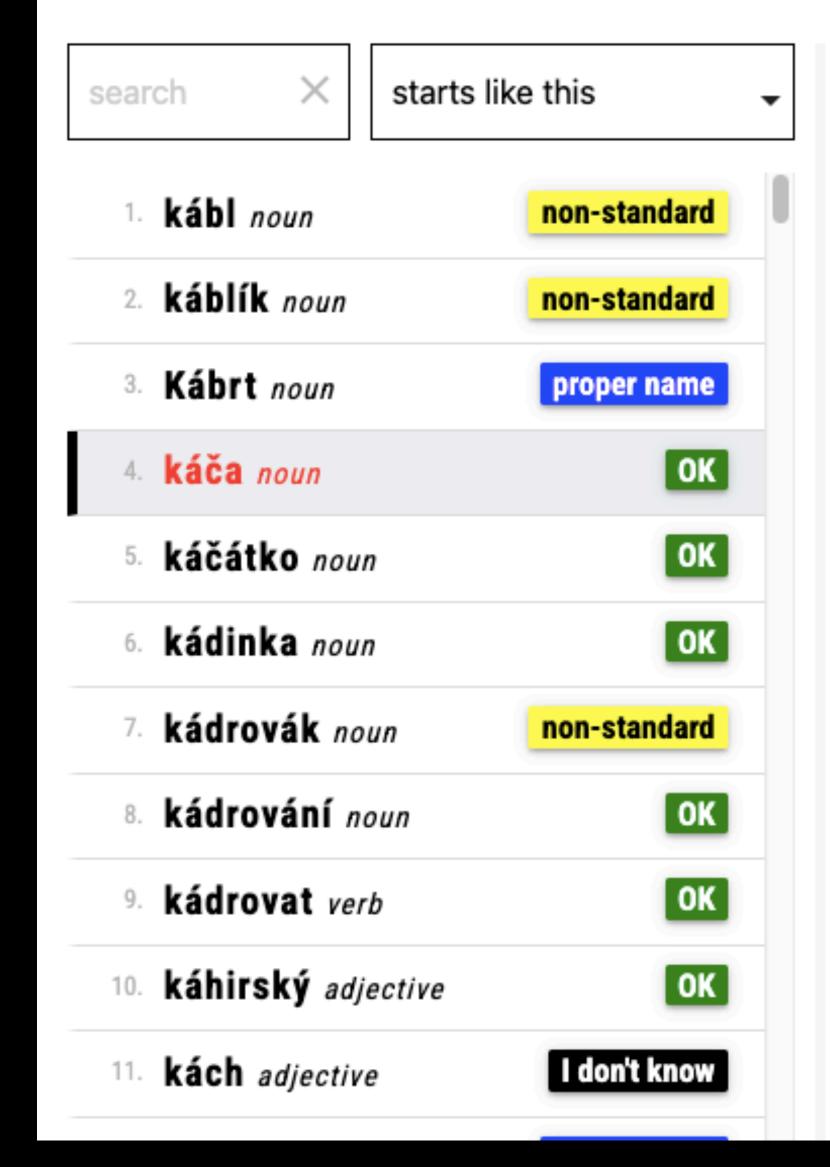
#### **Dictionary Express**

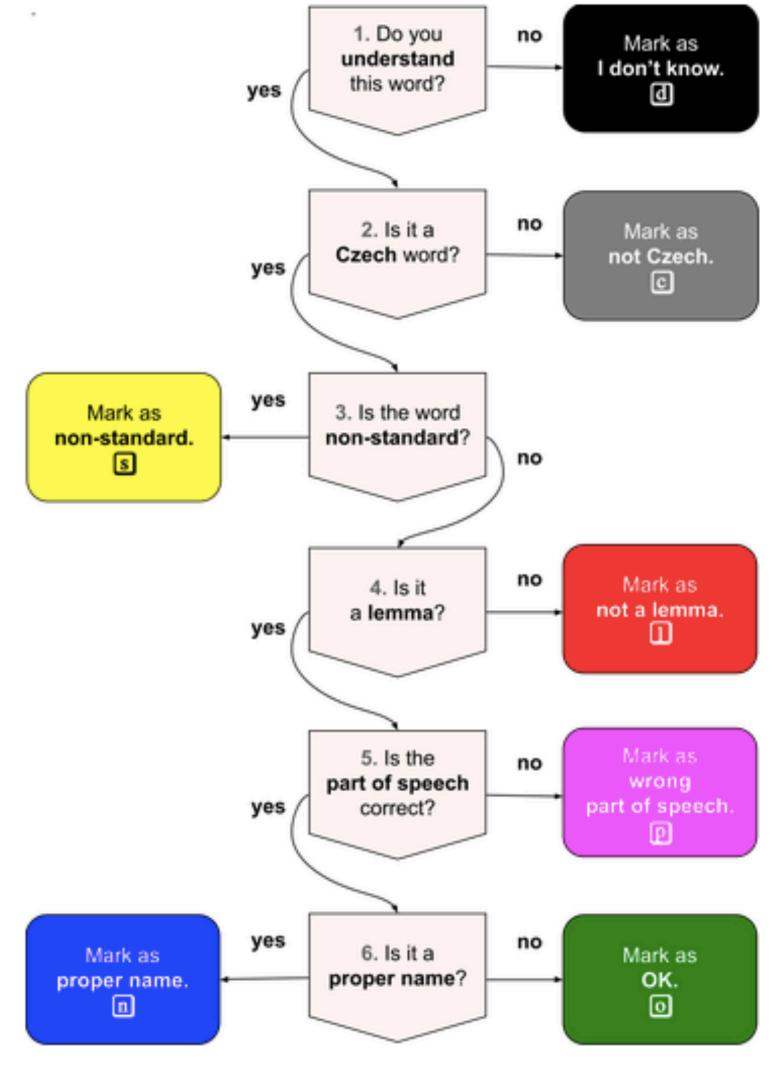
Words from corpora

+ opinions of people.

## CZECH HEADWORDS 119 (VERONIKAS)

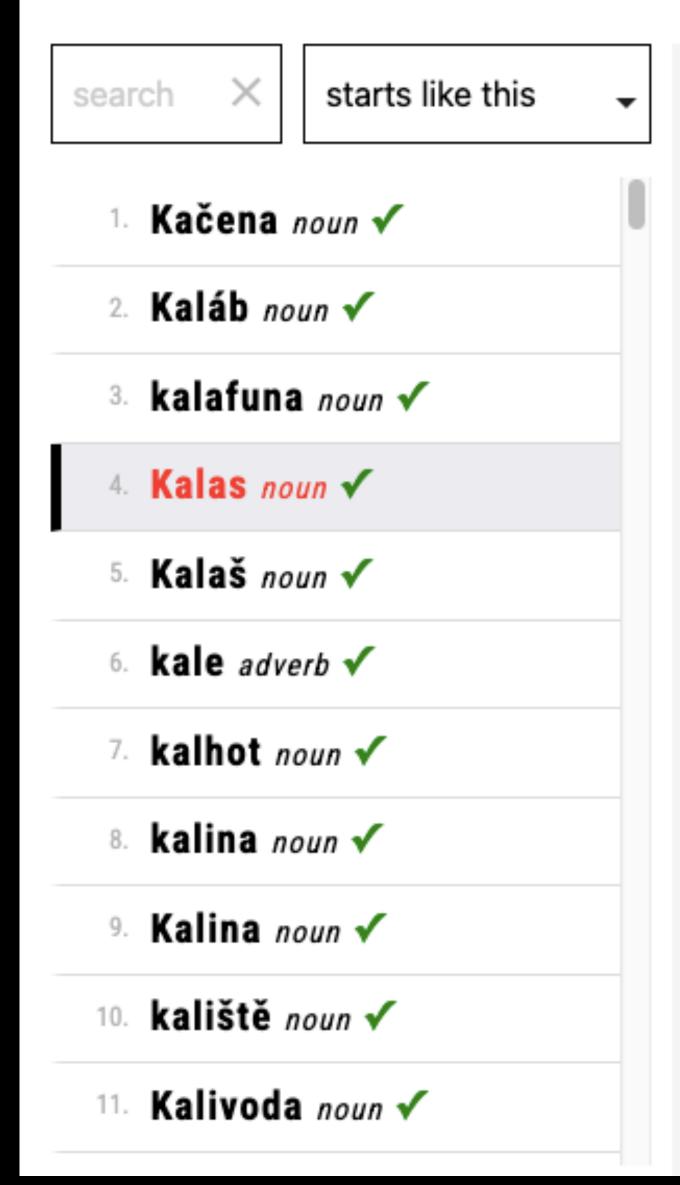
otal 1000 entries

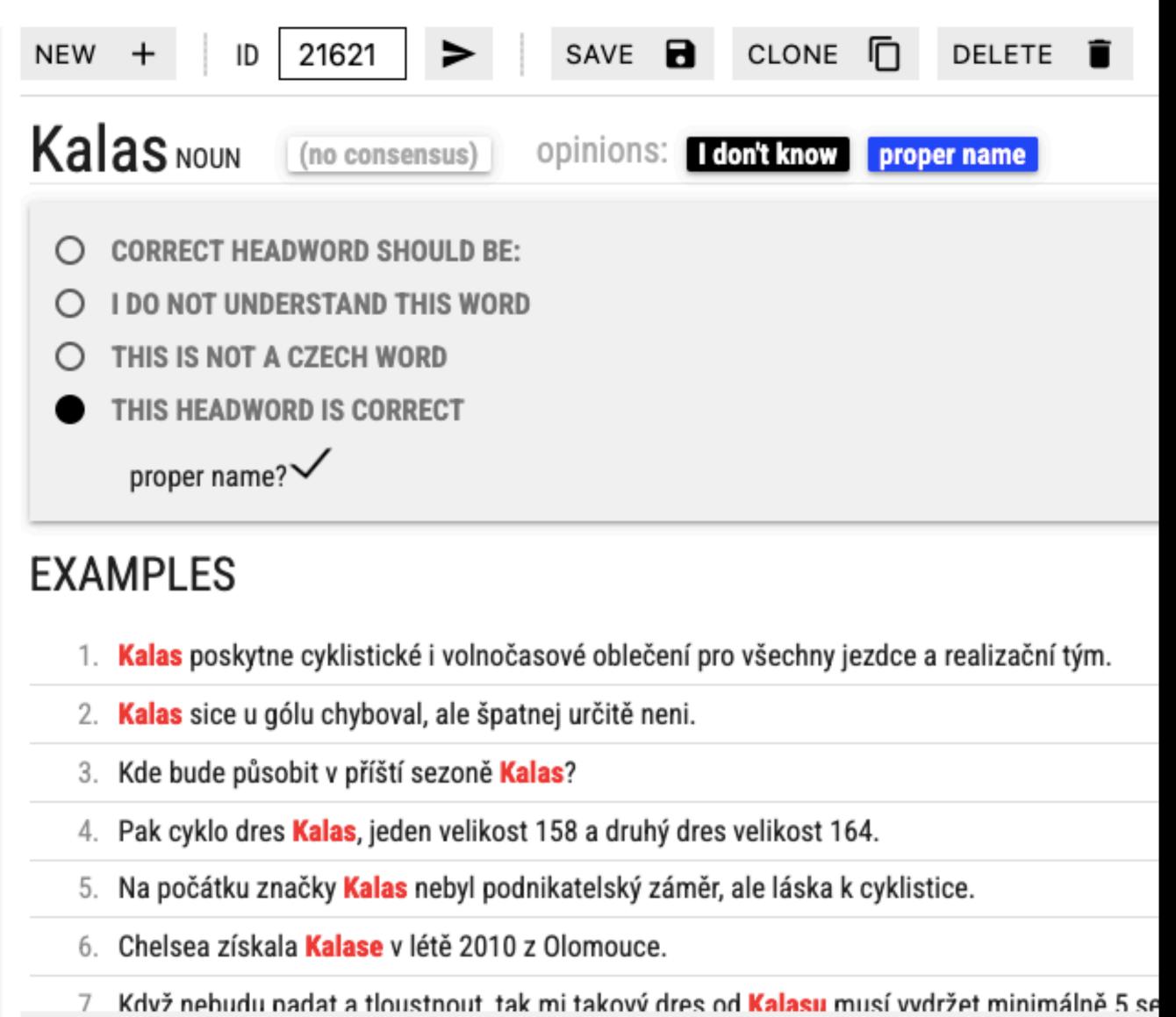




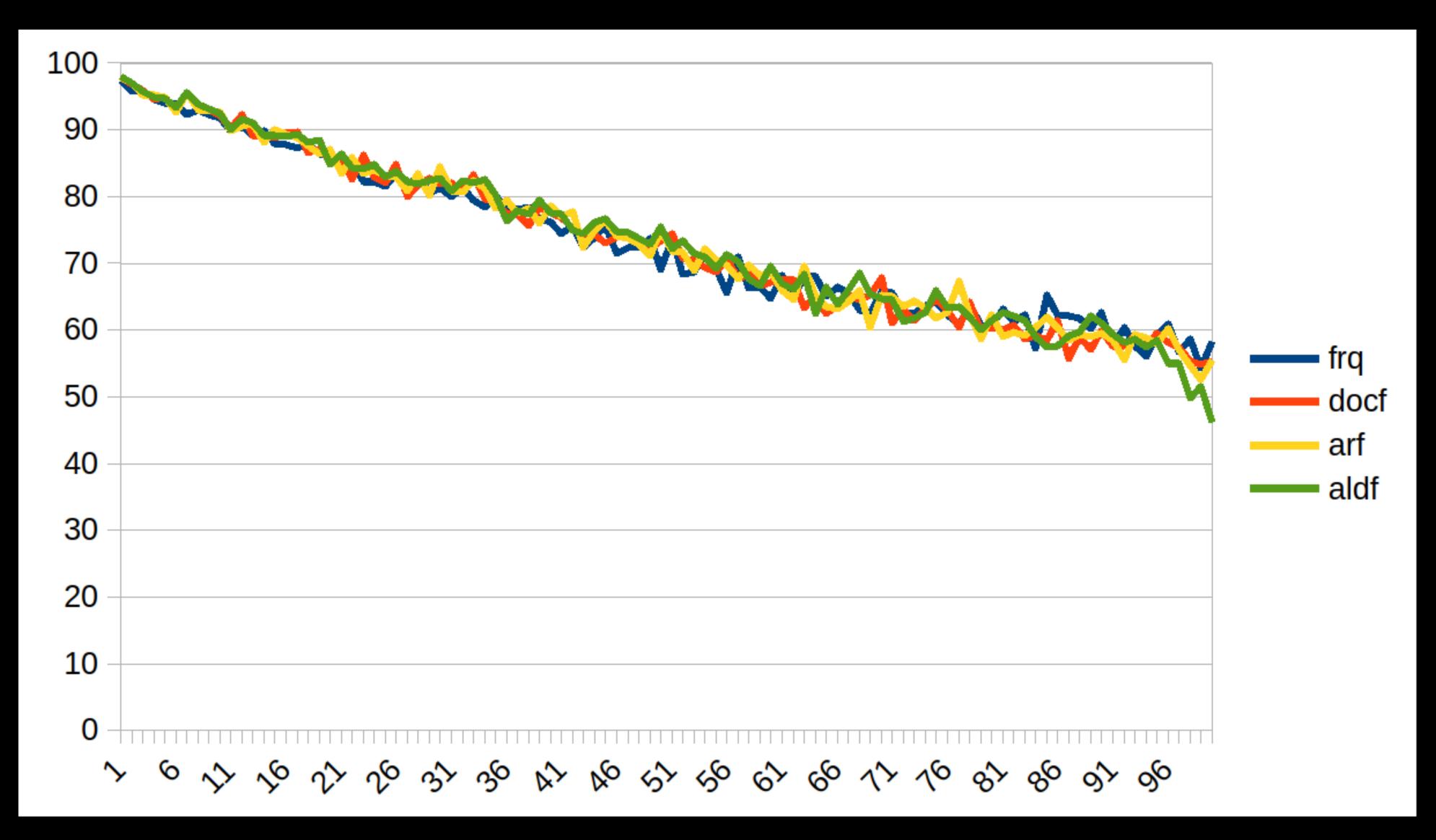
# CZECH REVISION 8 (MARTINR)

total 1000 entries





#### Frequent x "trouble-less"



#### Czech DExpress

100 000 corpus lemposes

=> cca 79 700 dictionary entries.

#### Frequency difference: a 100 000 DocF wordlist

AF: 4962

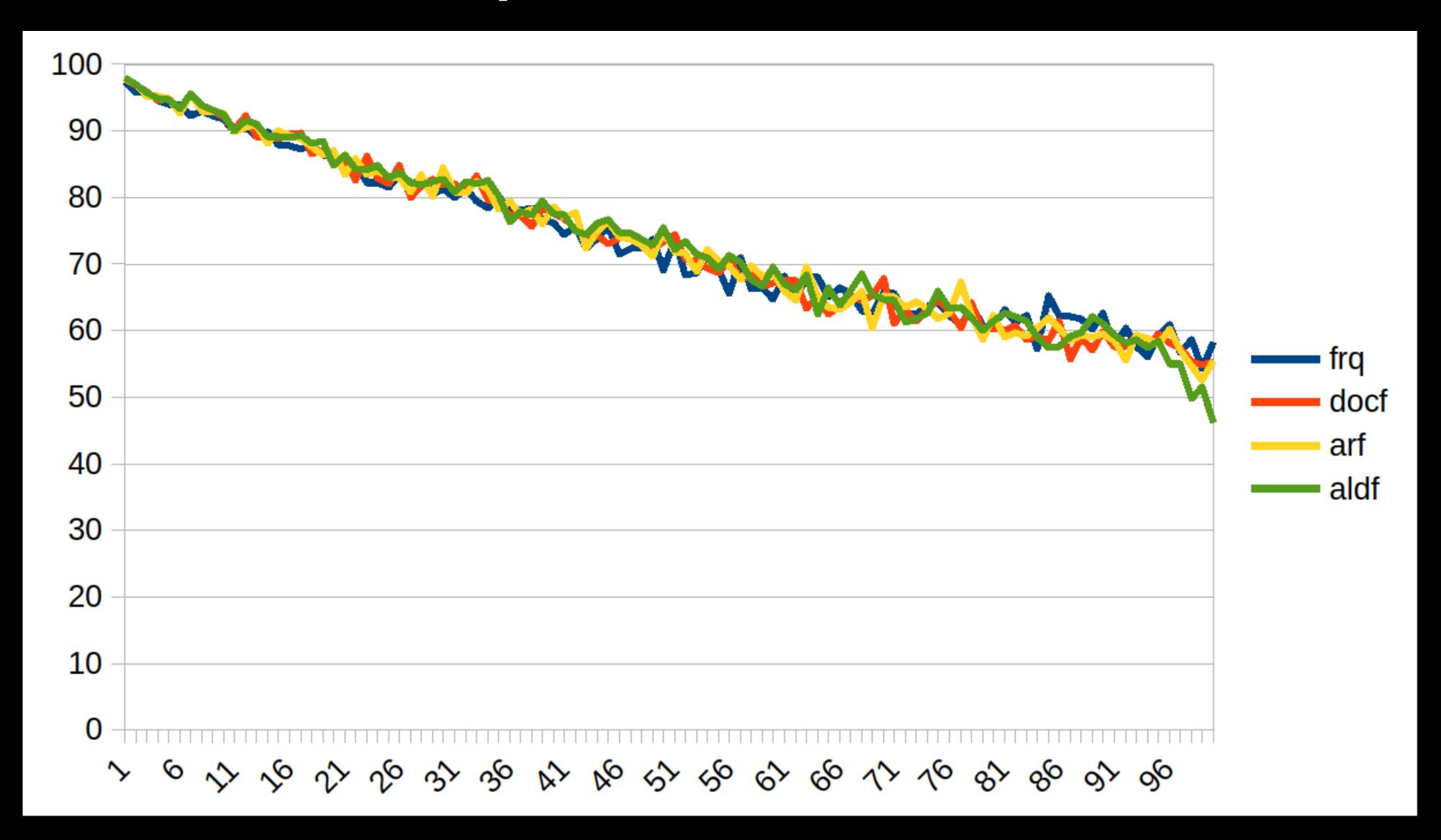
ARF: 1722

ALDF: 1927

For 100 000 entries?

Cca 150 000 corpus lemposes (assumption)

#### Frequent x "trouble-less"



### Which frequency?

	80 000	20 000	0
frq		76.78%	59.70%
docf		77.07%	58.13%
arf		77.29%	57.24%
aldf		77.62%	55.92%

#### Comparing 2 frequency wordlists

missing accepted

present accepted

missing rejected

present rejected

without annotation

#### Present accepted – absolute frequency

Vareni.cz, Echo24, Skyscanner, Ulož.to, ČSFD.cz

... lacks less frequent Czech words.

#### Present accepted – ARF and ALDF

ARF – more URLs and company names.

#### Present accepted – ARF and ALDF

ARF – more URLs and company names.

ALDF has the most "normal" Czech words.

#### Present accepted – ARF and ALDF

ARF – more URLs and company names.

ALDF has the most "normal" Czech words.

=> ALDF could be the best freq. type.

#### Incomplete research

... we still need to examine document frequency.

#### Wrapping up

Not many differences in 100k headwords.

Another 100k could be more interesting.