

N-gram 1-6, dansk (NST)

Disse n-grammene er laget ved Uni Research AS av Knut Hofland. Utgangspunkt er tekstene i det danske tekstkorpuset til Nordisk språkteknologi holding AS, som ble overført til Nasjonalbiblioteket i 2011.

N-grammene kan benyttes fritt. Brukere oppfordres til å informere Språkbanken om hvilken nytte de gjør av denne ressursen.

Eventuelle spørsmål og tilbakemeldinger kan rettes til sprakbanken@nb.no.

1. Innhold

N-grammene er basert på følgende materiale av dansk nyhetstekst fra siste halvdel av 1990-tallet:

Kilde	Tidsperiode	Antall ord*
Berlingske Tidende	1995-1999	132
Ekstrabladet	1995-1999	51
Politiken	1995-1999	107
Totalt		290

* Millioner ord

- <s> og </s> markerer setningsgrenser.
- Skilletegn er separert med blanke.
- Setninger som ikke er avsluttet med stort skilletegn (!?;:) er filtrert vekk (overskifter ol.).
- Setninger med mye tall er filtrert vekk. Dette er ofte resultatlistor, oppramsinger o.l.

2. Antall linjer

1-gram..... 3.274.757
2-gram..... 35.423.015
3-gram..... 114.114.098
4-gram..... 194.789.520
5-gram..... 241.732.764
6-gram..... 256.862.606

3. Filoversikt

- Filarkivet **ngram_dan_1000.zip** (45 KB, utpakket 141 KB) inneholder de 1000 mest frekvente 1-gram, 2-gram, 3-gram, 4-gram, 5-gram og 6-gram, til sammen seks tekstfiler:

ngram1-1-topp1000.txt	ngram4-1-topp1000.txt
ngram2-1-topp1000.txt	ngram5-1-topp1000.txt
ngram3-1-topp1000.txt	ngram6-1-topp1000.txt

- Filarkivet **ngram_dan.tar** (8 GB, utpakket 35 GB) inneholder følgende filer, alle rene tekstfiler:

Filnavn	Innhold
ngram1.srt ngram1-1.frk ngram1-1.srt ngram1-1-topp1000.txt	Alle 1-gram, alfabetisk liste 1-gram, frekvenssortert (fallende), frekvens > 1 1-gram, alfabetisk liste, frekvens > 1 De 1000 mest frekvente 1-gram
ngram2.srt ngram2-1.frk ngram2-1.srt ngram2-1-topp1000.txt	Alle 2-gram, alfabetisk liste 2-gram, frekvenssortert (fallende), frekvens > 1 2-gram, alfabetisk liste, frekvens > 1 De 1000 mest frekvente 2-gram
ngram3.srt ngram3-1.frk ngram3-1.srt ngram3-1-topp1000.txt	Alle 3-gram, alfabetisk liste 3-gram, frekvenssortert (fallende), frekvens > 1 3-gram, alfabetisk liste, frekvens > 1 De 1000 mest frekvente 3-gram
ngram4.srt ngram4-1.frk ngram4-1.srt ngram4-1-topp1000.txt	Alle 4-gram, alfabetisk liste 4-gram, frekvenssortert (fallende), frekvens > 1 4-gram, alfabetisk liste, frekvens > 1 De 1000 mest frekvente 4-gram
ngram5.srt ngram5-1.frk ngram5-1.srt ngram5-1-topp1000.txt	Alle 5-gram, alfabetisk liste 5-gram, frekvenssortert (fallende), frekvens > 1 5-gram, alfabetisk liste, frekvens > 1 De 1000 mest frekvente 5-gram
ngram6.srt ngram6-1.frk ngram6-1.srt ngram6-1-topp1000.txt	Alle 6-gram, alfabetisk liste 6-gram, frekvenssortert (fallende), frekvens > 1 6-gram, alfabetisk liste, frekvens > 1 De 1000 mest frekvente 6-gram