

Korpus verkkoteksteistä

- kiellentunnistus

Korpus verkkoteksteistä

- Mikä korpus?
- Mikä ”verkkokorpus”?
- Miksi?
- Miten?
 - Hakukoneet
 - Esi-/jälkikäsitteily
 - Www-sivujen lataus/siivous/analysointi
- Miten tekstin kieli tunnustetaan?

Kielentunnistus

- Tekniikat
 - N-grammit
 - HTTP-protokolla (Html-header), pelkkä URL?
 - Yleisten sanojen lista
 - Uniikit merkkijhdistelmät (huono?)
- Kaksi/monikieliset tekstit
 - Esimerkiksi kirjanorja/uusnorja, baski/espanja
 - Miten käsitellään/erotetaan
 - Kielentunnistus kappale- (jopa lause-) tasolla

.Lähteitä

- Web as Corpus:

- [The World Wide Web as Linguistic Corpus](#)
- [Can we beat Google?](#)
- [Corpus Linguistics and the Web](#)

- Kielen tunnistus (yleinen ja www)

- Natural language identification using corpus-based models
- [Language Identification in Web Pages](#)
- [Web Page Language Identification Based on URLs](#)
- [Language Identification for Multilingual Documents](#)