# Studio Mind Oy

## CEO, Pasi Karhu

Suurten tekstimassojen hyödyntäminen kontekstitietokantojen luonnissa
Utilizing large text masses (corpora) in creation of context databases

Kites Symposium 30.10.2014

STUDIOmind

# What is Context?

- Merriam-Webster online dictionary:

  : the words that are used with a certain word or phrase and that help to explain its meaning

  : the situation in which something happens
  : the group of conditions that exist where and when something happens

- Linguistic context
- Social context

STUDIOmind

# What is Context?

- My pragmatic approach in computer linguistics:

  : the words that are used with a certain word or phrase

- Statistical analysis of large bodies of text capture this

- Applying this statistics cleverly we can get computers to help to explain its meaning

- It is plausible to think that our brain does this kind of linguistical "statistics" (gathered from all our past experiences)

STUDIOmind

# Why should computers know about context?

- Without context awareness computers are stuck with the simpler but often not adequate understanding through keywords, synonym lists and thesauri

- Lexical disambiguation of polysemy, homonymy and synonymy need context awareness

- Proper context awareness opens road to advanced application like natural language communication with computers, more intelligent searches and better language translations

STUDIOmind

## Does Google do context?

- They have huge textual (and research) resources!

- Yet they don't offer contextual search

- What about Google translate:
"First World War shell explodes at former Ypres battlefield"
=> "Ensimmäisen maailmansodan kuori räjähtää entisillä Ypres taistelukenttä"
"Kiltti on miehen hame Skotlannissa"
=> "Kind is a man's skirt in Scotland"
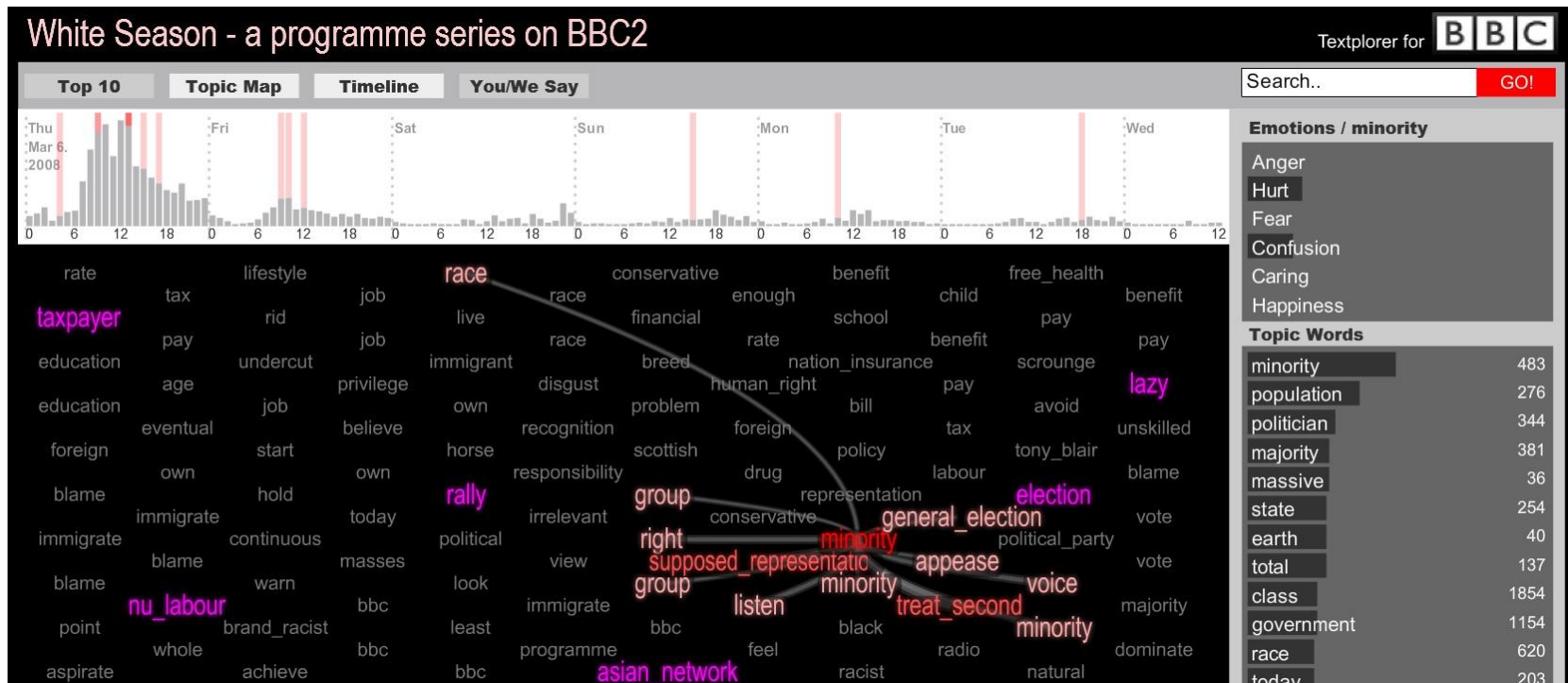"Kiltti hame Skotlannissa"
=> "Kilt Skirt in Scotland"

STUDIOmind

# Ontologies vs Statistical Context Databases?

- Ontologies are laborious to create and maintain

- Ontologies may have human bias

- Context info from large texts can be automatically created and maintained

- Statistical analysis of large text masses of interest captures all contextual relations of words as they are actually used in the texts (no human bias)

- Ontologies are more formal and precise

- Statistical context database are fuzzy by nature

- Maybe best results by combining both approaches!?

pasi.karhu@studiomind.fi          STUDIOmind

# How much text do you need?

- Some results already with 3.5 MB – example from 7000 comments online (BBC – Have Your Say)

# How much text do you need?

- With Wikipedia content (many GB) you can build a large generic context database (demo: 22 000 words)

**Some things you can do with a context database:**

- Identify most important words, sentences and paragraphs in document
  => automatic tagging (e.g. for Semantic Web)
  => creating headlines and "readers digest"

- Expand search terms

- Find correct senses of words (e.g. for translations)

- Detect different topics in a document collection or inside one document

- Classify, group and help visualize documents according to their context

- …

STUDIOmind

# Do we need more than linguistic context?

Shakespeare's Sonnet 18:

" Shall I Compare Thee to a Summer's Day? "

STUDIOmind

# THANK YOU FOR YOUR ATTENTION! QUESTIONS?

Should you need context awareness in your projects, you can reach me at:

pasi.karhu@studiomind.fi

+358 50 377 4875

Cooperation with other language technology companies also welcomed