# Integrating Treebank Annotation and User Activity Data in Translation Research

Michael Carl, Henrik Høeg Müller, Bartolomé Mesa-Lao

{mc, hhm,bm}.isv@cbs.dk

Copenhagen Business School
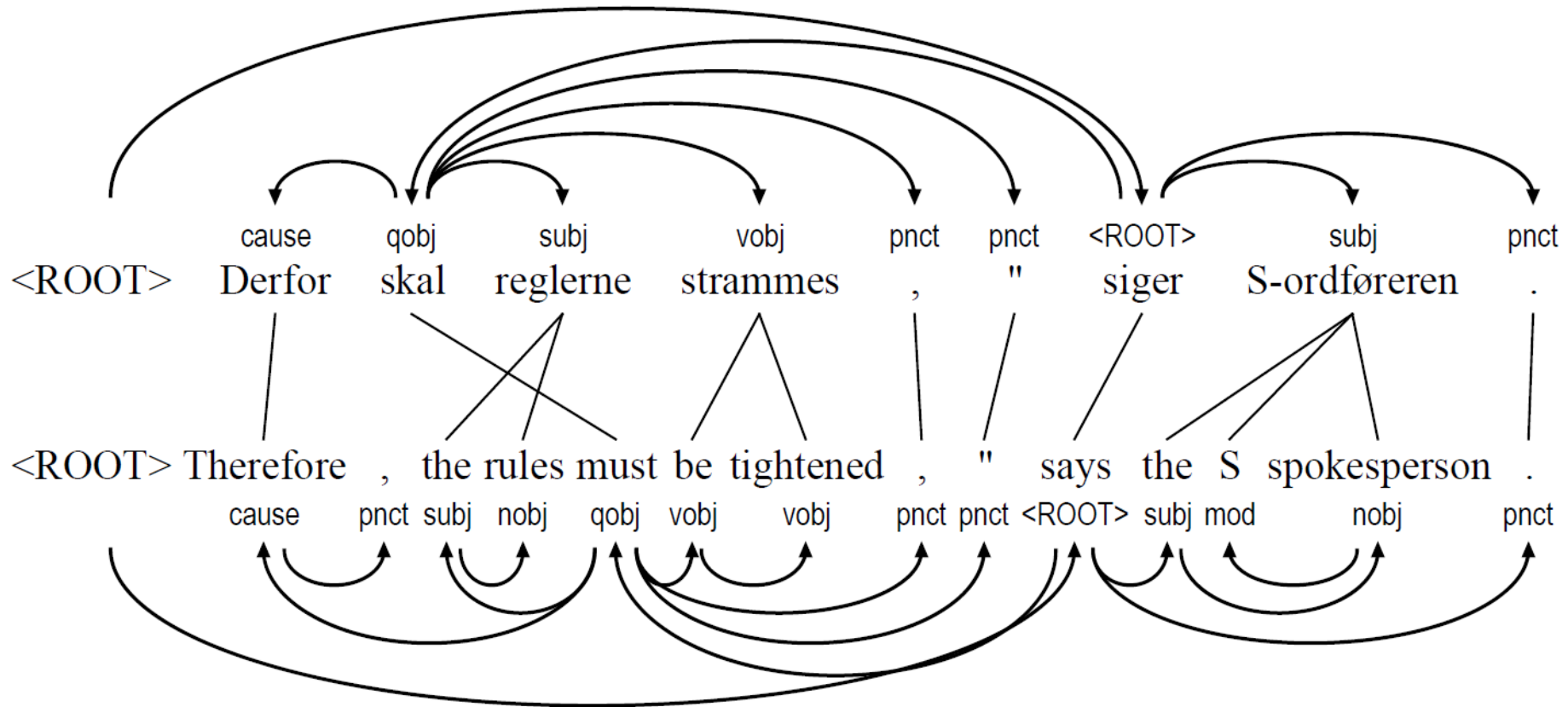Department of International Language Studies and Computational Linguistics
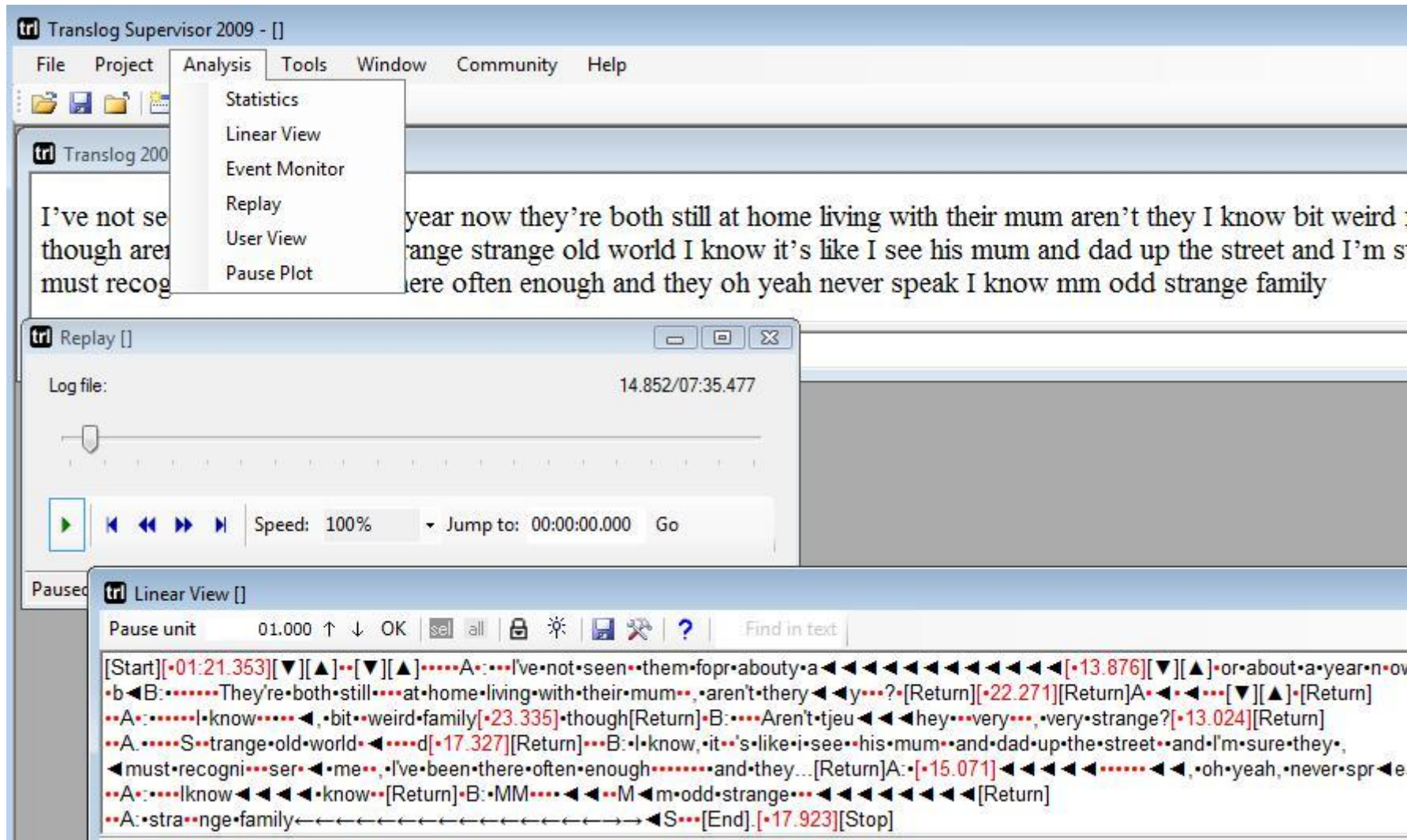
# Overview – **CRITT centre**

1. Focus on methodology for translation process research
2. Experiments with eyetracking reading
3. Keystroke background: Translog software (keylogging)
4. Integration of eyetracking and keylogging in Translog

❑ Data representation and analysis (User Activity Data)
❑ Constant development of Translog-II

Copenhagen
**Business School**
HANDELSHØJSKOLEN

# Overview – **Copenhagen Dependency Treebank**

1. NLP resource: parallel treebanks for Danish, English, German, Italian, and Spanish with 60,000 words in each language

2. Annotation levels:
   - MORPHOLOGY
   - SYNTAX
   - SEMANTICS
   - DISCOURSE
   - ANAPHORA

3. Unified account of morphology, syntax and discourse

4. Excellent basis for automatic parsers and MT-systems

# Copenhagen Dependency Treebank

# Logging behaviour with Translog-II

- ❑ *Experimental* approach - laboratory view of translation
- ❑ Data derived from *naturalistic*, but not fully natural translation events
- ❑ Translation is happening, but we can only record data occurring on the outside of the black box
- ❑ We get a very detailed, *microscopic* view of the translation process, but we do not get an inside view
- ❑ What goes on 'inside', the cognitive activity, can still only be inferred

# CRITT

## A bottom-up process: building a taxonomy of micro-behaviour (UAD) and recording dynamic interaction

❑ A bottom-up taxonomy where all the data is 'physical' (eye and finger movements in time)

❑ Method: tracing what user activity (UAD) went before a certain action and what UAD followed it (the 'history' of a process)

❑ Aim: by tracking the way in which different processes succeed one another or interact dynamically and seeing what linguistic material they operate on, we aim to build a dynamic model of translation

Copenhagen
Business School
HANDELSHØJSKOLEN

7

# Some basic assumptions

- ❑ What we record (the 'outside') correlates somehow with the 'inside'
  - ▪ E.g. what the eyes fixate is what the mind is currently processing ('the eye-mind hypothesis')
  - ▪ (eye movements are a window on the mind, but not necessarily a very clean and fully transparent window)
- ❑ Low-level 'outside' data co-vary predictably (probabilistically) with the linguistic material which is being processed
  - ▪ E.g. longer pausing (less keystroke activity) before infrequent (less familiar ) lexical items
- ❑ Higher-level, e.g. syntactic, processing phenomena can be identified from low-level data ('bottom-up')
  - ▪ E.g. high probability of regressive eye movements to antecedents of anaphors

CRITT

Copenhagen
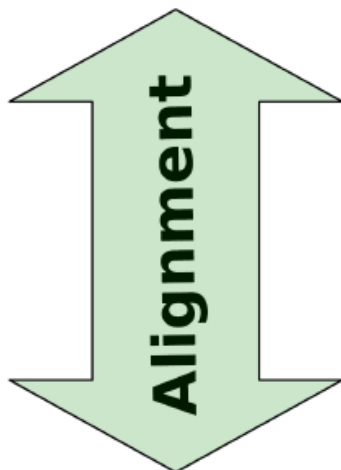Business School
HANDELSHØJSKOLEN

Although developing countries are understandably reluctant to compromise their chances of

achieving better standards of living for the poor, action on climate change need not threaten

economic development. Incentives must be offered to encourage developing countries to go the

extra green mile and implement clean technologies, and could also help minimise emissions from

deforestation. Some of the most vulnerable countries of the world have contributed the least to

climate change, but are bearing the brunt of it. Developing countries, in particular, need to adapt to

the effects of climate change. Adaptation and mitigation efforts must therefore go hand in hand.
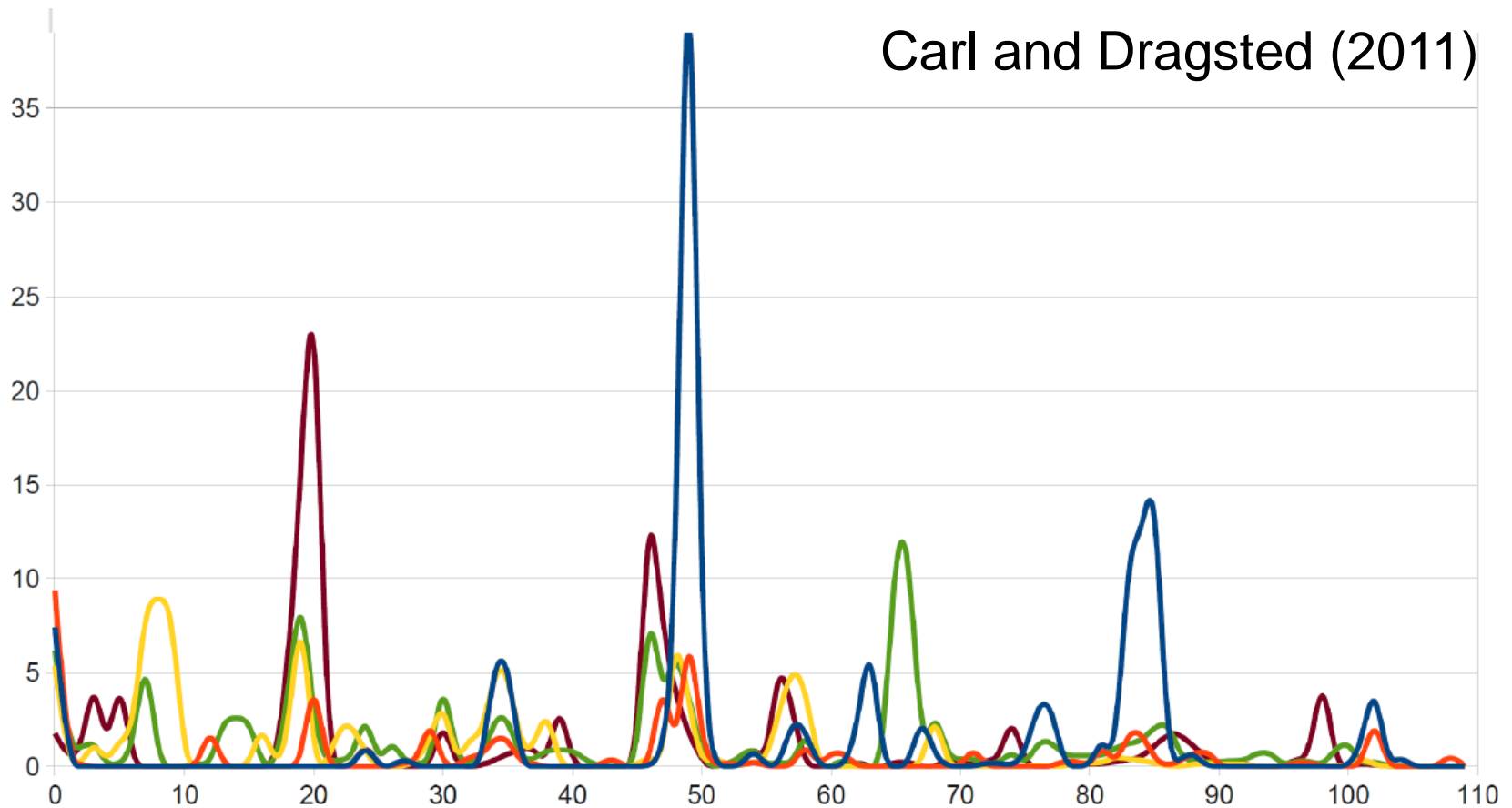
Selvom udviklingslandene forståeligt nok ikke er ivrige for at

# Visualization - Progression graphs

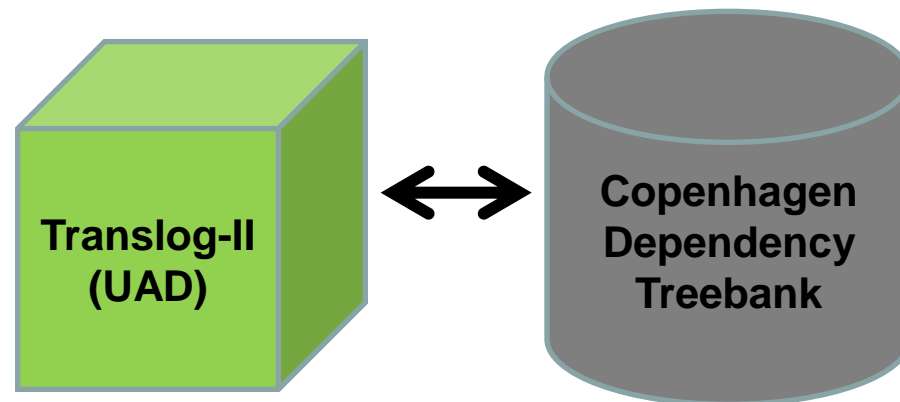# Predictability of translators' behaviour

Carl and Dragsted (2011)

# Future perspectives

❑ Using Translog-II:

  - Analysis of the typing process (pauses & typing)
  - Analysis of eye movements in the reading process
  - Analysis of coordination of reading and typing

ENRICHING THE COPENHAGEN DEPENDENCY TREEBANK WITH **USER ACTIVITY DATA** (UAD)

**Translog-II (UAD)** ↔ **Copenhagen Dependency Treebank**

# Future perspectives

❑ High expectations on the combination of UAD and CDT data to facilitate inquiries into isomorphism between:

  ▪ PAUSES & GAZE LOCATION/DURATION
  ▪ LINGUISTIC LAYERS IN THE BROADEST SENSE

❑ New insights about source text decoding, memory retrieval and encoding of textual segments

# Conclusions

❑ By mapping UAD to a parallel treebank with multilevel linguistic annotation, the possibilities of systematically analyzing correlations between **gaze fixations** / **keystrokes** and underlying **linguistic structure** of the texts are promising to uncover translation processes.

❑ We assume that the integration of these two NLP resources would allow us to correlate patterns of UAD with patterns of morphological, syntactic or discourse structure.

# Integrating Treebank Annotation and User Activity Data in Translation Research

Michael Carl, Henrik Høeg Müller, Bartolomé Mesa-Lao

{mc, hhm,bm}.isv@cbs.dk

Copenhagen Business School
Department of International Language Studies and Computational Linguistics