

## The Penn Corpora

*Albadji Jallow, Jan Reimer, Georg Hartisch*

13.06.2017

Penn Parsed Corpora of Historical English (<http://www.ling.upenn.edu/hist-corpora/>)

<b>PPCEME (Penn-Helsinki Parsed Corpus of Early Modern English)</b>	Anthony Kroch, Beatrice Santorini, and Lauren Delfs. 2004. <i>The Penn-Helsinki Parsed Corpus of Early Modern English (PPCEME)</i> . Department of Linguistics, University of Pennsylvania. CD-ROM, first edition, ( <a href="http://www.ling.upenn.edu/hist-corpora/">http://www.ling.upenn.edu/hist-corpora/</a> )
<b>PPCME2 (Penn-Helsinki Parsed Corpus of Middle English 2<sup>nd</sup> edition)</b>	Anthony Kroch and Ann Taylor. 2000. <i>The Penn-Helsinki Parsed Corpus of Middle English (PPCME2)</i> . Department of Linguistics, University of Pennsylvania. CD-ROM, second edition, ( <a href="http://www.ling.upenn.edu/hist-corpora/">http://www.ling.upenn.edu/hist-corpora/</a> ).
<b>PPCEEC (Parsed Corpus of Early English Correspondence)</b>	<i>Parsed Corpus of Early English Correspondence, tagged version</i> . 2006. Annotated by Arja Nurmi, Ann Taylor, Anthony Warner, Susan Pintzuk, and Terttu Nevalainen. Compiled by the CEEC Project Team. York: University of York and Helsinki: University of Helsinki. Distributed through the Oxford Text Archive.
<b>PPCMBE2 (Penn Parsed Corpus of Modern British English 2<sup>nd</sup> edition)</b>	Anthony Kroch, Beatrice Santorini, and Lauren Delfs. 2010. <i>The Penn-Helsinki Parsed Corpus of Modern British English (PPCMBE)</i> . Department of Linguistics, University of Pennsylvania. CD-ROM, first edition. ( <a href="http://www.ling.upenn.edu/hist-corpora/">http://www.ling.upenn.edu/hist-corpora/</a> )

What is it?

PPCEME2: 56 text samples, 1.2 million words, 1150-1500

PPCEME: 448 text samples, 1.7 million words, 1500-1720

PPCEEC: 84 letter collections, 2.2 million words = 4970 letters, 1410-1695

PPCMBE2: 2.8 million words, 1<sup>st</sup> edition 101 text samples, 2<sup>nd</sup> edition 275 text samples

Difficulties:

- It is not lemmatized → very long queries

What do we use it for?

- To analyse how grammatical phenomena have changed over time (e.g. *going to*)

## Sub-Periods of the Helsinki-Corpus

Period designation	Composition date	Manuscript date
<b>MX1</b>	unknown	1150-1250
<b>M1</b>	1150-1250	1150-1250
<b>M2</b>	1250-1350	1250-1350
<b>M23</b>	1250-1350	1350-1420
<b>M24</b>	1250-1350	1420-1500
<b>M3</b>	1350-1420	1350-1420
<b>M34</b>	1350-1420	1420-1500
<b>MX4</b>	unknown	1420-1500
<b>M4</b>	1420-1500	1420-1500
<b>E1</b>	1500-1569	
<b>E2</b>	1570-1639	
<b>E3</b>	1640-1720	

### PPCME2 set PrintStructures...

text\_id  
text\_year  
text\_period  
text\_genre  
text\_dialect\_reg # region  
text\_dialect\_dist # district (where  
available)  
text\_dialect\_hels # code as documented  
in Helsinki corpus

### PPCEME set PrintStructures...

text\_id  
text\_year  
text\_period  
text\_genre

### PPCEEC set PrintStructures...

text\_id  
letter\_id  
letter\_auth  
letter\_auth\_sex  
letter\_relToRec  
letter\_auth\_dob  
letter\_auth\_age  
letter\_rec  
letter\_rec\_sex  
letter\_relToAuth  
letter\_rec\_dob  
letter\_rec\_age  
letter\_period  
letter\_date  
letter\_authenticity