# Digital Preservation for Machine-Scale Access and Analysis

Lisa Green

Digital Preservation 2013
24 July 2013

#digpres13
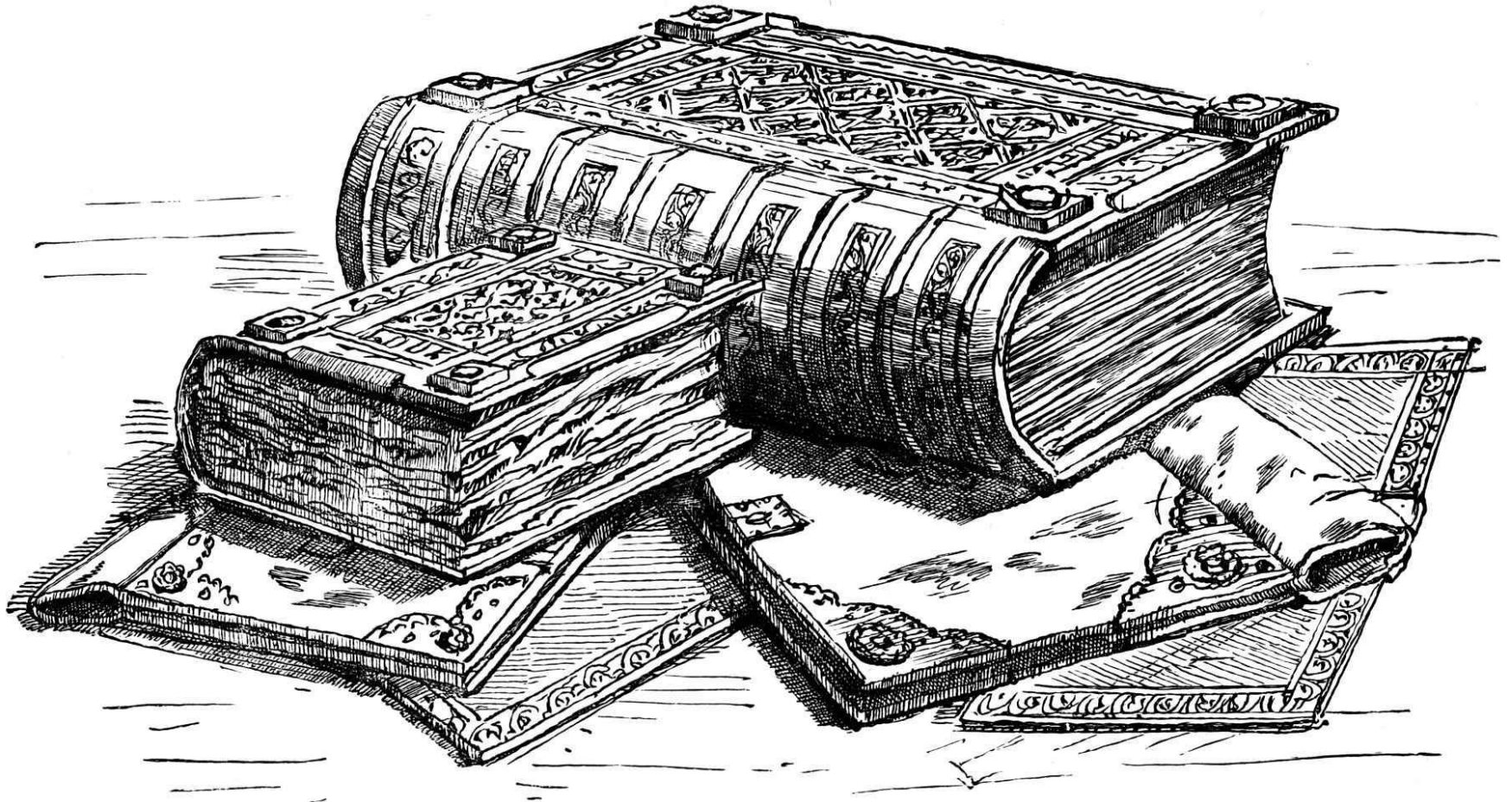
NDIIPP

#digpres13

Bundesarchiv, B 145 Bild-F031434-0006
Foto: Gathmann, Jens | 23. März 1970

NDIIPP

#digpres13

Photo license: Public Domain  Origin: http://en.wikipedia.org/wiki/File:Floppy_disk_2009_G1.jpg

#digpres13

NDIIPP

# Library of Congress Mission

To make its resources available and useful to Congress and the American people and to sustain and preserve a universal collection of knowledge and creativity for future generations.

#digpres13

# Library of Congress Mission

**To make its resources available and useful** to Congress and the American people and to sustain and preserve a universal collection of knowledge and creativity for future generations.

#digpres13

#digpres13

#digpres13

Natural Language Processing

Sentiment Analysis

Algorithms

Cluster Analysis

Machine Learning

Artificial Neural Networks

Statistics

#digpres13

PLOS | ONE

# The Expression of Emotions in 20th Century Books

Alberto Acerbi[1,2*], Vasileios Lampos[3], Philip Garnett[4], R. Alexander Bentley[1]

1 Department of Archaeology and Anthropology, University of Bristol, Bristol, United Kingdom, 2 Centre for the Study of Cultural Evolution, Stockholm University, Stockholm, Sweden, 3 Department of Computer Science, University of Sheffield, Sheffield, United Kingdom, 4 Department of Anthropology, Durham University, Durham, United Kingdom

## Abstract

We report here trends in the usage of "mood" words, that is, words carrying emotional content, in 20th century English language books, using the data set provided by Google that includes word frequencies in roughly 4% of all books published up to the year 2008. We find evidence for distinct historical periods of positive and negative moods, underlain by a general decrease in the use of emotion-related words through time. Finally, we show that, in books, American English has become decidedly more "emotional" than British English in the last half-century, as a part of a more general increase of the stylistic divergence between the two variants of English language.

#digpres13

# Figure 1. Historical periods of positive and negative moods.

#digpres13

#digpres13

#digpres13

#digpres13

# Library of Congress Mission

To make its resources available and useful to Congress and the American people and to sustain and preserve a universal collection of knowledge and creativity for future generations.

NDIIPP

# We must enable machine-scale access and analysis.

NDIIPP

# Thank you!

www.commoncrawl.org

lisa@commoncrawl.org

@commoncrawl

@boudicca

#digpres13