

Data-intensive Research and eScience

A New Possibility Space for Historical and Cultural Research

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

Dept. of History & SDU eScience Center
University of Southern Denmark

May 31, 2018

Outline

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

① Data-intensive research

- Data deluge
- Need for analytics

Data-intensive research
Data deluge
Need for analytics

② Data analytics

- Workflows
- Human analytics

Data analytics
Workflows
Human analytics
Applications
Latent semantics
Persistent dynamics
Behavioral annotation

③ Applications

- Latent semantics
- Persistent dynamics
- Behavioral annotation

Data-intensive research

Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

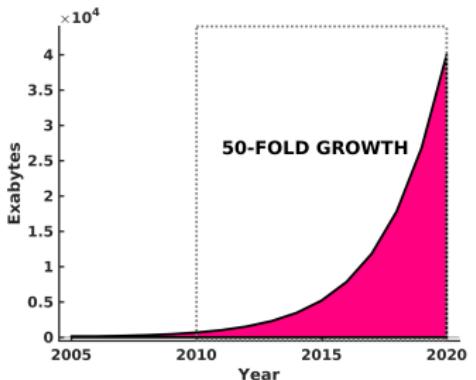
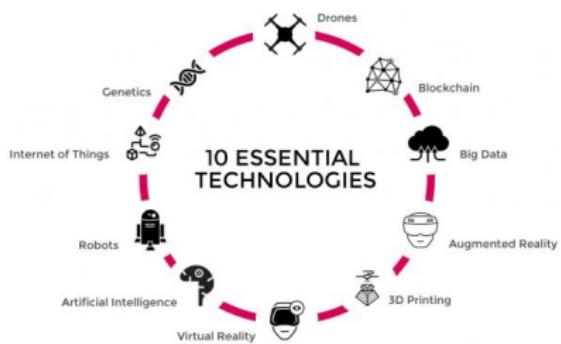
Applications
Latent semantics
Persistent dynamics
Behavioral annotation

DATA-INTENSIVE RESEARCH

The data deluge

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io



- the **data deluge** is transforming knowledge discovery and understanding in every domain of human inquiry

a large part of these data are soft and unstructured

- to get value from these data, humanities and social sciences must utilize advanced analytics solutions

Data-intensive research
Data deluge
Need for analytics
Data analytics
Workflows
Human analytics
Applications
Latent semantics
Persistent dynamics
Behavioral annotation

Humanities need analytics

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io



Figure 1: The increase in research publications & databases alone requires computational literacy. Publications related to Gospel of Marc (KJV) > 50K, ~ 16,500 words in 16 chp. on 11 p.

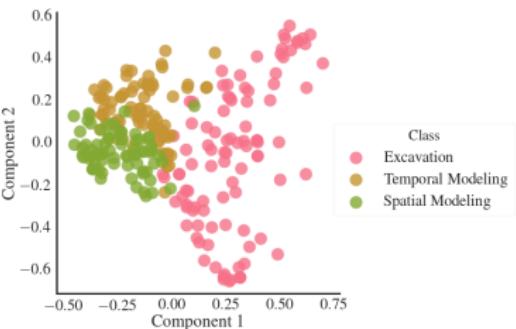


Figure 2: Advanced (human) analytics can merge, aggregate and project heterogeneous data into lower dimensional spaces that allow visual manipulation

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications
Latent semantics
Persistent dynamics
Behavioral annotation

Data-intensive research
Data deluge
Need for analytics

Data analytics

Workflows
Human analytics

Applications

Latent semantics
Persistent dynamics
Behavioral annotation

DATA ANALYTICS

SSH default workflow

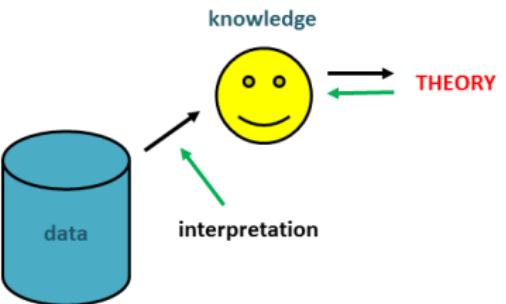
Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
[knielbo.github.io](https://github.com/knielbo)

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

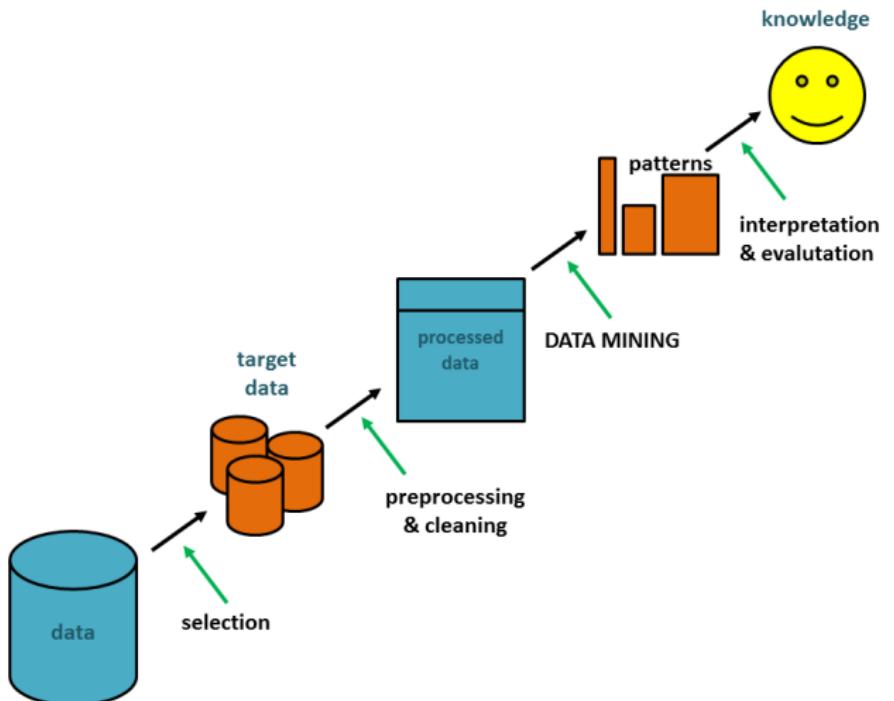
Applications
Latent semantics
Persistent dynamics
Behavioral annotation



KDD workflow

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
[knielbo.github.io](https://github.com/knielbo)



Data-intensive research
Data deluge
Need for analytics
Data analytics
Workflows
Human analytics
Applications
Latent semantics
Persistent dynamics
Behavioral annotation

Human analytics - computationally empowering humanities

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications
Latent semantics
Persistent dynamics
Behavioral annotation

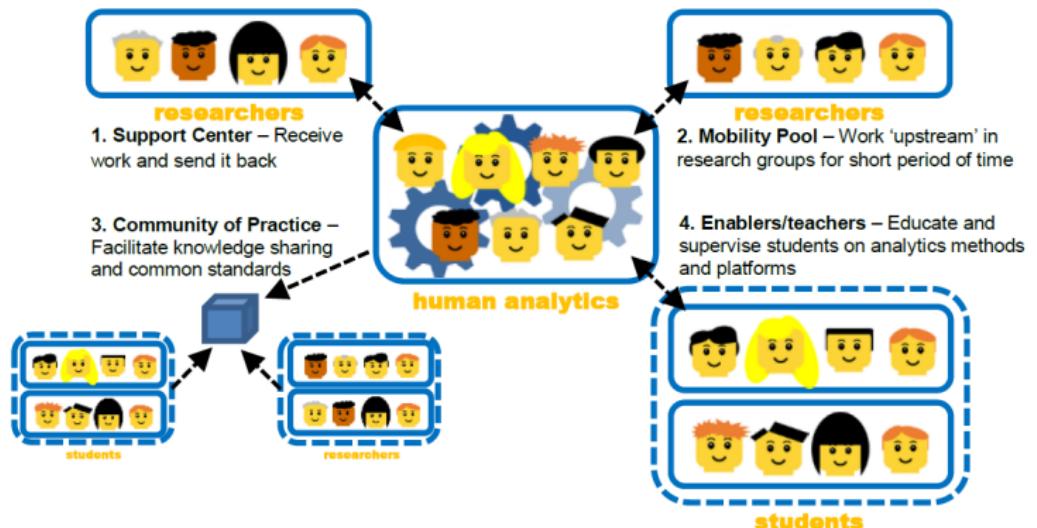


Figure 3: © SDU we have established an analytics group (2.5 FTEs) with four modes of empowerment: **support**, **mobility**, **community of practice**, and **enablers**

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications

Latent semantics
Persistent dynamics
Behavioral annotation

APPLICATIONS

Chinese philosophy

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications

Latent semantics
Persistent dynamics
Behavioral annotation

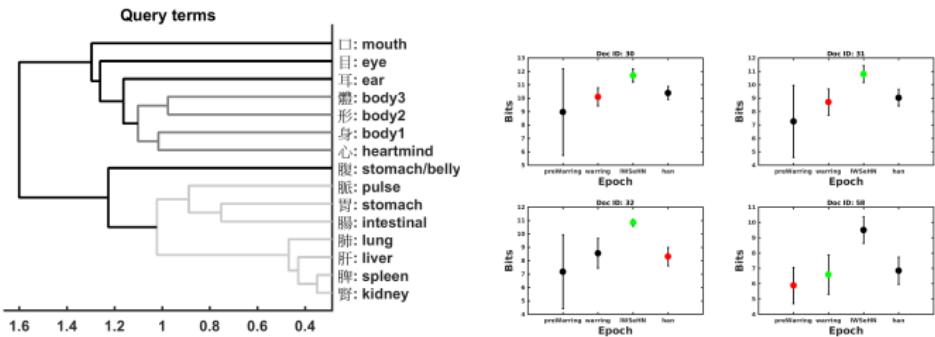


Figure 4: Philosophers and sinologists have been debating the existence of **mind-body dualism** in classical Chinese philosophy. With domain experts, we identified a **hierarchical dualistic space** based on latent semantic models. One model (LDA) was further utilized to predict class of origin for **controversial texts** slices.

Slingerland, E., Nichols, R., Nielbo, K., & Logan, C. (2017). The Distant Reading of Religious Texts: A Big Data Approach to Mind-Body Concepts in Early China. *Journal of the American Academy of Religion*, 85(4), 985–1016.

Nichols, R., Slingerland, E., Nielbo, K., Bergeton, U., Logan, C., & Kleinman, S. (2018). Modeling the Contested Relationship between Analects, Mencius, and Xunzi: Preliminary Evidence from a Machine-Learning Approach. *The Journal of Asian Studies*, 77(01), 19–57.

Medieval history and novelty detection

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

- historians debate historical transitions
- Saxo's *Gesta Danorum* c. 1200 CE history of the Danish royal dynasty
- transition between book 8 or 9?
- transition point or gradual?
- traditional word-level representation is ambivalent
- latent semantic model was trained over sentence windows
- change detection and recurrence plot used to identify phase transition centered in book 9

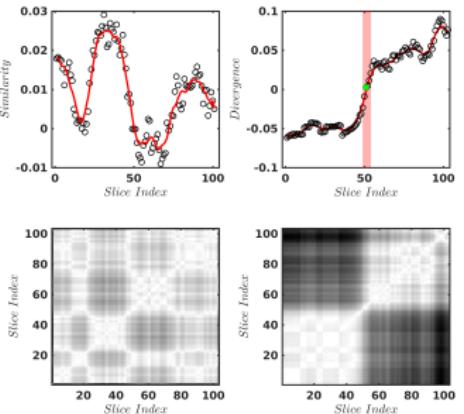


Figure 5: Cosine distance and KLD for TD high-rank vector space and guided LDA model respectively.

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications
Latent semantics
Persistent dynamics
Behavioral annotation

%domestic events indicates populism

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

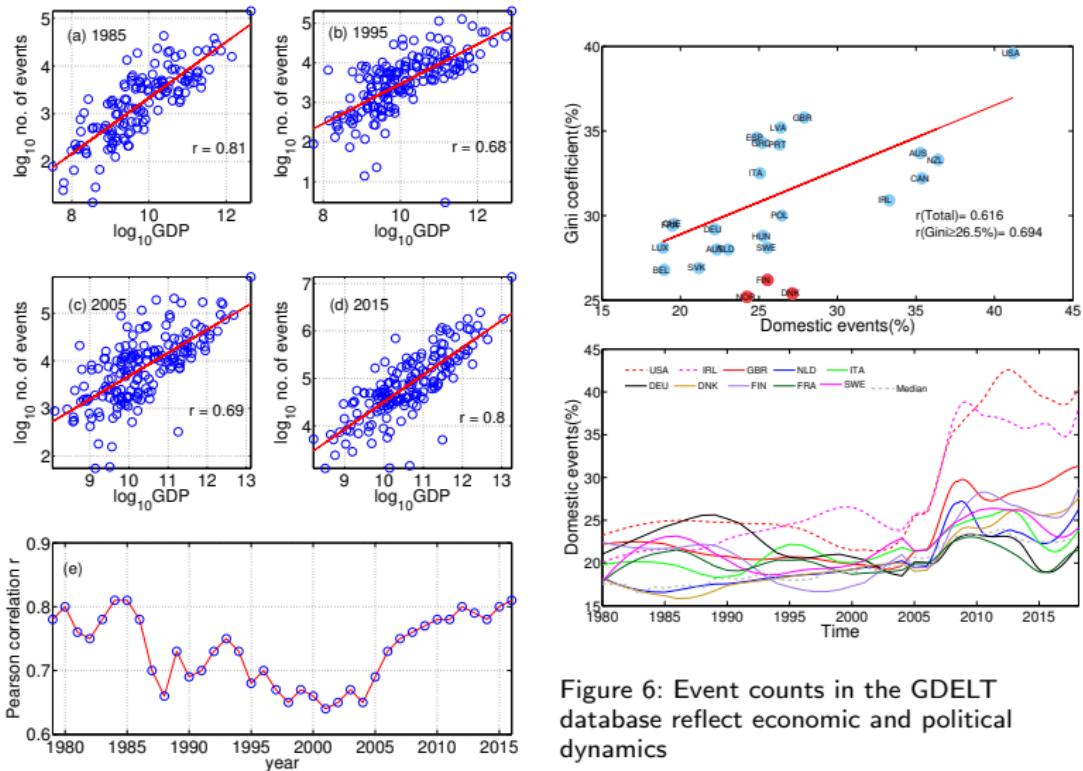


Figure 6: Event counts in the GDELT database reflect economic and political dynamics

Consumer goods and cultural memory

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications
Latent semantics
Persistent dynamics
Behavioral annotation

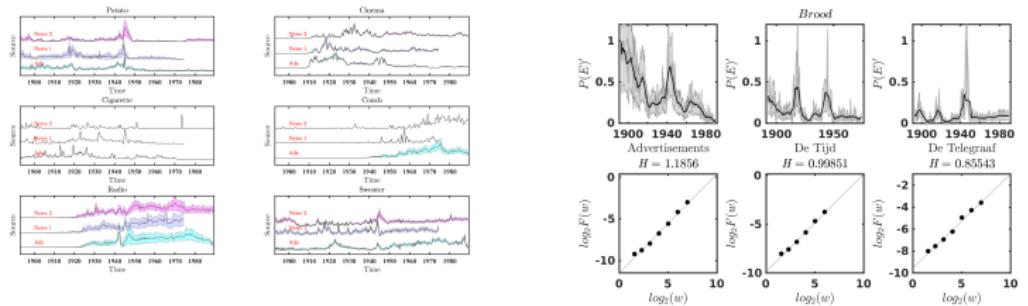


Figure 7: Historians and media researchers theorize about the causal dependencies between public discourse and advertisement. Causal modeling of keyword frequencies (from seedlists) indicated that for some categories 'ads shape society', while other categories merely 'reflect'. Advertisements show a faster decay (on-off intermittent behavior) than public discourse (long-range dependencies) - a proxy for cultural memory.

Literary optimality

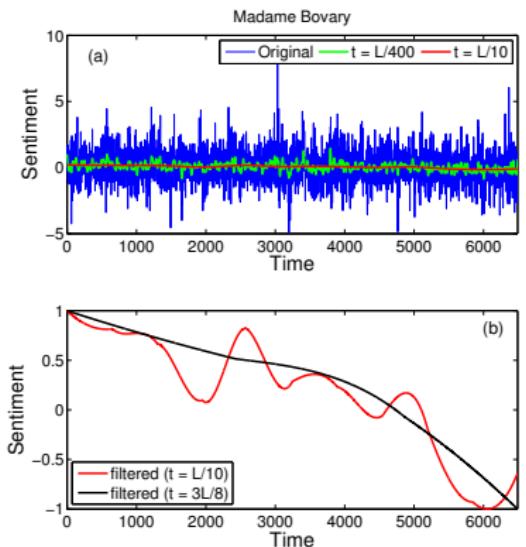
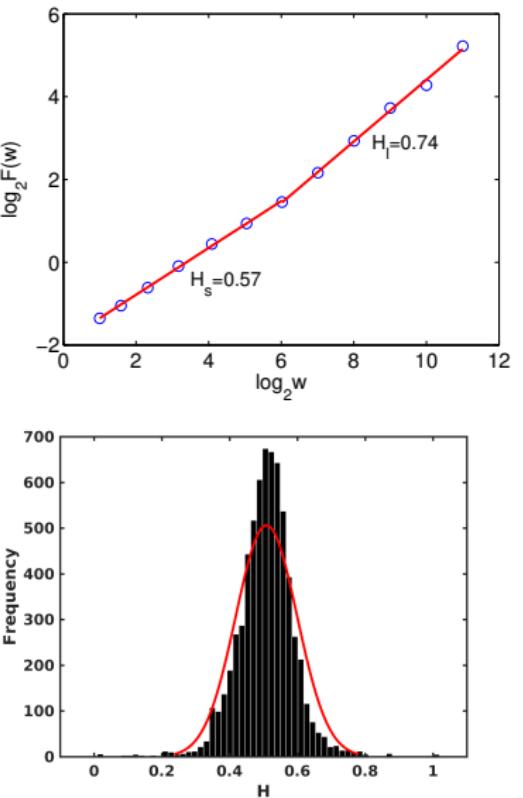


Figure 8: sentiment analysis and adaptive filtering reconstructs narrative vectors that reflect the reader experience. Particular fractal scaling-range, $0.6 < H \leq 0.8$, indicates literary optimality.



Author change points

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
[knielbo.github.io](https://github.com/knielbo)

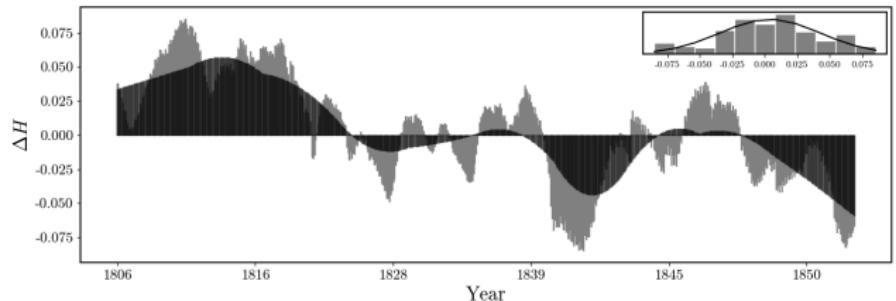


Table 1: Dominant Dynamic in the Phases of N.F.S. Grundtvig's Writings

Time period	Age of onset	$H(X)$	Behavior	Profile
1806-1826	23	$H > 0.5$	<i>persistent</i>	theoretician
1826-1839	43	$H \approx 0.5$	<i>short memory</i>	pragmatic
1839-1845	56	$H < 0.5$	<i>anti-persistent</i>	breakthrough
1845-1848	62	$H \approx 0.5$	<i>short memory</i>	disease
1849-1872	65	$H < 0.5$	<i>anti-persistent</i>	politician

Data-intensive research
Data deluge
Need for analytics
Data analytics
Workflows
Human analytics
Applications
Latent semantics
Persistent dynamics
Behavioral annotation

Dynamic author profiling

Data-intensive
Research and eScience

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications
Latent semantics
Persistent dynamics
Behavioral annotation

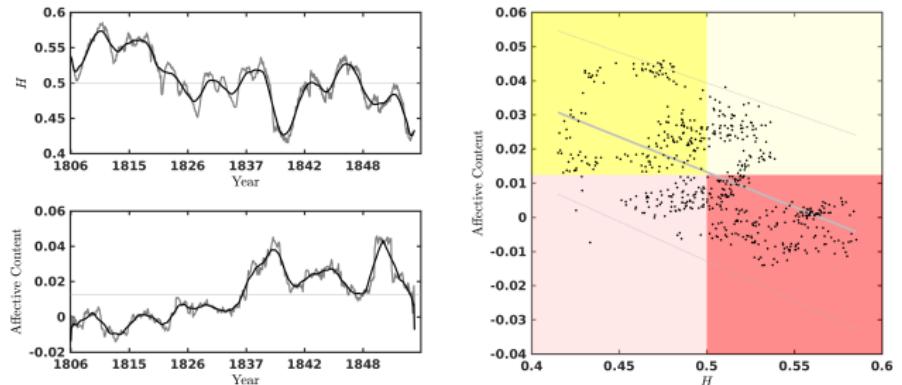


Figure 9: Combining persistent entropic trends with sentiment analysis and causal modeling, we can study “**the tormented artist**” phenomena in intellectual history.

Decision support for OCD

Data-intensive
Research and eScience

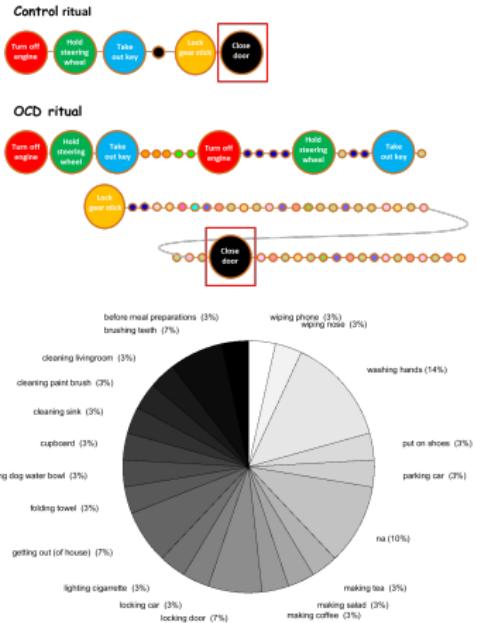


Figure 10: Event logging database annotated with Observer XT for OCD, comorbid, and control.¹

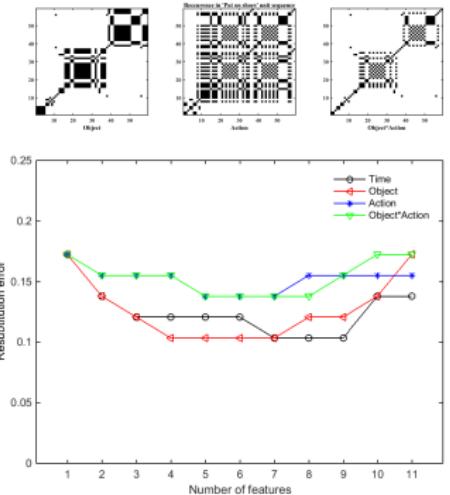


Figure 11: Binomial classifier (OCD vs. control) on unseen data.

¹Zor, R., Hermesh, H., Szechtman, H., & Eilam, D. (2009). Turning order into chaos through repetition and addition of elementary acts in obsessive-compulsive disorder (OCD). World Journal of Biological Psychiatry, 10(4.2), 480–487.

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

Data-intensive research
Data deluge
Need for analytics
Data analytics
Workflows
Human analytics
Applications
Latent semantics
Persistent dynamics
Behavioral annotation

The structure of dance

Data-intensive
Research and eScience

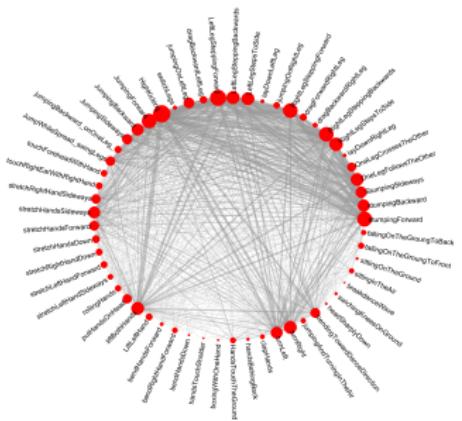


Figure 12: Event logging database annotated with Observer XT for ritualistic dance during **Zulu weddings** in two generations with male/female participants.

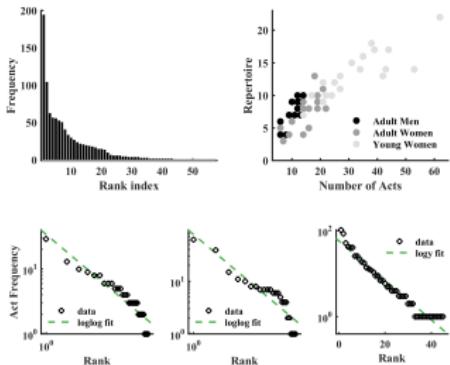


Figure 13: Movement inventory follow Heaps' and Zipf's laws known from other **communicative systems**.

Data-intensive research
Data deluge
Need for analytics

Data analytics
Workflows
Human analytics

Applications
Latent semantics
Persistent dynamics
Behavioral annotation

Kristoffer L Nielbo
knielbo@sdu.dk
knielbo.github.io

THANK YOU

knielbo@sdu.dk

knielbo.github.io

slides: http://knielbo.github.io/files/our_muse.pdf

& credits to

Max R. Echardt and Katrine F. Baunvig, datakube, University of Southern Denmark, DK

Jianbo Gao and Bin Liu, Institute of Complexity Science and Big Data, Guangxi University, CHN

Melvin Wevers, DH Lab, KNAW Humanities Cluster, NL

Culture Analytics @ Institute of Pure and Applied Mathematics, UCLA, US