# Fortschritt-Berichte VDI

VDI
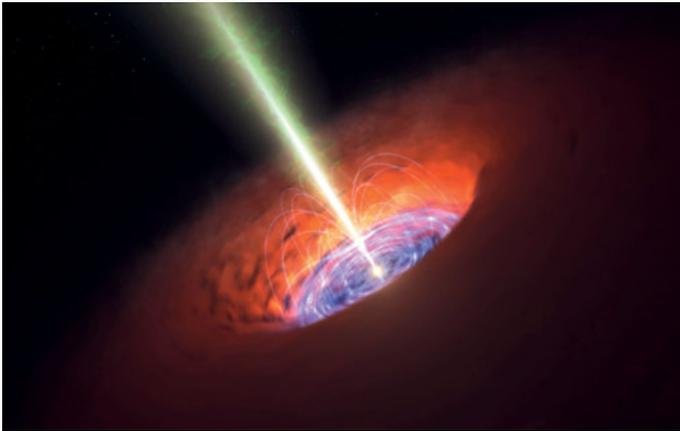
Herwig Unger,
Mario M. Kubek

# Theory and Application of Text-representing Centroids

**FernUniversität in Hagen**

**Schriften zur Informations-
und Kommunikationstechnik**

# Theory and Application of Text-representing Centroids

Herwig Unger and Mario M. Kubek

Email: kn.wissenschaftler@fernuni-hagen.de
Website: https://www.fernuni-hagen.de/kn/

# Fortschritt-Berichte VDI

## Theory and Application of Text-representing Centroids

Herwig Unger, Mario M. Kubek
**Theory and Application of Text-representing Centroids**
Fortschr.-Ber. VDI Reihe 10 Nr. 863. Düsseldorf: VDI Verlag 2019.
152 Seiten, 49 Bilder, 10 Tabellen.
ISBN 978-3-18-386310-5, ISSN 0178-9627,
€ 57,00/VDI-Mitgliederpreis € 51,30.

Centroid terms are single, descriptive words that semantically and topically characterise text documents and thus can act as their very compact representation in automated text processing tasks that strongly rely on the semantic similarity of texts. Algorithms to classify and cluster them make use of this information. In this book, the novel, brain- and physicsinspired concept of centroid terms is introduced and deeply discussed. Furthermore, their unique properties and practical usage in major natural language processing and text mining tasks are covered. In this regard, a new graph-based method for their fast calculation is presented as well. In contrast to methods relying on the bag-of-words model, the derived centroid distance measure can uncover a topical relationship between texts even when their wording differs. As centroid terms can also represent short texts, the presented first fully integrated, P2P-based web search engine, called "WebEngine", therefore makes heavy use of centroid terms when interpreting queries and forwarding them to peers with matching documents, represented by their own centroid terms.

Schriften zur Informations- und Kommunikationstechnik
Herausgeber:
Wolfgang A. Halang, ehemaliger Lehrstuhl für Informationstechnik
Herwig Unger, Lehrstuhl für Kommunikationstechnik
FernUniversität in Hagen

# Contents