

Betsy Van der Veer Martens

Keywords In and Out of Context

Synthesis Lectures on Information Concepts, Retrieval, and Services

Series Editor

Gary Marchionini, School of Information and Library Science, The University of North Carolina at Chapel Hill, Chapel Hill, NC, USA

This series publishes short books on topics pertaining to information science and applications of technology to information discovery, production, distribution, and management. Potential topics include: data models, indexing theory and algorithms, classification, information architecture, information economics, privacy and identity, scholarly communication, bibliometrics and webometrics, personal information management, human information behavior, digital libraries, archives and preservation, cultural informatics, information retrieval evaluation, data fusion, relevance feedback, recommendation systems, question answering, natural language processing for retrieval, text summarization, multimedia retrieval, multilingual retrieval, and exploratory search.

Betsy Van der Veer Martens

Keywords In and Out of Context

Betsy Van der Veer Martens
University of Oklahoma
Norman, OK, USA

ISSN 1947-945X ISSN 1947-9468 (electronic)
Synthesis Lectures on Information Concepts, Retrieval, and Services
ISBN 978-3-031-32529-8 ISBN 978-3-031-32530-4 (eBook)
<https://doi.org/10.1007/978-3-031-32530-4>

© The Editor(s) (if applicable) and The Author(s), under exclusive license to Springer Nature
Switzerland AG 2023

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

To my husband Jon, with “love” as keyword

Preface

Why this book about the “keyword”? Today, keywords and their multiple uses serve to bridge between the humanities and technology, between librarianship and information science, between the searcher and the web resource, between the search engine and the advertiser, between the social influencer and the audience, between the political strategist and the voter turnout, and between our contemporary world and that of our ancient predecessors. The keyword in its various guises (key word, concept symbol, hashtag, and search term) can point not only to text and other forms of media, but to associated ways of thinking and acting based on specific words that we may consider “key.”

This project is an effort to explore the rich history of the keyword from its earliest manifestations (long before it appeared anywhere in Google Trends or library cataloging textbooks) in order to illustrate its implicit and explicit mediation of human cognition and communication processes, from its deictic origins in primate and proto-speech communities, through its semiotic and symbolic instantiation in various physical artifacts and structures, through its development within oral traditions, through its initial appearances in numerous graphical forms, through its workings over time within a variety of indexing traditions and technologies, to its role in search engine optimization and social media strategies, to its potential as an element in the slowly emerging semantic web as well as in multiple voice search applications. The purpose of the book is to synthesize different perspectives on the significance of this often-invisible intermediary, both in and out of the library and information science context, and to understand how it has come to be so embedded in our daily life.

Norman, USA

Betsy Van der Veer Martens

Contents

1	Representation, Reference, Relevance, and Retention	1
1.1	Defining the Keyword	2
1.2	Grounding the Keyword	2
1.3	Representation	5
1.4	Reference	6
1.5	Relevance	7
1.6	Retention	8
1.7	Grounds for Key Concepts	8
	References	10
2	Signals, Semiotics	15
2.1	Introduction to Semiotics	15
2.2	Studying Signals	18
2.3	Gestural Repertoires	20
	References	23
3	Proto-Signs, Proto-Words	29
3.1	Proto-World	30
3.2	Theories of Language Evolution	32
3.3	Communication-Centric Theories of Language Evolution	32
3.4	Culture-Centric Theories of Language Evolution	35
	References	41
4	Philologies, Philosophies, Pragmatics	47
4.1	Philology and Plato	47
4.2	Historical Linguistics	49
4.3	Reference and Relevance	52
4.4	The Problems of Pragmatics	55
4.5	Expanding the World of Words	57
	References	59

5	Rites, Religions	65
5.1	Origins of Symbols	66
5.2	Oral Traditions	67
5.3	Ritual Writings	70
	References	74
6	Writing, Indexing	79
6.1	The Rise of Writing	79
6.2	The Importance of Interpretation	82
6.3	The Literacy Revolution	84
6.4	The Language of Science and the Science of Language	88
	References	91
7	Progress, Public	97
7.1	The Growth of Governments and the Push Toward Progress	97
7.2	Cooperation and Competition	100
	References	106
8	Discovery, Retrieval	111
8.1	Indexing Essentials	111
8.2	Library Discovery Systems	112
8.3	Information Discovery Systems	114
8.4	The TREC Conferences	121
	References	122
9	Databases, Search Engines	127
9.1	Commercial Databases	127
9.2	Search Engines	128
9.3	Effects on LIS Practice and Research	132
	References	136
10	Algorithms, Users	141
10.1	Search Engine Results Pages	141
10.2	Computational Advertising	144
10.3	Search Engines and the Public Sphere	146
10.4	The Future of Search	148
	References	150



Representation, Reference, Relevance, and Retention

1

Like many of the words that matter most, that tell us most about our intellectual and material life and about our cognitive and perceptual habits, 'keyword' hides in plain view. —Michael Leja, "Keyword" (2009)

Abstract

The long history of keywords and their predecessors as semiotic, symbolic, and semantic pointers to key concepts over time is introduced. This chapter describes current findings on four sensory specifics that are generally not considered as being aspects of library and information science but that are keywords which ground the discipline both physically and conceptually: that is, vision for representation, voice for reference, hearing for relevance, and memory for retention.

While words in any natural language can serve as symbolic and semantic tools in individual and social cognition, due to their ability to mobilize both abstract and concrete concepts in representation and reference and even to immobilize these in various ways for retention and retrieval over time, some of these verbal tools appear to be particularly useful as keys to communication, maintaining and retaining their significance for larger groups of people and for longer periods of time. However, as Leja (2009) observes above, the presence of these “key” words, especially in their functions as “keywords,” is largely taken for granted.

1.1 Defining the Keyword

The Oxford English Dictionary defines “keyword” as either “a. A word that serves as the key to a cipher or code” or “b. A word or idea that serves as a solution or explanation for something; a word, expression, or concept of particular importance or significance.” Rosenberg (2021) in discussing these two definitions suggests that they are best represented today by the prevalence of keyword searching in search engines and by the prevalence of polarizing terms in social discourse. His genealogy of the modern trajectory for both definitions can be traced to 1958, the year in which IBM engineer Hans-Peter Luhn (1958) published his method for automatically extracting and indexing “significant” words from scientific and technological articles and in which cultural theorist Raymond Williams (1958) published his initial analysis of how particular words as used by individual writers were the key to analyzing changes in cultural and social mores. These two trajectories, the technological and the cultural, have coincided in today’s information environment, to the point that, as Rosenberg says, “What makes *keyword* such a powerful idea is precisely the ambiguity of the relationship that it mediates between what is informative and what is significant, a conundrum of our time if there ever was one” (Rosenberg 2021, p. 121).

Bernard, in his history of the hashtag, which he calls the latest incarnation of the keyword, dismisses most of these earlier incarnations, saying that “without a doubt the category of ‘the keyword’ had occupied a rather inconspicuous place prior to the twenty-first century...Today, every Twitter feed and Instagram post provides further testimony to the collective indexing or ‘keywording’ of the world” (Bernard 2019, p. 2).

In their encyclopedic examination of “keyword” and its related terms (term, index term, free-text term, Uniterm, heading, subject heading, descriptor, concept symbol, tag, word, stopword, N-gram, and keyphrase), Lardera and Hjørland (2021) provide an in-depth intellectual background for the keyword in both library and information science (LIS) theory and practice. Stubbs (2010, p. 25) has done the same for the discipline of linguistics, arguing that “keyness is a textual matter” because “keywords are words which are significantly more frequent in a sample of text than would be expected, given their frequency in a large general reference corpus.”

Nevertheless, the keyword has its origins in a much longer history than proposed by these authors, so the rest of this chapter will explore the grounds on which our understanding of “key” words should begin.

1.2 Grounding the Keyword

Clearly, before there can be “keywords,” there must be a conceptual and communicative infrastructure in which any such coded or clear reference to or perceived relevance of “something” (or, indeed, “anything”) can be meaningful enough to be memorable. Bruner

referred to all of these as “routes to reference” (Bruner 1998). It is well accepted by now that evolutionary approaches to communication and information are necessarily intertwined, as both survival and reproduction require some successful internal and external communication of information at every level of taxa, from the lower bacteriological levels (Lyon 2015) to the higher zoological ones (Hoffecker 2013). In particular, the question of the development of the human “faculty of language” (Hauser et al. 2002) can no longer be considered in isolation from epigenetic (Gokhman et al. 2016), genetic (Graham et al. 2015), neural (Konopka and Roberts 2016), and other environmental (Greenhill 2016) factors.

Changaux and his colleagues (2021) have proposed that seemingly minor changes in the human genome since our fairly recent evolution from nonhuman primates can explain fundamental features of human brain connectivity, especially the tripling in size of the global neural architecture within the original primate brain, resulting in a larger number of neurons and areas and the increased modularity, efficiency, and differentiation of cortical connections. “The combination of these features with the developmental expansion of upper cortical layers, prolonged postnatal brain development, and multiplied nongenetic interactions with the physical, social, and cultural environment gives rise to categorically human-specific cognitive abilities including the recursivity of language. Thus, a small set of genetic regulatory events affecting quantitative gene expression may plausibly account for the origins of human brain connectivity and cognition” (2021, p. 2425).

Worden (2022) suggests that the traditional notion of language evolution through natural selection alone cannot account for the fact that the energy expenditure for our language-enabled brain (roughly 20% of our metabolic requirements) is much too high in comparison with that of a simpler brain with primitive language capabilities and smaller metabolic costs, which would be both more efficient and entirely adequate for our original survival needs. He proposes instead that both natural selection and sexual selection played a role in the evolution of language and intelligence, probably at different times. Specifically, he theorizes that early language was driven by natural selection to facilitate within-group collaboration. Early forms of information exchange, probably developed as sounds and gestures in various hunting and gathering activities over time, along with an emerging theory of mind, began to serve as markers for this superior intelligence, which ultimately played a critical role in sexual selection for early man, as the qualities of empathy and leadership it can embody are attractive to both peers and potential mates, thus precipitating the unique refinement of pragmatics, the development of spoken symbols, and the construction of syntax which was eventually to become modern language. As it does today, language acts as the main display mechanism for intelligence and also determines which “keywords” will become critically important in particular contexts, whether it is the term for a prey animal or the name of a political party.

Within library and information science, this is consistent with today’s “informational” turn in which, despite some skepticism, such information-oriented scholars as Bates (2022), Beynon-Davies (2011), Brier (2010), Madden (2004), O’Connor (1996), Shah

(2023), Spink (2010), and Stonier (1997) continue their interdisciplinary investigations of what may quite reasonably be termed the evolution of information research. These initiatives tend to stress the continuities and similarities among different forms of informationally oriented cognitive systems over space and time and to take a much wider perspective than usual. They also tend to support the utility of a broader approach to the information problematic, such as that posed by and through “keywords.”

Thinking about keywords as part of that problematic can raise central issues of representation (that is, presentation and organization of data or information) and reference (that is, meaning as intended by the speaker or writer), relevance (that is, meaning as understood by the listener or reader), and retention (that is, the varied forms of internal and external memory that may also be archived in both individual and social forms): all of these are fundamental aspects of LIS. In the absence of words, would any of these exist? Conversely, in their absence, might words still exist? This chapter explores some of the current findings on systems, symbols, and speech that must necessarily (though invisibly) ground any concept of “keywords.”

The biological systems necessary for representation, reference, relevance, and retention are usually taken for granted, as human beings are so used to our visual, oral, aural, and retentive processing capabilities that it seldom occurs to us to wonder how these are affecting what we see, say, hear, and remember, other than perhaps in thinking in terms of extending these senses through novel technologies. Nevertheless, a knowledge of some sensory specifics and their embodiment may be helpful in grounding this discussion, especially those related to the so-called “symbol grounding problem,” in which the question is how any one thing can be connected to a meaningful interpretation of that thing. This has been neatly expressed both by Searle (1980, p. 424) who observed “Of course the brain is a digital computer. Since everything is a digital computer, brains are too. The point is that the brain’s causal capacity to produce intentionality cannot consist in its instantiating a computer program, since for any program you like it is possible for something to instantiate that program and still not have any mental states. Whatever it is that the brain does to produce intentionality, it cannot consist in instantiating a program since no program, by itself, is sufficient for intentionality”) and by Harnad (1990, p. 335) who queried “How can the semantic interpretation of a formal symbol system be made *intrinsic* to the system, rather than just parasitic on the meanings in our heads?”

As Barsalou (2016, p. 1129) put it, “To a large extent, grounding concerns itself with the grounding problem raised initially by Searle and Harnad which asks how amodal symbols, specifically, and cognition, more generally, are linked to the modalities, body, and environment. In a review of research on grounding, Barsalou argued that researchers have attempted to ground concepts and cognition by establishing their relations with modality-specific systems, the body, the physical environment, and the social environment. ... Thus, at a general level, grounding simply refers to programmatically studying cognition in new ways. Rather than studying cognitive mechanisms in isolation, establish their relations with the contexts in which they are embedded and on which they depend. At more

specific levels, grounding refers to establishing specific accounts of how cognitive processes in the brain utilize the modalities, the body, and the environment. It does *not* mean reducing concepts and cognition to anything, including sensory-motor mechanisms.”

1.3 Representation

Vision, for instance, enables us to *represent* the world. Pylyshyn (2000, p. 197) explains that “Representations are the basic building blocks of cognitive explanations of human behavior.... [and] function in the same way as descriptions: they use the conceptual resources of the mind to encode properties of the world in much the same way as language uses words... [but] a conceptual description alone (what Bertrand Russell called a ‘definite description’) is inadequate for encoding certain types of [physical] knowledge... such as finding one’s way home. The most primitive contact that the visual system makes with the world (the contact that precedes the encoding of any sensory properties) is a contact with what have been termed visual objects or proto-objects.” Relatedly, Ballard and his colleagues (1997) found that very minute eye and hand movements are linked by processes underlying these elemental perceptual events through “deictic coding” to working memory, at time scales of approximately 1/3 of a second, and play an essential role in the brain’s symbolic computations of “embodied” representations. Pitcher and Ungerleider (2021) have proposed that on the lateral brain surface in the primate visual cortex, in addition to the ventral visual pathway, which computes the identity of an object and the dorsal visual pathway, which computes the location of an object and actions related to that object, there exists a third visual pathway which computes the actions of moving faces and bodies and is apparently specialized for the dynamic aspects of social perception.

Symbols are represented in our brains at different levels of complexity: at the initial, simplest level, as physical entities, in the corresponding primary and secondary sensory cortices, as conceptual ones. Symbols, however, no matter how simple their surface forms may appear, evoke higher order multifaceted representations that are implemented in distributed neural networks spanning a large portion of the cortex. These internal states that reflect our knowledge of the meaning of symbols are what we call semantic representations (Borghesani and Piazza 2017). Viganò and his colleagues (2021) showed that this categorization within the brain took place in at least three representational stages: first, the sensory regions process the features relevant for categorization, the left angular gyrus integrates the different sensory features into unique object identities, connecting them to the correct name, and the hippocampus encodes the abstract associative rule.

The question of whether there might be a “language of thought” cognitive coding that is separate from natural language, however, is still open, as opined by Mandelbaum and his colleagues (2022): “Recent advances in deep neural networks appear to suggest that there is no need for psychological models beyond ones that posit links between neuron-like nodes. But while Artificial Intelligence (AI) research has moved away from

transparently interpretable, richly structured internal representations, advances in many disparate areas of cognitive science suggest otherwise. Evidence from animal and infant cognition, Bayesian computational cognitive science, unconscious reasoning, and visual cognition suggests that the mind traffics in representations couched in an amodal code with a language-like structure.”

1.4 Reference

Like other mammals, our common way of sharing a representation is by making a vocal *reference* to it. (Pointing, the other referential method common in humans, is uncommon among mammals, even among the great apes, though a few have been occasionally observed using whole-hand gesturing to indicate a desired object from a human companion).

The mammalian voice production organ has three subsystems: the pulmonary system, which supplies power through the lungs, a sound generation system, typically the larynx, and a sound modifier system, the (pharyngeal, oral, and/or nasal) vocal tract (Herbst 2016). It was once believed that a descended larynx was uniquely human, but it has now been found in deer, for instance, though the human vocal tract still seems to be similar only to those of the Neanderthals and the Denisovans (Dediu et al. 2021). Regardless of languages and contexts, the amplitude modulation of the speech signal for humans consists of a rhythm that ranges between 3 and 8 Hz, while the vocalizations and facial expressions of monkeys and apes also have this rhythmic structure (Zhang and Ghazanfa 2020). Human infants are attuned to this rhythm even prior to birth, which helps to accelerate their process of language acculturation and accumulation (Ghio et al. 2021).

Although the abundance of sounds found in the world’s languages has been thought to have been fixed by biological constraints since the emergence of *Homo sapiens*, it has recently been proposed that post-Neolithic changes in bite configuration likely caused by diet changes gave rise to a new class of speech sounds, the labio-dentals, produced by positioning the lower lip against the upper teeth (Blasi et al. 2019). In general, vocalization is undergoing intense study at present, especially the questions of communicative exchange (Pika et al. 2018) and vocal learning (Vernes et al. 2021), as it becomes apparent that human vocalization seems to have more in common with other bird and animal sound emissions than earlier researchers believed. However, as Arbib (2021) notes, conveying “aboutness” in general is an apparently uniquely human capability.