1. Record Nr.          UNINA9910770258703321

   Titolo               Chinese Language Resources : Data Collection, Linguistic Analysis,
                        Annotation and Language Processing / / Chu-Ren Huang, Shu-Kai
                        Hsieh, and Peng Jin, editors

   Pubbl/distr/stampa   Cham, Switzerland : , : Springer, , [2023]
                        ©2023

   ISBN                 3-031-38913-1

   Edizione             [First edition.]

   Descrizione fisica   1 online resource (0 pages)

   Collana              Text, Speech and Language Technology Series ; ; Volume 49

   Disciplina           495.1072

   Soggetti             Chinese language - Data processing
                        Computational linguistics

   Lingua di pubblicazione   Inglese

   Formato              Materiale a stampa

   Livello bibliografico   Monografia

   Nota di bibliografia   Includes bibliographical references.

   Nota di contenuto    Intro -- Biography of Prof. Shiwen Yu -- A Chronological Biography of
                        Professor Shiwen Yu -- Acknowledgments -- Contents -- Editors and
                        Contributors -- Part I: Overview -- Chapter 1: Chinese Language
                        Resources Through One-Third of a Century -- 1.1 Headwater -- 1.2
                        Vision: Of Peaks and Giants -- 1.3 From the Great Mountains Long
                        Streams Flow -- 1.4 The Versatility of Language Resources  -- 1.5
                        Giving Shape to Water -- 1.6 Deriving Sharable and Versatile
                        Knowledge   -- 1.7 The Power of Language Data as Water  -- 1.8
                        Conclusion and Dedication   -- References -- Chapter 2: Chinese
                        Comprehensive Language Knowledge Base -- 2.1 Why Was the Chinese
                        Language Knowledge Base Constructed? -- 2.2 Cornerstone of the
                        CLKB: Grammatical Knowledge Base of Contemporary Chinese -- 2.3
                        Profile of the CLKB -- 2.3.1 PSKB -- 2.3.2 BPTC -- 2.4 What Was
                        Learned from the Development of the CLKB? -- 2.4.1 Fundamental
                        Research and Application Research -- 2.4.2 Theoretical Research and
                        Engineering Practices -- 2.4.3 Development Goals and Process
                        Monitoring -- 2.4.4 Balance of Scale and Quality -- 2.5 Conclusion --
                        References -- Chapter 3: Introduction to CKIP´s Language Resources
                        and Their Applications -- 3.1 Background -- 3.2 Language Resources
                        -- 3.2.1 Chinese Writing System Resources -- Database of Component
                        Parts of Chinese Characters -- Hantology -- 3.2.2 Lexical Databases