

The Unicode Cookbook for Linguists

Managing writing systems using
orthography profiles

Steven Moran & Michael Cysouw

Change dedication in localmetadata.tex

Preface

In the 1840s, several researchers and teachers began developing notation to accurately transcribe articulatory gestures and their acoustic correlates. Their goal was not only transcription, they wanted to use their notation for spelling reform (especially English), for teaching literacy and foreign languages, and in particular for developing writing systems for unwritten languages.

At first, the challenge of devising a universal phonetic notation seemed straightforward: normalize the relationship between graphemes and sounds (<c> in 'cake' -> /k/; <c> in 'cell' -> /s/) to remove inconsistencies introduced by the historical divergence of written and spoken forms (e.g. in English and French orthographies), and keep adding new symbol-sound mappings until all sounds are accounted for.

This approach was successful for the first few languages investigated. But as the early phoneticians were confronted with more and more articulations, they encountered a wealth of data that they had not expected. The in-depth exploration of articulatory phonetics led to the discovery that the number of possible spoken sounds is infinite. **Hockett1995** notes:

But as they observed articulation more and more closely, and as they took more and more languages into consideration, the number of sounds they could discriminate and the number of symbols needed to denote them both increased seemingly without limit. In fact, as early as 1848, Ellis asserted that in principle the number of different sounds is endless. With the development of instrumental phonetics, largely by L'Abbé Rousselot, it soon turned out that, given sufficiently accurate measuring devices, that is in fact the case.

The first instrumental measurements in the 1870-80s showed unequivocally that spoken language is comprised of sounds waves; it is dynamic and not discrete. Given the fluid nature of the speech stream, sounds differ from one utterance to the next and the same sound may differ in different phonetic environments. However, many of these alterations are predictable in the context of their environments and they need not necessarily form distinctive lexical contrasts.

Preface

These findings left the early phoneticians in a conundrum given the aims of developing precise phonetic notion, but also practically-minded alphabetic-based writing. It is in this vein that the phonetic and phonemic distinction was introduced very early into transcription. The etic vs emic principle took hold, in essence, even before the dichotomy was widely accepted and before terminology for linguistic distinctiveness was introduced and adopted.

A transcription system that showed both contrastive and non-contrastive elements emerged, taking form notably in the International Phonetic Alphabet (which itself has been updated from time-to-time over the last nearly 120 years). The association between phonetic and phonemic would be maintained in the International Phonetic Alphabet and through the various instantiations of phonological theory, which developed on the work by early phoneticians.

The model for phonemic contrast was alphabetic writing, in particular shallow orthography, i.e. one-to-one mappings between sound and symbol. Traditional alphabetic orthography, however, did not scale to precise phonetic notation, so the former was sprinkled with diacritics to denote non-contrastive phenomena. This dichotomy between phonemic and phonetic and its implementation in writing and consistent encoding causes many issues – linguistic, theoretic, and most recently technological – which we bring to light in this book.

Contents

Preface	v
1 Writing Systems	1
1.1 Introduction	1
1.2 Encoding	4
1.3 Linguistic terminology	7
1.4 The Unicode approach	11
2 Unicode pitfalls	17
2.1 Wrong it ain't	17
2.2 Pitfall: Characters are not glyphs	18
2.3 Pitfall: Characters are not graphemes	19
2.4 Pitfall: Missing glyphs	19
2.5 Pitfall: Faulty rendering	21
2.6 Pitfall: Blocks	22
2.7 Pitfall: Names	23
2.8 Pitfall: Homoglyphs	24
2.9 Pitfall: Canonical equivalence	25
2.10 Pitfall: Absence of canonical equivalence	27
2.11 Pitfall: File formats	29
2.12 Recommendations	31
3 The International Phonetic Alphabet	33
3.1 Brief history	33
3.2 Premises	34
3.3 Principles	36
3.4 Computational representation	38
3.5 SAMPA and X-SAMPA	41
3.6 The need for a single multilingual environment	42
3.7 Unicode and ISO 10646	45

Contents

4	IPA meets Unicode	49
4.1	The twain shall meet	49
4.2	Pitfall: No complete IPA code block	50
4.3	Pitfall: IPA homoglyphs in Unicode	51
4.4	Pitfall: Homoglyphs in IPA	54
4.5	Pitfall: Multiple encoding options in IPA	54
4.6	Pitfall: Tie bar	56
4.7	Pitfall: Ligatures and digraphs	58
4.8	Pitfall: Missing decomposition	58
4.9	Pitfall: Different notions of diacritics	60
4.10	Pitfall: No unique diacritic ordering	61
4.11	Pitfall: Revisions to the IPA	64
4.12	Additions to the IPA	66
4.13	Recommendations	67
5	Orthography profiles	77
5.1	Characterizing writing systems	77
5.2	Informal description	79
5.3	Formal specification	81
6	Implementation	87
6.1	Available implementations	87
6.2	Profiles and error reporting	90
6.3	Different ways to write a profile	91
6.4	Using an orthography profile skeleton	94
6.5	Rule ordering	96
6.6	Contextually specified graphemes	99
6.7	Profile skeleton with columns for editing	100
6.8	Formatting grapheme separation	102
6.9	Remaining issues	103
7	Use cases	107
7.1	Introduction	107
7.2	Tokenization and error checking	107
	List of references	109
	Index	112
	Name index	112

Language index	112
Subject index	112

1 Writing Systems

1.1 Introduction

Writing systems arise and develop in a complex mixture of cultural, technological and practical pressures. They tend to be highly conservative, in that people who have learned to read and write in a specific way—however impractical or tedious—are mostly unwilling to change their habits. Writers tend to resist spelling reforms. In all literate societies there exists a strong socio-political mainstream that tries to force unification of writing (for example by strongly enforcing “right” from “wrong” writing in schools). However, there is also a large community of users who take as many liberties in their writing as they can get away with.

For example, the writing of tone diacritics in Yoruba is often proclaimed to be the right way to write, although many users of Yoruba orthography seem to be perfectly fine with leaving them out. As pointed out by the proponents of the official rules, there are some homographs when leaving out the tone diacritics (Olúmúyìw 2013: 44). However, writing systems (and the languages they represent) are normally full of homographs (and homophones), which is normally not a problem at all for speakers of the language. More importantly, writing is not just a purely functional tool, but just as importantly it is a mechanism to signal social affiliation. By showing that you *know the rules* of expressing yourself in writing, others will more easily accept you as a worthy participant in their group. And that just as well holds for obeying to the official rules when writing a job application, as for obeying to the informal rules when writing an SMS to classmates in school. The case of Yoruba writing is an exemplary case, as even after more than a century of efforts to standardize the writing systems, there is still a wide range of variation of writing in daily use (Olúmúyìw 2013).

Formalizing orthographic structure

The resulting cumbersome and often illogical structure of writing systems, and the enormous variability of existing writing systems for the world’s languages is a fact of life scholars have to accept and should try to adapt to as good as possible.

Our plea in this book is a proposal how to do exactly that: formalize knowledge about individual writing systems in a form that is both easy to use for linguists in daily practice, and at the same time computer-readable for automatic processing.

When considering the worldwide linguistic diversity, including the many lesser-studied and endangered languages, there exist numerous different orthographies using symbols from the same scripts. For example, there are hundreds of orthographies using Latin-based alphabetic scripts. All of these orthographies use the same symbols, but these symbols differ in meaning and usage throughout the various orthographies. To be able to computationally use and compare different orthographies, we need a way to specify all orthographic idiosyncrasies in a computer-readable format (a process called *TAILORING* in Unicode parlance). We call such specifications *ORTHOGRAPHY PROFILES*. Ideally, these specifications have to be integrated into so-called Unicode locale descriptions, though we will argue that in practice this is often not the most useful solution for the kind of problems arising in the daily practice of linguistics. Consequently, a central goal of this book is to flesh out the linguistic challenges for locale descriptions, and work out suggestions to improve their structure for usage in a linguistic context. Conversely, we also aim to improve linguists' understanding and appreciation for the accomplishments of the Unicode Consortium in the development of the Unicode Standard.

The necessity to computationally use and compare different orthographies most forcefully arises in the context of language comparison. Concretely, the proper processing of orthographies and transcriptions becomes crucial for the development of quantitative methods for language comparison and historical reconstruction in order to investigate worldwide linguistic variation and to model the historical and areal processes that underlie linguistic diversity. In this work, it is crucial to be able to flexibly process across numerous resources with different orthographies. In many cases even different resources on the same language use different orthographic conventions. Another orthographic challenge that we encounter regularly in our linguistic practice is electronic resources on a particular language that claim to follow a specific orthographic convention (often a resource-specific convention), but on closer inspection such resources are almost always not consistently encoded. Thus a second goal of our orthography profiles is to allow for an easy specification of orthographic conventions, and use such profiles to check consistency and to report errors to be corrected.

A central step in our proposed solution to this problem is the tailored grapheme separation of strings of symbols, a process we call *GRAPHEME TOKENIZATION*. Basically, given some strings of symbols (e.g. morphemes, words, sentences) in

a specific source, our first processing step is to specify how these strings have to be separated into graphemes, considering the specific orthographic conventions used in a particular source document. Our experience is that such a graphemic tokenization can often be performed reasonably accurately without extensive in-depth knowledge about the phonetic and phonological details of the language in question. For example, the specification that <ou> is a grapheme of English is a much easier task than to specify what exactly the phonetic values of this grapheme are in any specific occurrence in English words. Grapheme separation is a task that can be performed relatively reliably and with limited availability of time and resources (compare, for example, the daunting task of creating a complete phonetic or phonological normalization).

Although grapheme tokenization is only one part of the solution, it is an important and highly fruitful processing step. Given a grapheme tokenization, various subsequent tasks become easier, like (a) temporarily reducing the orthography in a processing pipeline, e.g. only distinguishing high versus low vowels; (b) normalizing orthographies across sources (often including temporary reduction of oppositions), e.g. specifying an (approximate) mapping to the International Phonetic Alphabet; (c) using co-occurrence statistics across different languages (or different sources in the same language) to estimate the probability of grapheme matches, e.g. with the goal to find regular sound changes between related languages or transliterations between different sources in the same language.

Structure of this book

Before we deal with these proposals we will first discuss the theoretical background on text encoding, on the Unicode standard, and on the International Phonetic Alphabet. In the remainder of this chapter, we give an extended introduction to the notion of encoding (Section 1.2) and the principles of writing systems, both from a linguistic perspective (Section 1.3) and from the perspective of the Unicode Consortium (Section 1.4). We consider the Unicode Standard to be a breakthrough (and ongoing) development that fundamentally changed the way we look at writing systems, and we aim to provide here a slightly more in-depth survey of the many techniques that are available in the standard. A good appreciation for the solutions that the Unicode Standard provides, also allows for a thorough understanding of possible pitfalls that one might encounter when using it (Chapter 2). As an example of the current state-of-the-art, we discuss the International Phonetic Alphabet (IPA, Chapter 3) and its rather problematic marriage with the Unicode Standard (Chapter 4).

The second part of the book describes our proposals for how to deal with the Unicode Standard in the daily practice of linguistics. First, we discuss the challenges of characterizing a writing system. To solve these problems, we propose the notions of orthography profiles, closely related to Unicode locale descriptions (Chapter 5). We provide reference implementation of our proposals in R and in Python, available as open-source libraries. An introduction to these implementations of our proposals is described in Chapter 6. Finally, we discuss practical issues with actual examples in Chapter 7.

The following conventions are followed in this paper. All phonemic and phonetic representations are given in the International Phonetic Alphabet (IPA), unless noted otherwise (International Phonetic Association 2005). Standard conventions are used for distinguishing between graphemic < >, phonemic / / and phonetic [] representations. For character descriptions, we follow the notational conventions of the Unicode Standard (Unicode Consortium 2014). Character names are represented in small capital letters (e.g. LATIN SMALL LETTER SCHWA) and code points are expressed as U+n where *n* is a four to six digit hexadecimal number, e.g. U+0256, which can be rendered as the glyph <ə>.

1.2 Encoding

There are many in-depth histories of the origin and development of writing systems (e.g. Robinson 1995; Powell 2012), a story that we therefore will not repeat here. However, the history of turning writing into computer readable code is not so often told, so we decided to offer a short survey of the major developments of such encoding here.¹ This history turns out to be intimately related to the history of telegraphic communication.

Telegraphy

Writing systems have existed for roughly 6000 years, allowing people to exchange messages through time and space. Additionally, to quickly bridge large

¹ Because of the recent history as summarized in this section, we have used mostly rather ephemeral internet sources. When not referenced by traditional literature in the bibliography, we have used <http://www.unicode.org/history/> and various Wikipedia pages for the information presented here. A useful survey of the historical development of the physical hardware of telegraphy and telecommunication is Huurdeman (2003). Most books that discuss the development of encoding of telegraphic communication focus of cryptography, e.g. Singh (1999), and forego the rather interesting story of “open”, i.e. non-cryptographic, encoding that is related here.

geographic distances, telegraphic systems of communication (from Greek *τῆλε γράφειν* ‘distant writing’) have a long and widespread history since ancient times. The most common telegraphic systems worldwide are so-called whistled languages (Meyer 2015), but also drumming languages (Meyer, Dentel & Seifart 2012) and signalling by smoke, fire, flags or even change in water levels through hydraulic pressure have been used as forms of telegraphy.

Telegraphy was reinvigorated in the end of the eighteenth century through the introduction of so-called semaphoric systems by Claude Chapelle to convey messages over large distances. Originally, various specially designed contraptions were used to send messages. Today, descendants of these systems are still in limited use, for example utilizing flags or flashing lights. The “innovation” of those semaphoric systems was that all characters of the written language were replaced one-to-one by visual signals. Since then, all telegraphic systems have taken this principle, namely that any language to be transmitted first has to be turned into some orthographic system, which subsequently is encoded for transmission by the sender, and then turned back into orthographic representation at the receiver side.² This of course implies that the usefulness of any such telegraphic encoding completely depends on the sometimes rather haphazard structure of orthographic systems.

In the nineteenth century, electric telegraphy led to a next innovation in which written language characters were encoded by signals sent through a copper wire. Originally, BISIGNAL CODES were used, consisting of two different signals. For example, Carl Friedrich Gauss in 1833 used positive and negative current (Mania 2008: 282). More famous and influential, Samuel Morse in 1836 used long and short pulses. In those bisignal codes each character from the written language was encoded with a different number of signals (between one and five), so two different separators are needed: one between signals and one between characters. For example, in Morse-code there is a short pause between signals and a long pause between characters.³

² Sound and video-based telecommunication of course takes a different approach by ignoring the written version of language and directly encode sound waves or light patterns.

³ Actually, Morse-code also includes an extra long pause between words. Interestingly, it took a long time to consider the written word boundary—using white-space—as a bona-fide character that should simply be encoded with its own code point. This happened only with the revision of the Baudot-code (see below) by Donald Murray in 1901, in which he introduced a specific white-space code. This principle has been followed ever since.

Binary encoding

From those bisignal encodings, true BINARY CODES developed with a fixed length of signals per character. In such systems only a single separator between signals is needed, because the separation between characters can be established by counting until a fixed number of signals has passed.⁴ In the context of electric telegraphy, such a binary code system was first established by Émile Baudot in 1870, using a fixed combination of five signals for each written character.⁵ There are $2^5 = 32$ possible combination when using five binary signals; an encoding today designated as 5-bit. These codes are sufficient for all Latin letters, but of course they do not suffice for all written symbols, including punctuation and digits. As a solution, the Baudot code uses a so-called “shift” character, which signifies that from that point onwards—until shifted back—a different encoding is used, allowing for yet another set of 32 codes. In effect, this means that the Baudot code, and the INTERNATIONAL TELEGRAPH ALPHABET (ITA) derived from it, had an extra bit of information, so the encoding is actually 6-bit (with $2^6 = 64$ different possible characters). For decades, this encoding was the standard for all telegraphy and it is still in limited use today.

To also allow for different uppercase and lowercase letters and for a large variety of control characters to be used in the newly developing technology of computers, the American Standards Association decided to propose a new 7-bit encoding in 1963 (with $2^7 = 128$ different possible characters), known as the AMERICAN STANDARD CODE FOR INFORMATION INTERCHANGE (ASCII), geared towards the encoding of English orthography. With the ascent of other orthographies in computer usage, the wish to encode further variation of Latin letters (like German <ß> or various letters with diacritics, like <è>) led the Digital Equipment Corporation to introduce an 8-bit MULTINATIONAL CHARACTER SET (MCS, with $2^8 = 256$ different possible characters), first used with the introduction of the VT220 Terminal in 1983.

Because 256 characters were clearly not enough for the many different characters needed in the world’s writing systems, the ISO/IEC 8859 standard in 1987 extended the MCS to include 16 different 8-bit code pages. For example, part 5 was

⁴ Of course, no explicit separator is needed at all when the timing of the signals is known, which is the principle used in all modern telecommunication systems. An important modern consideration is also how to know where to start counting when you did not catch the start of a message, something that is known in Unicode as SELF SYNCHRONIZATION.

⁵ True binary codes have a longer history, going at least back to the Baconian cipher devised by Francis Bacon in 1605. However, the proposal by Baudot was the quintessential proposal leading to all modern systems.

used for Cyrillic characters, part 6 for Arabic, and part 7 for Greek.⁶ This system almost immediately was understood to be insufficient and impractical, so various initiatives to extend and reorganize the encoding started in the 1980s. This led, for example, to various proprietary encodings from Microsoft (e.g. Windows Latin 1) and Apple (e.g. Mac OS Roman), which one still sometimes encounters today.

More wide-ranging, various people in the 1980s started to develop true international code sets. In the United States, a group of computer scientists formed the UNICODE CONSORTIUM, proposing a 16-bit encoding in 1991 (with $2^{16} = 65,536$ different possible characters). At the same time in Europe, the INTERNATIONAL ORGANIZATION FOR STANDARDIZATION (ISO) was working on ISO 10646 to supplant the ISO/IEC 8859 standard. Their first draft of the UNIVERSAL CHARACTER SET (UCS) in 1990 was 31-bit (with theoretically $2^{31} = 2,147,483,648$ possible characters, but because of some technical restrictions only 679,477,248 were allowed). Since 1991, the Unicode Consortium and the ISO jointly develop the UNICODE STANDARD, or ISO/IEC 10646, leading to the current system including the original 16-bit Unicode proposal as the BASIC MULTILINGUAL PLANE, and 16 additional planes of 16-bit for further extensions (with in total $(1 + 16) \cdot 2^{16} = 1,114,112$ possible characters). The most recent version of the Unicode Standard (currently at version number 7.0) was published in June 2014 and it defines 112,956 different characters (Unicode Consortium 2014).

In the next section we provide a very brief overview of the linguistic terminology concerning writing systems before turning to the slightly different computational terminology in the subsequent section on the Unicode Standard.

1.3 Linguistic terminology

Linguistically speaking, a WRITING SYSTEM is a symbolic system that uses visible or tactile signs to represent language in a systematic way. The term writing system has two mutually exclusive meanings. First, it may refer to the way a particular language is written. In this sense the term refers to the writing system

⁶ In effect, because $16 = 2^4$, this means that ISO/IEC 8859 was actually an $8 + 4 = 12$ -bit encoding, though with very many duplicates by design, namely all ASCII codes were repeated in each 8-bit code page. To be precise, ISO/IEC 8859 used the 7-bit ASCII as the basis for each code page, and defined 16 different 7-bit extensions, leading to $(1 + 16) \cdot 2^7 = 2,176$ possible characters. However, because of overlap and not-assigned codes points the actual number of symbols was much smaller.

of a particular language, as, for example, in ‘*the Serbian writing system uses two scripts: Latin and Cyrillic.*’ Second, the term writing system may also refer to a type of symbolic system as used among the world’s languages to represent the language, as, for example, in ‘*alphabetic writing system.*’ In this latter sense the term refers to how scripts have been classified according to the way that they encode language, as in, for example, ‘*the Latin and Cyrillic scripts are both alphabetic writing systems.*’ To avoid confusion, this second notion of writing system would more aptly have been called SCRIPT SYSTEM.

Writing systems

Focussing on the first sense, we distinguish two different kinds of writing systems used for a particular language, namely transcriptions and orthographies.

First, TRANSCRIPTION is a scientific procedure (and also the result of that procedure) for graphically representing the sounds of human speech at the phonetic level. It incorporates a set of unambiguous symbols to represent speech sounds, including conventions that specify how these symbols should be combined. A transcription system is a specific system of symbols and rules used for transcription of the sounds of a spoken language variety. In principle, a transcription system should be language-independent, in that it should be applicable to all spoken human languages. The INTERNATIONAL PHONETIC ALPHABET (IPA) is a commonly used transcription system that provides a medium for transcribing languages at the phonetic level. However, there is a long history of alternative kinds of transcription systems (see Kemp 2006) and today various alternatives are in widespread use (e.g. X-SAMPA and Cyrillic-based phonetic transcription systems). Many users of IPA do not follow the standard to the letter, and many dialects based on the IPA have emerged, e.g. the Africanist and Americanist transcription systems. Note that IPA symbols are also often used to represent language on a phonemic level. It is important to realize that in this usage the IPA symbols are not a transcription system, but rather an orthography (though with strong links to the pronunciation). Further, a transcription system does not need to be as highly detailed as the IPA. It can also be a system of broad sound classes. Although such an approximative transcription is not normally used in linguistics, it is widespread in technological approaches (Soundex and variants, e.g. Knuth 1973: 391–392; Postel 1969; Beider & Morse 2008), and it is sometimes fruitfully used in automatic approaches to historical linguistics (Dolgopolsky 1986; List 2012; Brown, Holman & Wichmann 2013).

Second, an ORTHOGRAPHY specifies the symbols, punctuations, and the rules in which a specific language is written in a standardized way. Orthographies are of-

ten based on a phonemic analysis, but they almost always include idiosyncrasies because of historical developments (like sound changes or loans) and because of the widely-followed principle of lexical integrity (i.e. the attempt to write the same lexical root in a consistent way, also when synchronic phonemic rules change the pronunciation, as for example with final devoicing in many Germanic languages). All orthographies are language-specific (and often even resource-specific), although individual symbols or rules might be shared between languages. A PRACTICAL ORTHOGRAPHY is a strongly phoneme-based writing system designed for practical use by speakers. The mapping relation between phonemes and graphemes in practical orthographies is purposely shallow, i.e. there is mostly a systematic and faithful mapping from a phoneme to a grapheme. Practical orthographies are intended to jumpstart written materials development by correlating a writing system with the sound units of a language (cf. Meinhof & Jones 1928). Symbols from the IPA are often used by linguists in the development of such practical orthographies for languages without writing systems, though this usage of IPA symbols should not be confused with transcription (as defined above).

Further, a TRANSLITERATION is a mapping between two different orthographies. It is the process of “recording the graphic symbols of one writing system in terms of the corresponding graphic symbols of a second writing system” (Kemp 2006: 396). In straightforward cases, such a transliteration is simply a matter of replacing one symbol with another. However, there are widespread complications, like one-to-many or many-to-many mappings, which are not always easy, or even possible, to solve without listing all cases individually (cf. Moran 2012: Ch. 2).

Script systems

Different kinds of writing systems are classified into script systems. A SCRIPT is a collection of distinct symbols as employed by one or more orthographies. For example, both Serbian and Russian are written with subsets of the Cyrillic script. A single language, like Serbian or Japanese, can also be written using orthographies based on different scripts. Over the years linguists have classified script systems in a variety of ways, with the tripartite classification of logographic, syllabic, and alphabetic remaining the most popular, even though there are at least half a dozen different types of script systems that can be distinguished (Daniels 1990, 1996).

Breaking it down further, a script consists of GRAPHEMES, and graphemes consist of CHARACTERS. In the linguistic terminology of writing systems, a CHAR-

ACTER is a general term for any self-contained element in a writing system.⁷ Although in literate societies most people have a strong intuition about what the characters are in their particular orthography or orthographies, it turns out that the separation of an orthography into separate characters is far from trivial. The widespread intuitive notion of a character is strongly biased towards educational traditions, like the alphabet taught at schools, and technological possibilities, like the available type pieces in a printer's job case, the keys on a typewriter, or the symbols displayed in Microsoft Word's symbol browser. In practice, characters often consist of multiple building blocks, each of which could be considered a character in its own right. For example, although a Chinese character may be considered to be a single basic unanalyzable unit, at a more fine-grained level of analysis the internal structure of Chinese characters is often comprised of smaller semantic and phonetic units that should be considered characters (Sproat 2000). In alphabetic scripts, this problem is most forcefully exemplified by diacritics.

A DIACRITIC is a mark, or series of marks, that may be above, below, or through other characters (Gaultney 2002). Diacritics are sometimes used to distinguish homophonous words, but they are more often used to indicate a modified pronunciation (Daniels & Bright 1996: xli). The central question is whether, for example, <e>, <è>, <a> and <à> should be considered four characters, or different combinations of three characters. In general, multiple characters together can form another character, and it is not always possible to decide on principled grounds what should be the basic building blocks of an orthography.

For that reason, it is better to analyze an orthography as a collection of graphemes. A GRAPHEME is the basic, minimally distinctive symbol of a particular writing system, alike to the phoneme is an abstract representation of a distinct sound in a specific language. The term GRAPHEME was modeled after the term PHONEME and represents a contrastive graphical unit in a writing system (see Kohrt 1986 for a historical overview of the term grapheme). Most importantly, a single grapheme regularly consists of multiple characters, like <th>, <ou> and <gh> in English (note that each character in these graphemes is also a separate grapheme in English). Such complex graphemes are often used to represent single phonemes. So, a combination of characters is used to represent a single phoneme. Note that the opposite is also found in writing systems, in cases in which a single character represents a combination of two or more phonemes. For example, <x> in English orthography represents a combination of the phonemes /k/ and /s/.

⁷ There is a second interpretation of the term CHARACTER, i.e. a conventional term for a unit in the Chinese writing system (Daniels 1996). This interpretation will not be further explored in this paper.

Further, conditioned or free variants of a grapheme are called ALLOGRAPHS. For example, the distinctive forms of Greek sigma are conditioned, with <σ> being used word-internally and <ς> being used at the end of a word. In sum, there are many-to-many relationships between phonemes and graphemes as they are expressed in the myriad of language- and resource-specific orthographies.

This exposition of the linguistic terminology involved in describing writing systems has been purposely brief. We have highlighted some of the linguistic notions that are pertinent, yet sometimes confused with, the technological definitions developed for the computational processing of the world's writing systems, which we describe in the next section.

1.4 The Unicode approach

The conceptualization and terminology of writing systems was rejuvenated by the development of the Unicode Standard, with major input from Mark Davis, co-founder and long-term president of the Unicode Consortium. For many years, “exotic” writing systems and phonetic transcription systems on personal computers were constrained by the American Standard Code for Information Interchange (ASCII) character encoding scheme, based on the Latin script, which only allowed for a strongly limited number of different symbols to be encoded. This implied that users could either use and adopt the (extended) Latin alphabet or they could assign new symbols to the small number of code points in the ASCII encoding scheme to be rendered by a specifically designed font (Bird & Simons 2003). In this situation, it was necessary to specify the font together with each document to ensure the rightful display of its content. To alleviate this problem of assigning different symbols to the same code points, in the late 80's and early 90's the Unicode Consortium set itself the ambitious goal of developing a single universal character encoding to provide a unique number, a code point, for every character in the world's writing systems. Nowadays, the Unicode Standard is the default encoding of the technologies that support the World Wide Web and for all modern operating systems, software and programming languages.

The Unicode Standard

The Unicode Standard represents a massive step forward because it aims to eradicate the distinction between universal (ASCII) versus language-particular (font) by adding as much as possible language-specific information into the universal

standard. However, there are still language/resource-specific specifications necessary for the proper usage of Unicode, as will be discussed below. Within the Unicode structure many of these specifications can be captured by so-called `LOCALE DESCRIPTIONS`, so we are moving to a new distinction of universal (Unicode Standard) versus language-particular (locale description). The major gain is a much larger compatibility on the universal level (because Unicode standardizes a much larger portion of writing system diversity), and much better possibilities for automatic processing on the language-particular level (because locale descriptions are computer readable specifications).

Each version of the Unicode Standard (Unicode Consortium 2014, as of writing at version 7) consists of a set of specifications and guidelines that include (i) a core specification, (ii) code charts, (iii) standard annexes and (iv) a character database.⁸ The `CORE SPECIFICATION` is a book directed toward human readers that describes the formal standard for encoding multilingual text. The `CODE CHARTS` provide a humanly readable online reference to the character contents of the Unicode Standard in the form of PDF files. The `UNICODE STANDARD ANNEXES (UAX)` are a set of technical standards that describe the implementation of the Unicode Standard for software development, Web standards, and programming languages. The `UNICODE CHARACTER DATABASE (UCD)` is a set of computer-readable text files that describe the character properties, including a set of rich character and writing system semantics, for each character in the Unicode Standard. In this section, we introduce the basic Unicode concepts, but we will leave out many details. Please consult the above mentioned full documentation for a more detailed discussion. Further note that the Unicode Standard is exactly that, namely a standard. It normatively describes notions and rules to be followed. In the actual practice of applying this standard in a computational setting, a specific implementation is necessary. The most widely used implementation of the Unicode Standard is the `INTERNATIONAL COMPONENTS FOR UNICODE (ICU)`, which offers C/C++ and Java libraries implementing the Unicode Standard.⁹

⁸ All documents of the Unicode Standard are available at <http://www.unicode.org/versions/latest/>. For a quick survey of the use of terminology inside the Unicode Standard, their glossary is particularly useful, available at <http://www.unicode.org/glossary/>. For a general introduction to the principles of Unicode, Chapter 2 of the core specification, called `GENERAL STRUCTURE`, is particularly insightful. Different from many other documents of the Unicode Standard, this general introduction is relatively easy to read and illustrated with many interesting examples from various orthography traditions all over the world.

⁹ More information about the ICU is available here: <http://icu-project.org>.

Character encoding system

The Unicode Standard is a CHARACTER ENCODING SYSTEM which goal it is to support the interchange and processing of written characters and text in a computational setting. Underlyingly, the character encoding is represented by a range of numerical values called a CODE SPACE, which is used to encode a set of characters. A CODE POINT is a unique non-negative integer within a code space (i.e. within a certain numerical range). In the Unicode Standard character encoding system, an ABSTRACT CHARACTER, for example the LATIN SMALL LETTER P, is mapped to a particular code point, in this case the decimal value 112, normally represented in hexadecimal, which then looks in Unicode parlance as U+0070.¹⁰ That encoded abstract character is rendered on a computer screen (or printed page) as a GLYPH, e.g. <p>, depending on the FONT and the context in which that character appears.

In Unicode Standard terminology, an (abstract) CHARACTER is the basic encoding unit. The term CHARACTER can be quite confusing due to its alternative definitions across different scientific disciplines and because in general the word CHARACTER means many different things to different people. It is therefore often preferable to refer to Unicode characters simply as CODE POINTS, because there is a one-to-one mapping between Unicode characters and their numeric representation. In the Unicode approach, a character refers to the abstract meaning and/or general shape, rather than a specific shape, though in code tables some form of visual representation is essential for the reader's understanding. Unicode defines characters as abstractions of orthographic symbols, and it does not define visualizations for these characters (although it does presents examples). In contrast, a GLYPH is a concrete graphical representation of a character as it appears when rendered (or rasterized) and displayed on an electronic device or on printed paper. For example, <g g g g **g** **g**> are different glyphs of the same character, i.e. they may be rendered differently depending on the typography being used, but they all share the same code point. From the perspective of Unicode they are *the same thing*. In this approach, a FONT is then simply a collection of glyphs linked to code points. Allography is not specified in Unicode (except for a few exceptional cases, due to legacy encoding issues), but can be specified in a font as a CONTEXTUAL VARIANT (a.k.a. presentation form).

¹⁰ Hexadecimal (base-16) 0070 is equivalent to decimal (base-10) 112, which can be calculated by considering that $(0 \cdot 16^3) + (0 \cdot 16^2) + (7 \cdot 16^1) + (0 \cdot 16^0) = 7 \cdot 16 = 112$. Underlyingly, computers will of course treat this code point binary (base-2) as 11100000, as can be seen by calculating that $(1 \cdot 2^7) + (1 \cdot 2^6) + (1 \cdot 2^5) + (0 \cdot 2^4) + (0 \cdot 2^3) + (0 \cdot 2^2) + (0 \cdot 2^1) + (0 \cdot 2^0) = 64 + 32 + 16 = 112$.

Each code point in the Unicode Standard is associated with a set of CHARACTER PROPERTIES as defined by the Unicode character property model.¹¹ Basically, those properties are just a long list of values for each character. For example, code point U+0047 has the following properties (among many others):

- Name: LATIN CAPITAL LETTER G
- Alphabetic: YES
- Uppercase: YES
- Script: LATIN
- Extender: NO
- Simple_Lowercase_Mapping: 0067

These properties contain the basic information of the Unicode Standard and they are necessary to define the correct behavior and conformance required for interoperability in and across different software implementations (as defined in the Unicode Standard Annexes). The character properties assigned to each code point is based on each character's behavior in the real-world writing traditions. For example, the corresponding lowercase character to U+0047 is U+0067 (though note that the relation between uppercase and lowercase is in many situations much more complex than this, and Unicode has further specifications for those cases). Another use of properties is to define the script of a character.¹² In practice, script is simply defined for each character as the explicit SCRIPT property in the Unicode Character Database.

One frequently referenced property is the BLOCK property, which is often used in software applications to impose some structure to the large number of Unicode characters. Each character in Unicode belongs to a specific block. These blocks are basically an organizational structure to alleviate the administrative burden of keeping Unicode up-to-date. Blocks consist of characters that in some way belong together, so that characters are easier to find. Some blocks are connected with a specific script, like the Hebrew block or the Gujarati block. However, blocks are predefined ranges of code points, and often there will come a point

¹¹ The character property model is described in <http://www.unicode.org/reports/tr23/>, but the actual properties are described in <http://www.unicode.org/reports/tr44/>. A simplified overview of the properties is available at <http://userguide.icu-project.org/strings/properties>. The actual code tables listing all properties for all Unicode code points are available at <http://www.unicode.org/Public/UCD/latest/ucd/>.

¹² The Unicode Glossary defines the term SCRIPT as a “collection of letters and other written signs used to represent textual information in one or more writing systems. For example, Russian is written with a subset of the Cyrillic script; Ukrainian is written with a different subset. The Japanese writing system uses several scripts.”

after which the range is completely filled. Any extra characters will have to be assigned somewhere else. There is, for example, a block ARABIC, which contains most Arabic symbols. However, there is also a block ARABIC SUPPLEMENT, ARABIC PRESENTATION FORMS-A and ARABIC PRESENTATION FORM-B. The situation with Latin symbols is even more extreme. In general, the names for block should not be taken as a definitional statement. For example, many IPA symbols are not located in the aptly-names block IPA EXTENSIONS, but in other blocks (see Section 4.2).

Grapheme clusters

There are many cases in which a sequence of characters (i.e. a sequence of more than one code point) represents what a user perceives as an individual unit in a particular orthographic writing system. For this reason the Unicode Standard differentiates between ABSTRACT CHARACTER and USER-PERCEIVED CHARACTER. Sequences of multiple code points that correspond to a single user-perceived characters are called grapheme clusters in Unicode parlance. Grapheme clusters come in two flavors: (default) grapheme clusters and tailored grapheme clusters.

The (default) GRAPHEME CLUSTERS are locale-independent graphemes, i.e. they always apply when a particular combination of characters occurs independent of the writing system in which they are used. These character combinations are defined in the Unicode Standard as functioning as one TEXT ELEMENT.¹³ The simplest example of a grapheme cluster is a base character followed by a letter modifier character. For example, the sequence <n> + <~> (i.e. LATIN SMALL LETTER N at U+006E, followed by COMBINING TILDE at U+0303) combines visually into <ñ>, a user-perceived character in writing systems like that of Spanish. So, what the user perceives as a single character actually involves a multi-code-point sequence. Note that this specific sequence can also be represented with a single so-called PRECOMPOSED code point, LATIN SMALL LETTER N WITH TILDE at U+00F1, but this is not the case for all multi-code point character sequences. The problem that there multiple encodings possible for the same text element has been acknowledged early on in the Unicode Standard (e.g. for <ñ>, the sequence U+006E U+0303 should in all situations be treated identically to the precomposed U+00F1), and a system of CANONICAL EQUIVALENCE is available for such situations. Basically, the Unicode Standard offers different kind of normalizations to either decompose all precomposed characters (called NFD, NORMALIZATION

¹³ The Unicode Glossary defines text element as: “A minimum unit of text in relation to a particular text process, in the context of a given writing system. In general, the mapping between text elements and code points is many-to-many.”

FORM CANONICAL DECOMPOSITION), or compose as much as possible combinations (called NFC, NORMALIZATION FORM CANONICAL COMPOSITION). In current practice of software development, NFC seems to be preferred in most situations and is widely proposed as the preferred canonical form.

More difficult for text processing, because less standardized, is what the Unicode Standard terms TAILORED GRAPHEME CLUSTERS. Tailored grapheme clusters are locale-dependent graphemes, i.e. such combination of characters do not function as text elements in all situations. For example, the sequence <c> + <h> for the Slovak digraph <ch> or the sequence <ky> in the Sisaala practical orthography (pronounced as IPA /tʃ/, Moran 2006). These grapheme clusters are TAILORED in the sense that they must be specified on a language-by-language or writing-system-by-writing-system basis. The Unicode Standard provides technological specifications for creating locale specific data in so-called UNICODE LOCALE DESCRIPTIONS, i.e. a set of specification that defines a set of language-specific elements (e.g. tailored grapheme clusters, collation order, capitalization-equivalence), as well as other special information, like how to format numbers, dates, or currencies. Locale descriptions are saved in the COMMON LOCALE DATA REPOSITORY (CLDR),¹⁴ a repository of language-specific definitions of writing system properties, each of which describes specific usages of characters. Each locale can be encoded in a document using the LOCALE DATA MARKUP LANGUAGE (LDML). LDML is an XML format and vocabulary for the exchange of structured locale data. Unicode Locale Descriptions allow users to define language- or even resource-specific writing systems or orthographies.¹⁵ However, there are various drawbacks of locale descriptions for the daily practice of linguistic work in a multilingual setting.

¹⁴ More information about the CLDR can be found here: <http://cldr.unicode.org/>.

¹⁵ The Unicode Glossary defines WRITING SYSTEM only very loosely, as it is not a central concept in the Unicode Standard. A writing system is, “A set of rules for using one or more scripts to write a particular language. Examples include the American English writing system, the British English writing system, the French writing system, and the Japanese writing system.”

2 Unicode pitfalls

2.1 Wrong it ain't

In this chapter we describe some of the most common pitfalls that we have encountered when using the Unicode Standard in our own work, or in discussion with other linguists. This section is not meant as a criticism of the decisions made by the Unicode Consortium; on the contrary, we aim to highlight where the technological aspects of the Unicode Standard diverge from many users' intuitions. What have sometimes been referred to as problems or inconsistencies in the Unicode Standard are mostly due to legacy compatibility issues, which can lead to unexpected behavior by linguists using the standard. However, there are also some cases in which the Unicode Standard has made decisions that theoretically could have been taken differently, but for some reason or another (mostly very good reasons) were accepted as they are now. We call such behavior that executes without error but does something different than the user expected—often unknowingly—a `PITFALL`.

In this context, it is important to realize that the Unicode Standard was not developed to solve linguistic problems per se, but to offer a consistent computational environment for written language. In those cases in which the Unicode Standard behaves differently as expected, we think it is important not to dismiss Unicode as “wrong” or “deficient”, because our experience is that in almost all cases the behavior of the Unicode Standard has been particularly well thought through. The Unicode Consortium has a more wide-ranging view of matters and often examines important practical use-cases that from a linguistic point of view are normally not considered. Our general guideline for dealing with the Unicode Standard is to accept it as it is, and not to battle windmills. Alternatively, of course, it is possible to actively engage in the development of the standard itself, an effort that is highly appreciated by the Unicode Consortium.

2.2 Pitfall: Characters are not glyphs

A central principle of Unicode is the distinction between character and glyph. A character is the abstract notion of a symbol in a writing system, while a glyph is the concrete drawing of such a symbol. In practice, there is a complex interaction between characters and glyphs. A single Unicode character may of course be rendered as a single glyph. However, a character may also be a piece of a glyph, or vice-versa. Actually, all possible relations between glyphs and characters are attested.

First, a single character may have different contextually determined glyphs. For example, characters in writing systems like Hebrew and Arabic have different glyphs depending on where they appear in a word. Some letters in Hebrew change their form at the end of the word, and in Arabic, primary letters have four contextually-sensitive variants (isolated, word initial, medial and final). Second, a single character may be rendered as a sequence of multiple glyphs. For example, in Tamil one Unicode character may result in a combination of a consonant and vowel, which are rendered as two adjacent glyphs by fonts that supports Tamil. Third, a single glyph may be a combination of multiple characters. For example, the ligature <fi>, a single glyph, is the result of two characters, <f> and <i>, that have undergone glyph substitution by font rendering (see also Section 2.5). Like contextually-determined glyphs, ligatures are (intended) artefacts of text processing instructions. Finally, a single glyph may be a part of a character, as exemplified by diacritics like the diaeresis <¨> in <ë>.

Further, the rendering of a glyph is dependent on the font being used. For example, the Unicode character LATIN SMALL LETTER G appears as <g> and <g> in the Linux Libertine and Courier fonts, respectively, because their typefaces are designed differently. Furthermore, font face may change the visual appearance of a character, for example Times New Roman two-story <a> changes to a single-story glyph in italics <a>. This becomes a real problem for some phonetic typesetting (see Section 4.3).

In sum, character-to-glyph mappings are complex technological issues that the Unicode Consortium has had to address in the development of the Unicode Standard, but for the lay user they can be utterly confusing because visual rendering does not (necessarily) indicate logical encoding.

2.3 Pitfall: Characters are not graphemes

The Unicode Standard is a character encoding system, and not a writing system encoding system. This most forcefully becomes clear with the notion of grapheme. From a linguistic point of view, graphemes are the basic building blocks of a writing system (see Section 1.3). It is extremely common, up to the point of being universally attested, that writing systems use combinations of multiple symbols as a single grapheme, like <sch>, <th> or <ei>. There is no possibility to encode such complex graphemes using the Unicode Standard.

The Unicode Standard deals with complex graphemes only inasmuch they consist of base characters with diacritics (see Section 4.9 for a discussion of the notion of diacritic). The Unicode Standard calls such combination “grapheme clusters”. Complex graphemes consisting of multiple base characters, like <sch>, are called “tailored grapheme clusters” in Unicode parlance (see Section 1.4).

Inspecting the Unicode Standard, there appear to be special Unicode characters to “glue” together characters into larger tailored grapheme clusters, specifically the ZERO WIDTH JOINER at U+200D and the COMBINING GRAPHEME JOINER at U+034F. However, these characters are confusingly named (cf. Section 2.7). Both codepoints actually do not join characters, but explicitly separate them. The zero-width joiner (ZWJ) can be used to solve special problems related to ordering (called “collation” in Unicode parlance). The combining grapheme joiner (CGJ) can be used to separate characters that are not supposed to form ligatures.

To solve the issue of tailored grapheme clusters, Unicode offers some assistance in the form of the Unicode Locale Descriptions. However, in the practice of linguistic research, this is not a real solution. For that reason we propose to use orthography profiles (see Chapter 5). Basically, both orthography profiles and locale descriptions offer a way to specify tailored grapheme clusters. For example, for English one could specify that <sh> is such a cluster. Consequently, this sequence of characters is then always interpreted as a complex grapheme. For cases in which this is not the right decision, like in the English word *mishap*, the ZERO WIDTH JOINER at U+200D has to be entered between <s> and <h>.

2.4 Pitfall: Missing glyphs

The Unicode Standard is often praised (and deservedly so) for solving many of the perennial problems with the interchange and display of the world’s writing systems. However, a common complaint from users is that, while the praise may be true, they mostly just see some boxes on their screen instead of those promised

symbols. The problem of course is that users' computers do not have any glyphs installed matching the Unicode code points in the file they are trying to inspect. It is important to realize that internally in the computer everything still works as expected: any handling of Unicode code points works independently of how they are displayed on the screen. So, although a user might only see boxes being displayed, this user should be assured that everything is still in order.

The central problem behind the missing glyphs is that designing actual glyphs includes a lot of different considerations and it is a time-consuming process. Many traditional expectations of how specific characters should look like have to be taken into account when designing glyphs. Those expectations are often not well documented, and it is mostly up to the knowledge and experience of the font designer to try and conform to them as good as possible. Therefore, most designers produce fonts only including glyphs for certain parts of the Unicode Standard, namely for those characters they feel comfortable with. At the same time, the number of characters defined by the Unicode Standard is growing with each new version, so it is neigh impossible for any designer to produce glyphs for all characters. The result of this is that, almost necessarily, each font only includes glyphs for a subset of the characters in the Unicode Standard.

The simple solution to missing glyphs is thus to install additional fonts providing additional glyphs. For the more exotic characters there is often not much choice. There are a few particularly large fonts that might be considered. First, there is the *EVERSON MONO* font made by Michael Everson, which currently includes 9,756 different glyphs (not including Chinese) updated up to Unicode 7.0.¹ Already a bit older is the *TITUS CYBERBIT BASIC* font made by Jost Gippert and Carl-Martin Bunz, which includes 10,044 different glyphs (not including Chinese), but not including newer characters added after Unicode 4.0.²

Further, we suggest to always install at least one so-called *FALL-BACK FONT*, which provides glyphs that at least show the user some information about the underlying encoded character. Apple Macintoshes have such a font (which is invisible to the user), which is designed by Michael Everson and made available for other systems through the Unicode Consortium.³ Further, the *GNU UNIFONT* is a clever way to produce bitmaps approximating the intended glyph of each available character, updated to Unicode 7.0.⁴ Finally, the Summer Institute of Linguistics provides a *SIL UNICODE BMP FALLBACK FONT*, currently available up

¹ Everson Mono is available as shareware at <http://www.evertype.com/emono/>.

² Titus Cyberbit Basic is available at <http://titus.fkidg1.uni-frankfurt.de/unicode/tituut.asp>.

³ The Apple/Everson fallback font is available for non-Macintosh users at http://www.unicode.org/policies/lastresortfont_eula.html.

⁴ The GNU Unifont is available at <http://unifoundry.com/unifont.html>.

to Unicode version 6.1. This font does not even attempt to show a real glyph, but only shows the hexadecimal code inside a box for each character, so a user can at least see the Unicode codepoint of the character to be displayed.⁵

2.5 Pitfall: Faulty rendering

A similar complaint to missing glyphs, discussed previously, is that while there might be a glyph being displayed, it does not look right. There are two reasons for unexpected visual display, namely automatic font substitution and faulty rendering. Like missing glyphs, any such problems are independent from the Unicode Standard. The Unicode Standard only includes very general information about characters and leaves the specific visual display to others to decide on. Any faulty display is thus not to be blamed on the Unicode Consortium, but on a complex interplay of different mechanisms happening in a computer to turn Unicode codepoints into visual symbols. We will only sketch a few aspects of this complex interplay here.

Most modern software applications (like Microsoft Word) offer some approach to AUTOMATIC FONT SUBSTITUTION. This means that when a text is written in a specific font (e.g. Times New Roman) and an inserted Unicode character does not have a glyph within this font, then the software application will automatically search for another font to display the glyph. The result will be that this specific glyph will look slightly different from the others. This mechanism works differently depending on the software application, and mostly only limited user influence is expected and little feedback is given, which might be rather frustrating to font-aware users.⁶

The other problem with visual display is related to the so-called FONT RENDERING. Font rendering refers to the process of the actual positioning of Unicode characters on a page of written text. This positioning is actually a highly complex problem, and many things can go wrong in the process. Well-known rendering problems, like proportional glyph size or ligatures are reasonably well understood. In contrast, the positioning of multiple diacritics relative to a base

⁵ The SIL Unicode BMP Fallback Font is available at <http://scripts.sil.org/UnicodeBMPFallbackFont>.

⁶ For example, Apple Pages does not give any feedback that a font is being replaced, and the user does not seem to have any influence on the choice of replacement (except by manually marking all occurrences). In contrast, Microsoft Word does indicate the font replacement by showing the name in the font menu of the font replacement. However, Word simply changes the font completely, so any text written after the replacement is written in a different font as before. Both behaviors leave much to be desired.

character is still a widespread problem, even within the Latin script. Especially when more than one diacritic is supposed to be placed above (or below) each other, this often leads to unexpected effects in many modern software applications. The problems arising in Arabic and in many southeast Asian scripts (like Devanagari or Burmese) are even more complex.

To understand where any problems arise it is important to realize that there are basically three different approaches to font rendering. The most widespread is Adobe's and Microsoft's `OPENTYPE` system. This approach makes it relatively easy for font developers, as the font itself does not include all details about the precise placement of individual characters. For those details, additional script-descriptions are necessary. All of those systems can lead to unexpected behavior.⁷ Alternative systems are `APPLE ADVANCED TYPOGRAPHY (AAT)` and the open-source `GRAPHITE` system from the Summer Institute of Linguistics (SIL).⁸ In both of these systems, a larger burden is placed on the description inside the font.

There is mostly no real solution to problems arising from faulty font rendering. Switching to another software application that offers better handling is the only real alternative, but this is normally not an option for daily work. The experience with rendering on the side of the software industry is developing quickly, so we can expect the situation only to get better. In the meantime one can try to correct faulty layout by tweaking baseline and/or kerning (when such option are available).

2.6 Pitfall: Blocks

The Unicode code space is subdivided into blocks of contiguous code points. For example, the block called `CYRILLIC` runs from `U+0400` till `U+04FF`. These blocks arose as an attempt at ordering the enormous amount of characters in Unicode, but the ideas of blocks very quickly ran into problems. First, the size of a block is fixed, so when a block is full, a new block will have to be instantiated somewhere further in the code space. For example, this led to the blocks `CYRILLIC SUPPLEMENT`, `CYRILLIC EXTENDED-A` (both of which are also already full) and `CYRILLIC`

⁷ For more details about OpenType, see <http://www.adobe.com/products/type/opentype.html> and <http://www.microsoft.com/typography/otspec/>. Additional systems for complex text layout are, among others, Microsoft's DirectWrite <https://msdn.microsoft.com/library/dd368038.aspx> and the open-source project HarfBuzz <http://www.freedesktop.org/wiki/Software/HarfBuzz/>.

⁸ More information about AAT can be found at <https://developer.apple.com/fonts/>. SIL's Graphite is described in detail at http://scripts.sil.org/cms/scripts/page.php?site_id=projects&item_id=graphite_home.

EXTENDED-B. Second, when a specific character already exists, then it is not duplicated in another block, although the name of the block might indicate that a specific symbol should be available there. In general, names of blocks are just an approximate indication of the kind of characters that will be in the block.

The problem with blocks arises because finding the right character among the thousands of Unicode characters is not easy. Many software applications present blocks as a primary search mechanism, because the block names suggest where to look for a particular character. However, when a user searches for an IPA character in the block IPA EXTENSIONS, then many IPA characters will not be found there. For example, the velar nasal <ŋ> is not part of the block IPA EXTENSIONS because it was already included as LATIN SMALL LETTER ENG at U+014B in the block LATIN EXTENSIONS-A.

In general, finding a specific character in the Unicode Standard is often not trivial. The names of the blocks can help, but they are not (and never were supposed to be) a foolproof structure. It is not the goal nor aim of the Unicode Consortium to provide a user interface to the Unicode Standard. If one often encounters the problem of needing to find a suitable character, there are various other useful services for end-users available.⁹

2.7 Pitfall: Names

The names of characters in the Unicode Standard are sometimes misnomers and should not be misinterpreted as definitions. For example, the COMBINING GRAPHEME JOINER at U+034F does not join characters into larger graphemes (see Section 2.3) and the LATIN LETTER RETROFLEX CLICK U+01C3 is actually not the IPA symbol for a retroflex click, but for an alveolar click (see Section 4.3). In a sense, these names can be seen as “errors.” However, it is probably better to realize that such names are just convenience labels that are not going to be changed. Just like the block names (Section 2.6), the character names are often helpful, but they are not supposed to be definitions.

The actual intended “meaning” of a Unicode codepoint is a combination of the name, the block and the character properties (see Section 1.4). Further details

⁹ The Unicode website offers a basic interface to the code charts at <http://www.unicode.org/charts/index.html>. As a more flexible interface, we particularly like PopChar from Macility, available for both Macintosh and Windows. There are also various free websites that offer search interfaces to the Unicode code tables, like <http://unicode-search.net> or <http://unicode-search.net>. A further useful approach for searching characters using shape matching is <http://shapecatcher.com>.

about the underlying intentions with which a character should be used are only accessible by perusing the actual decisions of the Unicode Consortium. All proposals, discussions and decisions of the Unicode Consortium are publicly available. Unfortunately there is not (yet) any way to easily find everything that is ever proposed, discussed and decided in relation to a specific codepoint of interest, so many of the details are often somewhat hidden.¹⁰

2.8 Pitfall: Homoglyphs

Homoglyphs are visually indistinguishable glyphs (or highly similar glyphs) that have different code points in the Unicode Standard and thus different character semantics. As a principle, the Unicode Standard does not specify how a character appears visually on the page or the screen. So in most cases, a different appearance is caused by the specific design of a font, or by user-settings like size or boldface. Taking an example already discussed in Section 2.8, the following symbols <g g g g g> are different glyphs of the same character, i.e. they may be rendered differently depending on the typography being used, but they all share the same code point (viz. LATIN SMALL LETTER G at U+0067). In contrast, the symbols <AAAAAA[?]> are all different code points, although they look highly similar—in some cases even sharing exactly the same glyph in some fonts. All these different A-like characters include the following code points in the Unicode Standard:

- <A> LATIN CAPITAL LETTER A, at U+0041
- <A> CYRILLIC CAPITAL LETTER A, at U+0410
- <A> GREEK CAPITAL LETTER ALPHA, at U+0391
- <A> CHEROKEE LETTER GO, at U+13AA
- <A> CANADIAN SYLLABICS CARRIER GHO, at U+15C5
- <A> LATIN SMALL LETTER CAPITAL A, at U+1D00
- <A> LISU LETTER A, at U+A4EE
- <Q> CARIAN LETTER A, at U+102A0
- <[?]> MATHEMATICAL SANS-SERIF CAPITAL A, U+1D5A0
- <[?]> MATHEMATICAL MONOSPACE CAPITAL A, at U+1D670

¹⁰ All proposals and other documents that are the basis of Unicode decisions are available at <http://www.unicode.org/L2/all-docs.html>. The actual decisions that make up the Unicode Standard are documented in the minutes of the Unicode Technical Committee, available at <http://www.unicode.org/consortium/utc-minutes.html>.

The existence of such homoglyphs is partly due to legacy compatibility, but for the most part these characters are simply different characters that happen to look similar.¹¹ Yet, they are supposed to behave different from the perspective of a font designer. For example, when designing a Cyrillic font, the <A> will have different aesthetics and different traditional expectation compared to a Latin <A>.

Such homoglyphs are a widespread problem for consistent encoding. Although for most users it looks like the words <voces> and <voces> are almost identical, in actual fact they do not even share a single code point.¹² For computers these two words are completely different entities. Commonly, when users with Cyrillic or Greek keyboards have to type some Latin-based orthography, they mix similar looking Cyrillic or Greek characters into their text, because those characters are so much easier to type. Similarly, when users want to enter an unusual symbol, they normally search by visual impression in their favorite software application, and just pick something that looks reasonably alike to what they expect the glyph to look like.

It is really easy to make errors at text entry and add characters that are not supposed to be included. Our proposals for orthography profiles (see Chapter 5) are a method for checking the consistency of any text. In situations in which interoperability is important, we consider it crucial to add such checks in any workflow.

2.9 Pitfall: Canonical equivalence

For some characters, there is more than one possible encoding in the Unicode Standard. This is a possible pitfall, as this would mean that for the computer there exist multiple different entities that for a user are the same. This would, for example, lead to problems with searching, as the computer would search for specific encodings, and not find all expected characters. As a solution, the Unicode Standard includes a notion of `CANONICAL EQUIVALENCE`. Different encodings are explicitly declared as equivalent in the Unicode Standard code tables. Further, to harmonize all encodings in a specific piece of text, the Unicode Standard proposes a mechanism of `NORMALIZATION`.

Consider for example the characters and following Unicode code points:

¹¹ A particularly nice interface to look for homoglyphs is <http://shapecatcher.com>, based on the principle of recognizing shapes (Belongie, Malik & Puzicha 2002).

¹² The first words consists completely of Latin characters, namely U+0076, U+006F, U+0063, U+0065 and U+0073, while the second is a mix of Cyrillic and Greek characters, namely U+03BD, U+03BF, U+0041, U+0435 and U+0455.

2 Unicode pitfalls

- <Å> LATIN CAPITAL LETTER A WITH RING ABOVE U+00C5
- <Å> ANGSTROM SIGN U+212B
- <Å> LATIN CAPITAL LETTER A U+0041 + COMBINING RING ABOVE U+030A

The character, represented here by glyph <Å>, is encoded in the Unicode Standard in the first two examples by a single-character sequence; each is assigned a different code point. In the third example, the glyph is encoded in a multiple-character sequence that is composed of two character code points. All three sequences are CANONICALLY EQUIVALENT, i.e. they are strings that represent the same abstract character and because they are not distinguishable by the user, the Unicode Standard requires them to be treated the same in regards to their behavior and appearance. Nevertheless, they are encoded differently. For example, if one were to search an electronic text (with software that does not apply Unicode Standard normalization) for ANGSTROM SIGN (U+212B), then the instances of LATIN CAPITAL LETTER A WITH RING ABOVE (U+00C5) would not be found.

In other words, there are equivalent sequences of Unicode characters that should be normalized, i.e. transformed into a unique Unicode-sanctioned representation of a character sequence called a NORMALIZATION FORM. Unicode provides a Unicode Normalization Algorithm, which essentially puts combining marks into a specific logical order and it defines decomposition and composition transformation rules to convert each string into one of four normalization forms. We will discuss here the two most relevant normalization forms: NFC and NFD.

The first of the three characters above is considered the NORMALIZATION FORM C (NFC), where C stands for composition. When the process of NFC normalization is applied to the character sequences in 2 and 3, both sequences are normalized into the PRE-COMPOSED character sequence in 1. Thus all three canonical character sequences are standardized into one composition form in NFC. The other central Unicode normalization form is the NORMALIZATION FORM D (NFD), where D stands for decomposition. When NFD is applied to the three examples above, all three, including importantly the single-character sequences in 1 and 2, are normalized into the DECOMPOSED multiple-sequence of characters in 3. Again, all three are then logically equivalent and therefore comparable and syntactically interoperable.

As illustrated, some characters in the Unicode Standard have alternative representations (in fact, many do), but the Unicode Normalization Algorithm can be used to transform certain sequences of characters into canonical forms to test for equivalency. To determine equivalence, each character in the Unicode Standard is associated with a combining class, which is formally defined as a character property called CANONICAL COMBINING CLASS which is specified in the Unicode

Character Database. The combining class assigned to each code point is a numeric value between 0 and 254 and is used by the Unicode Canonical Ordering Algorithm to determine which sequences of characters are canonically equivalent. Normalization forms, as very briefly described above, can be used to ensure character equivalence by ordering character sequences so that they can be faithfully compared.

It is very important to note that any software applications that is Unicode Standard compliant is free to change the character stream from one representation to another. This means that a software application may compose, decompose or reorder characters as its developers desire; as long as the resultant strings are canonically equivalent to the original. This might lead to unexpected behavior for users. Various players, like the Unicode Consortium, the W3C, or the TEI recommend NFC in most user-directed situations, and some software applications that we tested indeed seem to automatically convert strings into NFC.¹³ This means in practice that if a user, for example, enters <a> and <ò>, i.e. LATIN SMALL LETTER A at U+0061 and COMBINING GRAVE ACCENT at U+0300, this might be automatically converted into <à>, i.e. LATIN SMALL LETTER A WITH GRAVE at U+00E0.¹⁴

2.10 Pitfall: Absence of canonical equivalence

Although in most cases canonical equivalence will take care of alternative encodings of the same character, there are some cases in which the Unicode Standard decided against equivalence. This leads to identical characters that are not equivalent, like <ø> LATIN SMALL LETTER O WITH STROKE at U+00F8 and <ø> a combination of LATIN SMALL LETTER O at U+006F with COMBINING SHORT SOLIDUS OVERLAY at U+0037. The general rule followed is that extensions of Latin characters that are connected to the base character are not separated as combining diacritics. For example, characters like <ŋ n ŋ> or <đ d> are obviously derived from <n> and <d> respectively, but they are treated like new separate characters in the Unicode Standard. Likewise, characters like <ø> and <ç> are not separated into a base character <o> and <c> with an attached combining diacritic.

¹³ See the summary of various recommendation here: http://www.win.tue.nl/~aeb/linux/uc/nfc_vs_nfd.html.

¹⁴ The behavior of software applications can be quite erratic in this respect. For example, Apple's TextEdit does not do any conversion on text entry. However, when you copy and paste some text inside the same document in rich text mode (i.e. RTF-format), it will be transformed into NFC on paste. Saving a document does not do any conversion to the glyphs on screen, but it will save the characters in NFC.

2 Unicode pitfalls

Interestingly, and somewhat illogically, there are three elements, which are directly attached to their base characters, but which are still treated as separable in the Unicode Standard. Such characters are decomposed (in NFD normalization) in a base character with a combining diacritic. However, it is these cases that should be considered the exceptions to the rule. These three elements are the following:

- <◌̣>: the COMBINING CEDILLA at U+0327
This diacritic is for example attested in the precomposed character <ç> LATIN SMALL LETTER C WITH CEDILLA at U+00E7. This <ç> will thus be decomposed in NFC normalization.
- <◌̣>: the COMBINING OGONEK at U+0328
This diacritic is for example attested in precomposed <ą> LATIN SMALL LETTER A WITH OGONEK at U+0105. This <ą> will thus be decomposed in NFC normalization.
- <◌̥>: the COMBINING HORN at U+031B
This diacritic is for example attested in precomposed <ø> LATIN SMALL LETTER O WITH HORN at U+01A1. This <ø> will thus be decomposed in NFC normalization.

There are further combinations that deserve specific care because it is actually possible to produce identical characters in different ways without them being canonically equivalent. In these situations, the general rule holds, namely that characters with attached extras are not decomposed. However, in the following cases the “extras” actually exist as combining diacritics, so there is also the possibility to construct a character by using a base character with those combining diacritics.

- First, there are the combining characters designated as “combining overlay” in the Unicode Standard, like <◌̃> COMBINING TILDE OVERLAY at U+0334 or <◌̂> COMBINING SHORT STROKE OVERLAY at U+0335. There are many characters that look like they are precomposed with such an overlay, for example <ł ʙ đ Ɔ> or <ł ı j Ɔ>, or also the example of <ø> given at the start of this section. However, they are not decomposed in NFD normalization.
- Second, the same situation also occurs with combining characters designated as “combining hook”, like <◌̨> COMBINING PALATALIZED HOOK BELOW at U+0321. This element seems to occur in precomposed characters like <ɸ ɹ ƶ ƶ>. However, they are not decomposed in NFD normalization.

To harmonize the encoding in these cases it is not sufficient to use Unicode normalization. Additional checks are necessary, for example by using orthography profiles (see Chapter 5).

2.11 Pitfall: File formats

Unicode is a character encoding standard, but these characters of course actually appear inside some kind of computer file. The most basic Unicode-based file format is pure line-based text, i.e. strings of Unicode-encoded characters separated by line breaks (note that these line breaks are what for most people intuitively corresponds to paragraph breaks). Unfortunately, even within this apparently basic setting there exist a multitude of variants. In general, these different possibilities are well-understood in the software industry, and nowadays they normally do not lead to any problems for the end user. However, there are some situations in which a user is suddenly confronted with cryptic questions in the user interface involving abbreviations like LF, CR, BE, LE or BOM. Most prominently this occurs with exporting or importing data in several software applications from Microsoft. Basically, there are two different issues involved. First, the encoding of line breaks and, second, the encoding of the Unicode characters into code units and the related issue of endianness.

Line breaks

The issue with LINE BREAKS originated with the instructions necessary to direct a printing head of a physical printer to a new line. This involves two movements, known as CARRIAGE RETURN (CR, returning the printing head to the start of the line on the page) and LINE FEED (LF, moving the printing head to the next line on the page). Physically, these are two different events, but conceptually together they form one action. In the history of computing, various encodings of line breaks have been used (e.g. CR+LF, LF+CR, only LF, or only CR). Currently, all Unix and Unix-derived systems use only LF as code for a line break, while software from Microsoft still uses a combination of CR+LF. Today, most software applications recognize both options, and are able to deal with either encoding of line breaks (until rather recently this was not the case, and using the wrong line breaks would lead to unexpected errors). Our impression is that there is a strong tendency in software development to standardize on the simpler “only LF” encoding for line breaks, and we suggest that everybody use this encoding whenever possible.

Code units

The issue with CODE UNITS stems from the question how to separate a stream of binary ones and zero, i.e. bits, into chunks representing Unicode characters. A code unit is the sequence of bits used to encode a single character in an encoding. Depending on different use cases, the Unicode Standard offers three different approaches, called UTF-32, UTF-16 and UTF-8.¹⁵ The details of this issue is extensively explained in section 2.5 of the Unicode Core Specification Unicode Consortium (2014).

Basically, UTF-32 encodes each character in 32 bits (32 *binary units*, i.e. 32 zeros or ones) and is the most disk-space-consuming variant of the three. However, it is the most efficient encoding processing-wise, because the computer simply has to separate each character after 32 bits.

In contrast, UTF-16 uses only 16 bits per character, which is sufficient for the large majority of Unicode characters, but not for all of them. A special system of SURROGATES is defined within the Unicode Standard to deal with these additional characters. The effect is a more disk-space efficient encoding (approximately half the size), while adding a limited computational overhead to manage the surrogates.

Finally, UTF-8 is a more complex system that dynamically encodes each character with the minimally necessary number of bits, choosing either 8, 16 or 32 bits depending on the character. This represents again a strong reduction in space (particularly due to the high frequency of data using erstwhile ASCII characters, which need only 8 bits) at the expense of even more computation necessary to process such strings. However, because of the ever growing computational power of modern machines, the processing overhead is in most practical situations a non-issue, while saving on space is still useful, particularly for sending texts over the Internet. As an effect, UTF-8 has become the dominant encoding on the World Wide Web. We suggest that everybody uses UTF-8 as their default encoding.

A related problem is a general issue about how to store information in computer memory, which is known as ENDIANNESS. The details of this issue go beyond the scope of this book. It suffices to realize that there is a difference between BIG-ENDIAN (BE) storage and LITTLE-ENDIAN (LE) storage. The Unicode Standard offers a possibility to explicitly indicate what kind of storage is used by starting

¹⁵ The letters UTF stand for UNICODE TRANSFORMATION FORMAT, but the notion of “transformation” is a legacy notion that does not have meaning anymore. Nevertheless, the designation UTF (in capitals) has become an official standard designation, but should probably best be read as simply “Unicode Format.”

a file with a so-called `BYTE ORDER MARK` (BOM). However, the Unicode Standard does not require the usage of BOM, preferring other non-Unicode methods to signal to computers which kind of endianness is used. This issue only arises with UTF-32 and UTF-16 encodings. When using the preferred UTF-8, using a BOM is theoretically possible, but strongly dispreferred according to the Unicode Standard. We suggest that everyone tries to prevent the inclusion of BOM in your data.

2.12 Recommendations

Summarizing the pitfalls discussed in this chapter, we propose the following recommendations:

- To prevent strange boxes instead of nice glyphs, always install a few fonts with a large glyph collection and at least one fall-back font (see Section 2.4).
- Unexpected visual impressions of symbols does not necessarily mean that the actual encoding is wrong. It is mostly a problem of faulty rendering (see Section 2.5).
- Do not trust the names of codepoints as a definition of the character (see Section 2.7). Also do not trust Unicode blocks as a strategy to find specific characters (see Section 2.6)
- To ensure consistent encoding of texts, apply Unicode normalization (NFC or NFD, see Section 2.9).
- To prevent remaining inconsistencies after normalization, for example stemming from homoglyphs (see Section 2.8) or from missing canonical equivalence in the Unicode Standard (see Section 2.10) use orthography profiles (see Chapter 5).
- To deal with “tailored” grapheme clusters (Section 2.3), use Unicode Locale Descriptions, or orthography profiles (see Chapter 5).
- As a preferred file format, use Unicode Format UTF-8 in Normalization Form Composition (NFC) with LF line endings, but without byte order mark (BOM), whenever possible (see Section 2.11). This last nicely cryptic recommendation has T-shirt potential:

I prefer it
UTF-8 NFC LF no BOM

3 The International Phonetic Alphabet

3.1 Brief history

Established in 1886, the INTERNATIONAL PHONETIC ASSOCIATION (henceforth *Association*) has long maintained a standard alphabet, the INTERNATIONAL PHONETIC ALPHABET (henceforth *IPA*), which is a common standard in linguistics to transcribe sounds of spoken languages. It was first published in 1888 as an international system of phonetic transcription for oral languages and for pedagogical purposes. It contained phonetic values for English, French and German. Diacritics for length and nasalization were already present in this first version, and the same symbols are still used today.

Originally, the IPA was a list of symbols with pronunciation examples using words in different languages. In 1900 the symbols were first organized into a chart and were given phonetic feature labels (e.g. for manner of articulation among others ‘plosives’, ‘nasales’, ‘fricatives’, for place of articulation among others ‘bronchiales’, ‘laryngales’, ‘labiales’ and for vowels e.g. ‘fermées’, ‘mi-fermées’, ‘mi-ouvertes’, ‘ouvertes’). Throughout the last century, the structure of the chart has changed with increases in phonetic knowledge. Thus, similar to notational systems in other scientific disciplines, the IPA reflects facts and theories of phonetic knowledge that have developed over time. It is natural then that the IPA is modified occasionally to accommodate scientific innovations and discoveries. In fact, this is part of the Association’s mandate. These changes are captured in the recurring revisions to the IPA.

Over the years there have been several revisions, many minor. For example, the articulation labels – what often are called ‘features’ even though the IPA deliberately avoids this term – have changed (e.g. terms like ‘lips’, ‘throat’ or ‘rolled’ are not used anymore). Phonetic symbol values have changed (e.g. voiceless is no longer marked by <h>). Symbols have been dropped (e.g. the caret diacritic denoting ‘long and narrow’ is no longer used). And many symbols have been added to reflect contrastive sounds found in the world’s very diverse phonological systems. The current construction and use of the IPA are guided by principles

outlined in the HANDBOOK OF THE INTERNATIONAL PHONETIC ASSOCIATION (The International Phonetic Association 1999), henceforth simply called *Handbook*.

Today, the IPA is designed to meet practical linguistic needs and is used to transcribe the phonetic or phonological structure of languages. It is also used increasingly as a foreign language learning tool, as a standard pronunciation guide and as a tool for creating practical orthographies of previously unwritten languages. The IPA suits many linguists' needs because:

- it is intended to be a set of symbols for representing all possible sounds in the world's (spoken) languages;
- its chart has a linguistic basis (specifically a phonological bias) rather than just being a general phonetic notation scheme;
- its symbols can be used to represent distinctive feature combinations;¹
- its chart provides a summary of linguists' agreed-upon phonetic knowledge.

Several styles of transcription with the IPA are possible, as illustrated in the *Handbook*, and they are all valid.² Therefore, there are different but equivalent transcriptions, or as Ladefoged (1990: 64) captures it, "perhaps now that the Association has been explicit in its eclectic approach, outsiders to the Association will no longer speak of *the* IPA transcription of a given phenomenon, as if there were only one approved style." Clearly not all phoneticians agree (or will likely ever agree) on all aspects of the IPA or on transcription practices. As noted above, there have been several revisions in the IPA's long history, but the current version (2005) is strikingly similar to the 1926 version, showing the viability of the IPA as a common standard for linguistic transcription.

3.2 Premises

Any IPA transcription is based on two premises: (i) that it is possible to describe the acoustic speech signal (sound waves) in terms of sequentially ordered discrete segments, and, (ii) that each segment can be characterized by an articulatory target.

¹ Although the chart uses traditional manner and place of articulation labels, the symbols can be used as a representation of any defined bundle of features, binary or otherwise, to define phonetic dimensions.

² For illustration of different possibilities, see the 29 languages and their transcriptions in the original and initial *Illustrations of the IPA* (The International Phonetic Association 2007: 41–154)

Once spoken language data are segmented, the IPA provides symbols to unambiguously represent phonetic details. However, since phonetic detail could potentially include anything, e.g. something like ‘deep voice’, the IPA restricts phonetic detail to linguistically relevant aspects of speech. Phonological consideration thus inextricably play a roll in transcription. In other words, phonetic observations beyond quantitative acoustic analysis are always made in terms of some phonological framework.

Today, the IPA chart reflects a linguistic theory grounded in principles of phonological contrast and in knowledge about the attested linguistic variation. This fact is stated explicitly in several places, including in the *Report on the 1989 Kiel convention* published in the *Journal of the International Phonetic Association* (Association et al. 1989: 67-68):

The IPA is intended to be a set of symbols for representing all the possible sounds of the world’s languages. The representation of these sounds uses a set of phonetic categories which describe how each sound is made. These categories define a number of natural classes of sounds that operate in phonological rules and historical sound changes. The symbols of the IPA are shorthand ways of indicating certain intersections of these categories.

and in the *Handbook* (The International Phonetic Association 1999: 18):

[...] a symbol can be regarded as a shorthand equivalent to a phonetic description, and a way of representing the contrasting sounds that occur in a language. Thus [m] is equivalent to ‘voiced bilabial nasal’, and is also a way of representing one of the contrasting nasal sounds that occur in English and other languages. [...] When a symbol is said to be suitable for the representation of sounds in two languages, it does not necessarily mean that the sounds in the two languages are identical.

From its earliest days the Association aimed to provide “a separate sign for each distinctive sound; that is, for each sound which, being used instead of another, in the same language, can change the meaning of a word” (The International Phonetic Association 1999: 27). Distinctive sounds became later known as PHONEMES and the IPA has developed historically into a notational devise with a strictly segmental phonemic view. A phoneme is an abstract theoretical notion derived from an acoustic signal as produced by speakers in the real world. Therefore the IPA contains a number of theoretical assumptions about speech and how to transcribe speech in written form.

3.3 Principles

Transcription has essentially two parts. The first is a text containing IPA symbols and the second is a set of conventions for interpreting those symbols (and their combinations).

First, the symbols of IPA crucially distinguish between letter-like symbols and diacritics. The use of the letter-like symbols to represent a language's sounds is guided by the principle of contrast; where two words are distinguishable by phonemic contrast, those contrasts should be transcribed with different letter symbols (not diacritics). Allophonic distinction falls under the rubric of diacritically-distinguished symbols, e.g. the difference in English between and unrealized /p/ in [stop] vs. an aspirated /p/ in [sp^hot] should be marked with a diacritic. In sum:

- Different letter-like symbols should be used whenever a language employs two sound contrastively.
- When two sounds in a language are not known to be contrastive, the same symbol should be used to represent these sounds. Diacritics may be used to distinguish different articulations when necessary.

Yet, in some situations diacritics can be used to mark phonemic contrasts. The *Handbook* recommends to limit the use of phonemic diacritics to the following situations:

- denoting length, stress and pitch;
- obviating the design of a (large) number of new symbols when a single diacritic suffices (e.g. nasalized vowels, aspirated stops).

The interpretation of the IPA symbols in specific usage is not trivial. Although the articulatory properties of the IPA symbols themselves are rather succinctly summarised by the normative description in the *Handbook*, it is common in practical applications that the transcribed symbols do not precisely represent the phonetic value of the sounds that they represent. So an IPA symbol <t> in one transcription is not always the same as an IPA <t> in another transcription (or even within a single transcription). The interpretation of any particular <t> is mostly a language-specific convention (or even resource-specific and possibly even context-specific), a fact which – unfortunately – is in most cases not made explicit by users of the IPA.

There are different reasons for this difficulty to interpret IPA symbols, all officially sanctioned by the IPA. An important principle of the IPA is that different

representations resulting from different underlying analyses are allowed. Because the IPA does not provide phonological analyses for specific languages, the IPA does not define a single “correct” transcription system. Rather, the IPA aims to provide a resource that allows users to express any analysis so that it is widely understood. Specifically, the IPA allows for both a ‘narrow’ phonetic transcription and a ‘broad’ phonological transcription. A ‘narrow’ phonetic transcription may freely use the all symbols in the IPA chart with direct reference to their phonetic value, i.e. the transcriber can indicate with the symbols <ɱ̥> that the phonetic value of the attested sound is a simultaneous labial and velar closure which is voiceless and contains nasal airflow, independently of the phonemic status of this sound in the language in question.

In contrast, the basic goal of a phonemic “broad” transcription is to distinguish all words in a language with a minimal number of transcription symbols (Abercrombie 1964). A phonemic transcription includes the conventions of a particular language’s phonological rules. These rules determine the realization of that language’s phonemes. For example, in the transcription of German, Dutch, English and French a symbol <t> might be used for the voiceless plosive in the alveolar area. This symbol is sufficient for a succinct transcription of these languages because there is no further phonemic subdivisions in this domain in either of these languages. However, the language-specific realization of this consonant is closer to [tʰ], [t], [tʰ] and [t̚], respectively. Similarly, the five vowels of Greek can be represented phonemically in IPA as /ieaou/, though phonetically they are closer to [ieaou]. Similarly, the Japanese five vowel system can also be transcribed in IPA as /ieaou/, while the phonetic targets in this case are closer to [ieaou].

Note also that there can be different systems of phonemic transcription for the same variety of a language, so two different resources on the ‘same’ language might use different symbols that represent the same sound. The differences may result from the fact that more than one phonetic symbol may be appropriate for a phoneme, or the differences may be due to different phonemic analyses, e.g. Standard German’s vowel system is arguably contrastive in length, tenseness or vowel quality. Finally, even within a single phonemic transcription a specific symbol can have different realizations because of allophonic variation which might not be explicitly transcribed.

In sum, there are three different reasons why phonemically-based IPA transcription is difficult to interpret:

- A symbol represents the phonemic status, and not necessarily the precise phonetic realization. So, different transcriptions might use the same symbol for different underlying sounds.
- Any symbol that is used for a specific phoneme is not necessarily unique, so different transcriptions might use different symbols for the same underlying sound.
- Allophonic variation can be disregarded in phonemic transcription, so the same symbol within a single transcription can represent different sounds.

Ideally, all such implicit conventions of a phonemic transcription would be explicitly codified. This could very well be performed by using an orthography profile (see Chapter 5), linking the selected phonemic transcription symbols to ‘narrow’ phonetic transcriptions, possibly also including specifications of contextual interpretation.

3.4 Computational representation

In 1989 an IPA revision convention was held in Kiel, Germany. As per other revisions, there was expansion and revisions to phonetic symbols in the IPA chart. Notably the marking of tone was extended with the addition of a second system for marking linguistic tones (the so-called ‘Chao’ tone letters). Importantly, however, this was also the first revision to address issues of computational representations for the IPA.

Prior to the Kiel Convention for the modern revision of the IPA in 1989, Wells (1987) collected and published practical approaches to coded representations of the IPA, which dealt mainly with the assignment of characters on the keyboard to IPA symbols. The process of assigning standardized “computer codes” to phonetic symbols was assigned to the WORKGROUP ON COMPUTER CODING (henceforth *working group*) at the Kiel Convention. This working group was assigned to consider the following tasks (Esling 1990; Esling & Gaylord 1993):

- determining how to represent the IPA numerically
- developing a set of numbers to refer to the IPA symbols unambiguously
- providing each symbol a unique name (intended to provide a mnemonic description of that character’s shape)

The identification of IPA symbols with unique identifiers was a first step in formalizing the IPA computationally because it would give each symbol an unambiguous numerical identifier called an IPA NUMBER. The numbering system was

to be comprehensive enough to support future revisions of the IPA, including symbol specifications and diacritic placement. The application of diacritics was also to be made explicit.

Although the Association had never officially approved a set of names for the IPA symbols, each IPA symbol received a unique IPA NAME. Many symbols already had an informal name (or two) used by linguists, but consensus on symbol names was growing due to the recent publication of the *Phonetic Symbol Guide* (Pullum & Ladusaw 1986). Thus most of the IPA symbol names were taken from that source (The International Phonetic Association 1999: 31).

The working group insightfully decided that the computing-coding convention for the IPA should be independent of computer environments or formats (e.g. ASCII), i.e. the IPA Number was not meant to be implemented directly in a computer encoding. The working group report's declaration includes the following explanatory remarks (Association et al. 1989: 82):

The recommendation of a 7-bit ASCII or 8-bit extended-ASCII coding system would be short-sighted in view of development towards 16-bit and 32-bit processors. In fact, any specific recommendations would tie the Association to a stage of technological development which is bound to be outdated long before the next revision of the handbook.

Thus the coding convention was not meant to address the engineering aspects of the actual encoding in computers (cf. Anderson 1984). However, it was meant to serve as a basis for a communication-interchange standard for creating mapping tables from various computer encodings, fonts, phonetic-character-set software, etc., to common IPA Numbers, and thus symbols.³

Furthermore, the assignment of computer codes to IPA symbols was meant to represent an unbiased formulation. The Association here played the role of an international advisory body and it stated that it should not recommend a particular existing system of encoding. In fact, during this time there were a number of coding systems used (see Section 1.2), but none of them had a dominant international position. The differences between systems were also either too great or too subtle to warrant an attempt at combining them (Association et al. 1989).

³ Remember, at this time in the late 1980s there was no stable multilingual computing environment. But some solution was needed because scholars were increasingly using personal computers for their research and many were quickly adopting electronic mail or discussion boards like Usenet as a medium for international exchanges. Most of these systems ran on 8-bit hardware systems using a 7-bit ASCII character encoding.

3 *The International Phonetic Alphabet*

The working group assigned each IPA symbol to a unique three-digit number, i.e. an IPA Number. Encoded in this number scheme implicitly includes information about the status of each symbol (see below). The IPA numbers were listed with the IPA symbols and they were also illustrated in IPA chart form (see Esling & Gaylord 1993: 84 or The International Phonetic Association 2007: App. 2). The numbers were assigned in linear order (e.g. [p] 101, [b] 102, [t] 103...) following the IPA revision of 1989 and its update in 1996. Although the numbering scheme still exists, in practice it is superseded by the Unicode codification of symbols.

The working group made the decision that no IPA symbol, past or present, could be ignored. The comprehensive inclusion of all IPA symbols was to anticipate the possibility that some symbols might be added, withdrawn, or reintroduced into current or future usage. For example, in the 1989 revision voiceless implosives < ɓ, ɗ, ɠ, ɡ > were added; in the 1993 revision they were removed. Ligatures like < ɥ, ɹ̥ > are included as formerly recognized IPA symbols; they are assigned to the “200 series” of IPA numbers as members of the group of symbols formerly recognized by the IPA. To ensure backwards compatibility, legacy IPA symbols would retain an IPA Number and an IPA Name for reference purposes. As we discuss below, this decision is later reflected in the Unicode Standard as many legacy IPA symbols still reside in the IPA EXTENSIONS block.

The IPA Number is simply expressed as a “three-digit number numerical directory of digit triples” The numbering scheme specifies three-digit codes, the first digit of which indicates the symbol’s category (Esling 1990; Esling & Gaylord 1993):

- 100s for accepted IPA consonants
- 200s for former IPA consonants and non-IPA symbols
- 300s for vowels
- 400s for segmental diacritics
- 500s for suprasegmental symbols
- 600s-800s for future specifications
- 900s for escape sequences

After a symbol is categorized, it is assigned a number sequentially, e.g. [i] 301, [e] 302, [ɛ] 303. The system allows for the addition of new symbols within the various series by appending them, e.g. [v] 184. Former non-IPA symbols (or often-used but non-official IPA symbols) for consonants, vowels and diacritics are numbered from x99 backwards. For example, the voiceless and voiced postalveolar affricates and fricatives < ɕ, ʝ, ʑ, ʒ > are assigned the IPA numbers 299, 298, 297 and 296, respectively, because they are not sanctioned IPA symbols.

The assignment of the IPA numbers to IPA symbols provided the basis for uniquely identifying the set of past and present IPA symbols as a type of computational representational standard of the IPA. Within each revision of the IPA, the coding defines a closed and clearly defined set of characters. The benefits of this standardization are clear in at least two ways: it is used in translation tables that reference ASCII representations of the IPA, and this early computational representation of the IPA became the basis X-SAMPA and for the inclusion of the IPA into the Unicode Standard version 1.0.

3.5 SAMPA and X-SAMPA

True to the working group's aim, the IPA numbers provided a mechanism for a communication interchange standard for creating mapping tables to various computer encodings. For example, the IPA coding system was used as a mapping system in the creation of SAMPA (Wells et al. 1992), an ASCII representation of the IPA symbols.

For a long time, linguists, like all other computer users, were limited to ASCII-encoded 7-bit characters, which only includes Latin characters, numbers and some punctuation and symbols. Restricted to these standard character sets that lacked IPA support or other language-specific graphemes that they needed, linguists devised their own solutions.⁴ For example, some chose to represent unavailable graphemes with substitutes, e.g. the combination of <ng> to represent <ŋ>. Tech-savvy linguists redefined selected characters from a character encoding by mapping custom made fonts to specific code points.⁵ However, one linguist's electronic text would not render properly on another linguist's computer without access to the same font. Furthermore, if two character encodings defined two character sets differently, then data could not be reliably and correctly displayed. This is a commonly encountered example of the non-interoperability of data and data formats.

One solution was the ASCII-ification of the IPA, which simply involved defining keyboard-able sequences consisting of ASCII combinations as an IPA symbol codings. For example, Wells 1987 provides an in-depth description of IPA codings from country-to-country. Later ASCII-IPAs include Kirshenbaum (cre-

⁴ Early work addressing the need for a universal computing environment for writing systems and their computational complexity is discussed in Simons (1989). A more recent survey of practical recommendations for language resources, including notes on encoding, can be found in Bird & Simons (2003)

⁵ For example, SIL's popular font SIL IPA 1990

ated in 1992 in a Usenet group and named after its lead developer who was at Hewlett-Packard Laboratories) and Worldbet (published in 1993 by Hieronymus 1993, who was at AT&T Laboratories). The most successful effort was SAMPA (Speech Assessment Methods Phonetic Alphabet), which was created between 1988–1991 in Europe to represent IPA symbols with ASCII character sequences (Wells 1987; Wells et al. 1992), using e.g. <p\> for [ɸ]. SAMPA was developed by a group of speech scientists from nine countries in Europe and it constituted the ASCII-IPA symbols needed for phonemic transcription of the principal European languages (Wells 1995). It is still widely used in language technology.

Two problems with SAMPA are that (i) it is only a partial encoding of the IPA and (ii) it encodes different languages in separate data tables, instead of using a universal alphabet, like IPA. SAMPA tables were developed as part of a European Commission-funded project to address technical problems like electronic mail exchange (what is now simply called email). SAMPA is essentially a hack to work around displaying IPA characters, but it provided speech technology and other fields a basis that has been widely adopted and often still used in code. So, SAMPA is a collection of tables to be compared, instead of a large universal table representing all languages.

An extended version of SAMPA, called X-SAMPA, set out to include every symbol, including all diacritics, in the IPA chart (Wells 1995). X-SAMPA is considered more universally applicable because it consists of one table that encodes all characters in IPA. In line with the principles of the IPA, SAMPA and X-SAMPA include a repertoire of symbols. These symbols are intended to represent phonemes rather than all allophonic distinctions. Additionally, both ASCII-ifications of IPA – SAMPA and X-SAMPA – are (reportedly) uniquely parsable (Wells 1995). However, like the IPA, X-SAMPA has different notations for encoding the same phonetic phenomena (cf. Section 4.5).

SAMPA and X-SAMPA have been widely used for speech technology and as an encoding system in computational linguistics. In fact, they are still used in popular software packages that require ASCII input, e.g. RuG/L04 and SplitsTree4.⁶

3.6 The need for a single multilingual environment

In hindsight it is easy to lose sight of how impactful 30 years of technological development have been on linguistics, from theory development using quantitative means to pure data collection and dissemination. But at the end of the 1970s,

⁶ See <http://www.let.rug.nl/kleiweg/L04/> and <http://www.splitsree.org/>, respectively

virtually no ordinary working linguist was using a personal computer (Simons 1996). Personal computer usage, however, dramatically increased throughout the 1980s. By 1990, dozens of character sets were in common use. They varied in their architecture and in their character repertoires, which made things a mess.

During the 1980's, it became increasingly clear that an adequate solution to the problem of multilingual computing environments was needed. Linguists were on the forefront of addressing this issue because they faced these challenges head-on by wishing to publish and communicate electronic text with phonetic symbols which were not included in basic ASCII. One only needs to look at facsimiles of older electronic documents to see exotic symbols written in by hand after the preparation of typed version.

A major benefit of the standardization of the IPA in a computational representation by the Kiel working group is that it provided the basis for a formal proposal to be submitted to various international standards organizations, several of which were trying to tackle (and in a sense 'win') the multilingual computing environment problem (cf. Section 1.2). Basically, everyone – from corporations to governments to language scientists – wanted a single unified multilingual character encoding set for all the world's writing systems, even if they did not understand or appreciate the challenges involved in creating and adopting a solution.

Industry was starting to tackle the issues involved in developing a single multilingual computing environment on a variety of fronts, including the then new technology of bitmap fonts and the creation of Font Manager and Script Manager by Apple (Computer 1985, 1986, 1988). As noted above, around this time linguists were developing work-arounds such as SAMPA, so that they could communicate IPA transcription and use ASCII-based software. Some linguists formalized the issues of multilingual text processing from a computational perspective (Anderson 1984; Becker 1984; Simons 1989). The study of writing systems was also being invigorated (Sampson 1985) by the computational challenges in making computers work in a multilingual environment. The engineering problems and solutions had been spelled out years before, e.g. a two-byte encoding for multilingual text (Anderson 1984). Although languages vary to an astounding extent (cf. Evans & Levinson 2009), writing systems are quite similar formally and the issue of formal representation of the world's orthographic systems had already been addressed (Simons 1989).

After the Kiel Convention in 1989, the working group assisted the International Phonetic Association in representing the IPA to the INTERNATIONAL STANDARDS ORGANISATION (ISO) and to the TEXT ENCODING INITIATIVE (TEI) (Esling & Gaylord 1993). The working group's formalization of the IPA, i.e. a full listing of

agreed upon computer codings for phonetic symbols, was used in developing writing system descriptions, which were at the time being solicited for scripts to be included in the new multilingual international character encoding standards. The working group for ISO/IEC 10646 and Unicode were two such initiatives.

In the historical context of the IPA being considered for inclusion in ISO/IEC 10646, it is important to realize that there were a variety of sources (i.e. not just from the Association) which submitted character proposals for phonetic alphabets. These proposals, including the one from the Association via the Kiel working group, were considered as a whole by the ISO working groups that were responsible for incorporating a phonetic script into the universal character set (UCS). The ISO working groups that were responsible for assigning a phonetic character set then made their own submissions as part of a review process by ISO for approval based on both “informatic” and phonetic criteria (Esling & Gaylord 1993: 86).

Character set ISO/IEC 10646 was approved by ISO, including the phonetic characters submitted to them in May 1993. The set of IPA characters were assigned UCS codes in 16 bit representation (in hexadecimal) and were published Tables 2 and 3 in Esling & Gaylord 1993, which include a graphical representation of the IPA Symbol, its IPA Name, phonetic description, IPA Number, UCS Code and AFII Code.⁷ When the character sets of ISO/IEC 10646 and the Unicode Standard later converged (see Section 1.4), the IPA proposal was included in the Unicode Standard Version 1.0 – largely as we know it today.

With subsequent revisions to the IPA, one might have expected that the Unicode Consortium would update the Unicode Standard in a way that is inline with the development linguistic insights. However, updates that go against the principle of maintaining backwards compatibility lose out, i.e. it is more important to deal with the pitfalls created along the way than it is to change the standard. Therefore, many of the pitfalls we encounter today when using Unicode IPA are historic relics that we have to come to grips with.

It was a long journey, but the goal of achieving a single multilingual computing environment has largely been accomplished. Current users, however, must cope with the pitfalls that were dug along the way, as will be discussed in the next chapter. As the Association foresightfully remarked about Unicode:

“When this character set is in wide use,

⁷ The ASSOCIATION FOR FONT INFORMATION INTERCHANGE (AFII) was an international database of glyphs created to promote the standardization of font data required to produce ISO/IEC 10646.

it will be the normal way to encode IPA symbols.”

(The International Phonetic Association 1999: 164).

3.7 Unicode and ISO 10646

In the late 1980s, a universal character set was being developed by what is now referred to as the Unicode Consortium (officially incorporated in January 1991). This consortium consisted largely, although not entirely, of major US corporations, with the aim of overcoming the inoperability of different coded character sets and their costly hinderance for developing multilingual software development and for internationalization efforts. Commercial importance of course drove the early inclusion of Latin, non-Latin, and some exotic scripts; see the table of commercial importance as measured by GDP of countries using certain writing systems (**unicode88**).

The original Unicode manifesto is **unicode88.pdf**⁸ Its aim was for a reliable international multilingual text encoding standard that would encompass all scripts of the world, or in the author’s own words, “a new, world-wide ASCII”. An in-depth history of Unicode, highlighting interesting facts like its first text prototypes at Apple and its incorporation into TrueType, is retold online.⁹

Unicode 88 provided the basic principles for the Unicode Standard’s design – pushing for 16 bit representations of characters with a clear distinction between characters and glyphs. Some of the contents of this status proposal of 1988 were reworked for inclusion in the early Unicode Standard pre-publication drafts and by August 1990, the proposal was in a (very) rough draft format. Its editors and the Unicode Working Group (the predecessor to the Unicode Technical Committee) worked together to lay out the the proposed standard’s structure and content. At this time, the proposal contained no code charts nor block descriptions.

The other major player in developing a universal character set was the ISO working group from the International Standards Organization (ISO), based in Europe, which was responsible for ISO/IEC 10646. This character set standard was composed in 1989 and a draft was published in 1990.¹⁰ The ‘Universal Multiple-Octet Coded Character Set’ or simply UCS was the first officially standardized

⁸ <http://www.unicode.org/history/unicode88.pdf>

⁹ <http://www.unicode.org/history/earlyyears.html>

¹⁰ http://www.iso.org/iso/catalogue_detail.htm?csnumber=56921

3 *The International Phonetic Alphabet*

character encoding with the aim of including all characters from all writing systems.¹¹

ISO/IEC 10646 is partly based on ISO/IEC 8859, a series of of ASCII-based standard character encodings published in 1987 that use a single bit 8-byte character set. Each part of the standard, e.g. 8859-1, 8859-5, 8859-6, encodes characters to support different languages' writing systems, e.g. Latin-1 Western European, Latin/Cyrillic, Latin/Arabic, respectively. Being a joint effort by the International Organization for Standardization (ISO) and the International Electrotechnical Commission (IEC), the aim of the standard is reliable information exchange. So, again, issues of phonetic symbol encoding, typography, etc., were ignored – or perhaps more properly put, not commercially driven at this early stage.

Intended for the major Western European languages, ISO/IEC 8859 was an extension of the ASCII character encoding standard, which included the English alphabet, numerals and computer control characters (e.g. beep, space, carriage return). By extending ASCII's 7-bit system to 8-bit, the character repertoire of each of ISO/IEC 8859 character set was doubled from 128 to 256 characters. Each character set defined a mapping between digital bit patterns and characters, which are visually rendered on screen as graphic symbols. ASCII was shared between ISO/IEC 8859 character sets, but the characters in the extra bit patterns were different. Thus an aim of the ISO working group responsible for ISO/IEC 10646 was to bring all characters in all writings systems into a single unified encoding.

In 1991, the Unicode Consortium and the ISO Working Group for ISO/IEC 10646 decided to create a single universal standard for encoding multilingual text.¹² The two character sets converged, resulting in mutually acceptable changes to both, and each group keeps versions of their respective character codes and encoding forms synchronized.¹³ Although each standard has its own form of reference and the terminology in each may differ slightly, the practical difference is that the Unicode Standard is a formal implementation of ISO/IEC 10646 and imposes additional constraints on its implementation. The Unicode Standard includes character data, algorithms and specifications, outside the scope of ISO/IEC 10646, which ensure, when properly implemented in software applications and platforms, that characters are treated uniformly.

The incorporation of the Unicode Standard into the international encoding standard ISO 10646 was approved by ISO as an International Standard in June

¹¹ <http://www.nada.kth.se/i18n/ucs/unicode-iso10646-oview.html>

¹² <http://unicode.org/book/appC.pdf>

¹³ <http://www.unicode.org/versions/>

1992.¹⁴ The joint Unicode and ISO/IEC 10646 standard has become *the* universal character set and it is a single multilingual environment for the majority of the world's written languages. Its formal implementation has also been vital to the rise of a multi-lingual Internet.

¹⁴ <http://www.unicode.org/versions/Unicode1.0.0/Notice.pdf>

4 IPA meets Unicode

4.1 The twain shall meet

The International Phonetic Alphabet (IPA) is a common standard in linguistics to transcribe sounds of spoken language into discrete segments using a Latin-based alphabet. Although IPA is reasonably easily adhered to with pen and paper, it is not trivial to encode IPA characters electronically. In this chapter we discuss various pitfalls with the encoding of IPA in the Unicode Standard. We will specifically refer to the 2005 version of the IPA (International Phonetic Association 2005) and the 7.0 version of Unicode (Unicode Consortium 2014).

As long as a transcription is only directed towards phonetically trained eyes, then all the details of the Unicode-encoding are unimportant. For a linguist reading an IPA transcription, many of the details that will be discussed in this chapter might seem like hair-splitting trivialities. However, if IPA transcriptions are intended to be used across resources (e.g. searching similar phenomena across different languages) then it becomes crucial that there are strict encoding guidelines. Our main goal in this chapter is to present the encoding issues and propose recommendations for a “strict” IPA encoding for situations in which cross-resource consistency is crucial.

There are several pitfalls to be aware of when using the Unicode Standard to encode IPA. As we have said before, from a linguistic perspective it might sometimes look like the Unicode Consortium is making incomprehensible decisions, but it is important to realize that the consortium has tried and is continuing to try to be as consistent as possible across a wide range of use cases, and it does place linguistic traditions above other orthographic choices. Furthermore, when we look at the history of how the IPA met Unicode, we see that many of the decisions for IPA symbols in the Unicode Standard come directly from the International Phonetic Association itself. Therefore, many pitfalls that we will encounter have their grounding in the history of the principles of the IPA, as well as in the technological considerations involved in creating a single multilingual encoding. In general, we strongly suggest to linguists to not complain about any decisions in the Unicode Standard, but to try and understand the rationale

of the International Phonetic Association and the Unicode Consortium (both of which are almost always well-conceived in our experience) and devise ways to work with any unexpected behavior.

Many of the current problems derive from the fact that the IPA is clearly historically based on the Latin script, but different enough from most other Latin-based writing systems to warrant special attention. This ambivalent status of the IPA glyphs (partly Latin, partly special) is unfortunately also attested in the treatment of IPA in the Unicode Standard. In retrospect, it might have been better to consider the IPA (and other transcription systems) to be a special kind of script within the Unicode Standard, and treat the obvious similarity to Latin glyphs as a historical relic. All IPA glyphs would then have their own code points, instead of the current situation in which some IPA glyphs have special code points, while others are treated as being identical to the regular Latin characters. Yet, the current situation, however unfortunate, is unlikely to change, so as linguists we must learn to deal with the specific pitfalls of IPA within the Unicode Standard.

4.2 Pitfall: No complete IPA code block

The ambivalent nature of IPA glyphs arises because, on the one hand, the IPA uses Latin-based glyphs like <a>, or <p>. From this perspective, the IPA seems to be just another orthographic tradition using Latin characters, all of which do not get a special treatment within the Unicode Standard (just like e.g. the French, German, or Danish orthographic traditions do not have a special status). On the other hand, the IPA uses many special symbols (like turned <v>, mirrored <ɹ> and/or extended <ɸ> Latin glyphs) not found in any other Latin-based writing system. For this reason a special block with code points, called IPA EXTENSIONS was already included in the first version of the Unicode Standard (Version 1.0 from 1991).

As explained in Section 2.6, the Unicode Standard code space is subdivided into character blocks, which generally encode characters from a single script. However, as is illustrated by the IPA, characters that form a single writing system may be dispersed across several different character blocks. With its diverse collection of symbols from various scripts and diacritics, the IPA is spread across 12 blocks in the Unicode Standard:¹

¹ This number of blocks depends on whether only IPA-sanctioned symbols are counted or if the phonetic symbols commonly found in the literature are also included, see Moran 2012: Appendix C. The 159 characters from 12 code blocks shown here are the characters proposed for “strict” IPA encoding, as discussed in Section 4.13.

- [illegible]

4.3 Pitfall: IPA homoglyphs in Unicode

Another problem is the large number of homoglyphs, i.e. different characters that have highly similar glyphs (or even completely identical glyphs, depending on the font rendering). For example, a user of a Cyrillic computer keyboard should ideally not use the <a> CYRILLIC SMALL LETTER A at code point U+0430 for IPA transcriptions, but instead use the <a> LATIN SMALL LETTER A at code point U+0061, although visually they are mostly indistinguishable, and the Cyrillic character is more easily typed on a Cyrillic keyboard. Some further prob-

lematic homoglyphs related to encoding IPA in the Unicode Standard are the following.

- The uses of the apostrophe has led to long discussions on the Unicode Standard email list. An English keyboard inputs the symbol <’> APOSTROPHE at U+0027, although the preferred Unicode apostrophe is the <◌’> RIGHT SINGLE QUOTATION MARK at U+2019.² However, the glottal stop-/glottalization/ejective marker is yet another completely different character, namely <◌ʔ>, i.e. the MODIFIER LETTER APOSTROPHE at U+02BC, which unfortunately looks mostly extremely similar to U+2019.
- Another problem is the <◌ʰ> MODIFIER LETTER REVERSED GLOTTAL STOP at U+02C1 vs. the <◌ʳ> MODIFIER LETTER SMALL REVERSED GLOTTAL STOP at U+02E4. Both appear in various resources representing phonetic data online. This is thus a clear example for which the Unicode Standard does not solve the linguistic standardization problem.
- Linguists are also unlikely to distinguish between the <ə> LATIN SMALL LETTER SCHWA at code point U+0259 and <ɐ> LATIN SMALL LETTER TURNED E at U+01DD.
- The alveolar click <!> at U+01C3 is of course often simply typed as <!> EXCLAMATION MARK at U+0021.³
- The dental click <|>, in Unicode known as LATIN LETTER DENTAL CLICK at U+01C0, is often simply typed as <|> VERTICAL LINE at U+007C.
- For the marking of length there is a special Unicode characters, namely <◌:̯> MODIFIER LETTER TRIANGULAR COLON at U+02D0. However, typing <◌:̯> COLON at U+003A is of course much easier.
- There are two mostly identical-looking Unicode characters for the superscript <^h>: the COMBINING LATIN SMALL LETTER H at U+036A and the MODIFIER LETTER SMALL H at U+02B0. Making the situation even more problematic is that they have different behavior (see Section 4.9). To harmonize the behavior of <^h> with other superscript letters, we propose to standardize on the modifier letter at U+02B0 (see Section 4.10).

Conversely, non-linguists are unlikely to distinguish any semantic difference between an open back unrounded vowel, which is encoded in IPA with a ‘single-story’ <q> LATIN SMALL LETTER ALPHA at U+0251, and the open front unrounded

² Note that many word processors (like Microsoft Word) by default will replace straight quotes by curly quotes, depending on the whitespace around it.

³ In the Unicode Standard the <!> at U+01C3 is labeled LATIN LETTER RETROFLEX CLICK, but in IPA that glyph is used for an alveolar or postalveolar click (not retroflex). This naming is probably best seen as an error in the Unicode Standard. For the “real” retroflex click, see Section 4.12.

vowel, which is encoded in IPA with a ‘double-story’ <a> LATIN SMALL LETTER A at U+0061, basically treating them as homoglyphs, although they are different phonetic symbols. But even among linguists this distinction leads to problems. For example, as pointed out by Mielke (2009), there is a problem stemming from the fact that about 75% of languages are reported to have a five-vowel system (Maddieson 1984). Historically, linguistic descriptions tend not to include precise audio recording and measurements of formants, so this may lead one to ask if the many <a> characters that are used in phonological description reflects a transcriptional bias. The common use of <a> in transcriptions could be in part due to the ease of typing the letter on an English keyboard (or for older descriptions, the typewriter). We found it to be exceedingly rare that a linguist uses <ɑ> for a low back unrounded vowel.⁴ They simply use <a> as long as there is no opposition to <ɑ>.

Making things even more problematic, there is an old typographic tradition that the double-story <a> uses a single-story <ɑ> in italics. This leads to the unfortunate effect that even in many well-designed fonts the italics of <a> and <ɑ> use the same glyph. For example, in Linux Libertine (the font of this book) the italics of these characters are highly similar <a> and <ɑ>, while in Charis SIL they are identical: <a> and <ɑ>. If this distinction has to be kept upright in italics, the only solution we can currently offer is to use SLANTED glyphs (i.e. artificially italicized glyphs) instead of real italics (i.e. special italics glyphs designed by a typographer).⁵

⁴ One example is Vidal2001a in which the author states: “The definition of Pilagá /a/ as [+back] results from its behavior in certain phonological contexts. For instance, uvular and pharyngeal consonants only occur around /a/ and /o/. Hence, the characterization of /a/ and /o/ as a natural class of (i.e., [+back] vowels), as opposed to /i/ and /e/.”

⁵ For example, the widely used IPA font Doulos SIL (http://scripts.sil.org/cms/scripts/page.php?item_id=DoulosSIL) does not have real italics. This leads some word-processing software, like Microsoft Word, to produce slanted glyphs instead. That particular combination of font and software application will thus lead to the desired effect distinguishing <a> from <ɑ> in italics. However, note that when the text is transferred to another font (i.e. one that includes real italics) and/or to another software application (like Apple Pages, which does not perform slanting), then this visual appearance will be lost. In this case we are thus still in the pre-Unicode situation in which the choice of font and rendering software actually matters. The ideal solution from a linguistic point of view would be the introduction of a new IPA code point for a different kind of <a>, which explicitly specifies that it should still be rendered as a double-story character when italicized. After informal discussion with various Unicode players, our impression is that this highly restricted problem is not sufficiently urgent to introduce even more <a> homoglyphs in Unicode (which already lead to much confusion, see Section 2.8).

4.4 Pitfall: Homoglyphs in IPA

Reversely, there are a few cases in which the IPA distinguishes different phonetic concepts, but the visual characters used by the IPA look very much alike. Such cases are thus homoglyphs in the IPA itself, which of course need different encodings.

- The dental click <|> and the indication of a minor group break <|> look almost the same in most fonts. For a proper encoding, the LATIN LETTER DENTAL CLICK at U+01C0 and the VERTICAL LINE at U+007C should be used, respectively.
- Similarly, the alveolar lateral click <||> should be encoded with a LATIN LETTER LATERAL CLICK at U+01C1, different from <||>, which according to the IPA is the character to be used for a major group break (by intonation), to be encoded by DOUBLE VERTICAL LINE at U+2016.
- The marking of primary stress <'> looks like an apostrophe, and is often typed with the same symbol as the ejective <◌'>. For a proper encoding, these two symbols should be typed as MODIFIER LETTER VERTICAL LINE at U+02C8 and MODIFIER LETTER APOSTROPHE at U+02BC, respectively.
- There are two different “dashed”-l characters in IPA, namely the <ł> LATIN SMALL LETTER L WITH MIDDLE TILDE at U+026B and the <ł̣> LATIN SMALL LETTER L WITH BELT at U+026C. These of course look highly similar, although they are different sounds. As a solution, we will actually propose to not use the middle tilde at all (see Section 4.5).

4.5 Pitfall: Multiple encoding options in IPA

It is not just the Unicode Standard that offers multiple options for encoding the IPA. Even the IPA specification itself offers some flexibility in how transcriptions have to be encoded. There are a few cases in which the IPA explicitly allows for different options of transcribing the same phonetic content. This is understandable from a transcriber’s point of view, but it is not acceptable when the goal is interoperability between resources written in IPA. We consider it crucial to distinguish between “valid” IPA, for which it is sufficient that any phonetically-trained reader is able to understand the transcription, and “strict” IPA, which should be standardized on a single unique encoding for each sound, so search will work across resources. We are aware of the following non-unique encoding options in the IPA, which will be discussed in turn below:

- The marking of tone
- The marking of <g>
- The marking of velarization and pharyngealization
- The placement of diacritics

The first case in which the IPA allows for different encodings is the question of how to transcribe tone. There is an old tradition to use diacritics on vowels to mark different tone levels, e.g. <èèéé>.⁶ The IPA also proposes the option of tone letters, e.g. <ᵀᵀᵀᵀ>, which are much less often used, but are more consistent for contours.⁷ Tone letters in the IPA have five different levels, and sequences of these letters can be used to indicate contours. Well-designed fonts will even merge a sequence of tone letters into a contour. For example, compare the font Linux Libertine, which does not merge tone letters <ᵀᵀᵀᵀ>, with the font CharisSIL, which merges this sequence of four tone letters into a single contour <ᵀᵀ>. For strict IPA encoding we propose to standardize on tone letters.

Second, we commonly encounter the use of <g> LATIN SMALL LETTER G at U+0067, instead of the Unicode Standard IPA character for the voiced velar stop <g> LATIN SMALL LETTER SCRIPT G at U+0261. One begins to question whether this issue is at all apparent to the working linguist, or if they simply use the U+0067 because it is easily keyboarded and thus saves time, whereas the latter must be clumsily inserted as a special symbol in most software. The International Phonetic Association has taken the stance that both the keyboard LATIN SMALL LETTER G and the LATIN SMALL LETTER SCRIPT G are valid input characters for the voiced velar plosive. Unfortunately, this decision further introduces ambiguity for linguists trying to adhere to a strict Unicode Standard IPA encoding. For strict IPA encoding we propose to standardize on the more idiosyncratic LATIN SMALL LETTER SCRIPT G at U+0261.

Third, the IPA has special markers for velarization <◌^v> and pharyngealization <◌^ɣ>. Confusingly, there is also a marker for “velarized or pharyngealized”, using the <◌̤> COMBINING TILDE OVERLAY at U+0334. The tilde overlay seems to be extremely rarely used. We suggest to try and avoid using the tilde overlay, though for reasons of backward compatibility we will allow it in valid-IPA.

⁶ To make things even more complex, there are at least two different Unicode homoglyphs for the low and high level tones, namely <◌̣> COMBINING GRAVE TONE MARK at U+0340 vs. <◌̤> COMBINING GRAVE ACCENT at U+0300 for low tone, and <◌̥> COMBINING ACUTE TONE MARK at U+0341 vs. <◌̦> COMBINING ACUTE ACCENT at U+0301 for high tone.

⁷ Not sanctioned by the IPA, but nevertheless widely attested, is the usage of superscript numbers for marking tones.

Finally, the IPA states that “diacritics may be placed above a symbol with a descender”. For example, for marking marking of voiceless pronunciation of voiced segments the IPA uses the ring diacritic. Originally, the ring should be placed below the base character, like in <ɱ̥>, using the COMBINING RING BELOW at U+0325. However, in letters with long descenders the IPA also allows to put the ring above the base, like in <ɱ̊>, using the COMBINING RING ABOVE at U+030A. Yet, proper font design does not have any problem with rendering the ring below the base character, like in <ɱ̥>, so for strict IPA encoding we propose to standardize on the ring below. As a principle, for strict IPA encoding only one options should be allowed for all diacritics.

The variable encoding as allowed by the IPA becomes even more troublesome for the tilde and diaeresis diacritics. In these cases, the IPA itself attaches different semantics to the symbols above and below a base characters. The tilde above a character (like in <ã>, using the COMBINING TILDE at U+0303) is used for nasalization, while the tilde below a character (like in <ḃ>, using the COMBINING TILDE BELOW at U+0330) indicates creaky voice. Likewise, the diaeresis above (like in <ä>, using the COMBINING DIAERESIS at U+0308) is used for centralization, while the diaeresis below a character (like in <Ḃ>, using the COMBINING DIAERESIS BELOW at U+0324) indicates breathy voice. These cases strengthen our plea to not allow diacritics to switch position for typographic convenience.

4.6 Pitfall: Tie bar

In the major revision of the IPA in 1932, affricates were represented by two consonants <tf>, ligatures <ɬ>, or with the tie-bar <ɬ̯>. In the 1938 revision the tie-bar’s semantics were broadened to indicate simultaneous articulation, as for example in labial velars such as <k̟p>. Thus, the tie-bar is a convenient diacritic for visually tokenizing input strings into chunks of phonetically salient groups, including affricates, doubly articulated consonants or diphthongs.

The tie bar can be placed above or below the base characters, e.g. <ts̯> or <ts̰>. IPA allows both options. The choice between the two symbols is purely for legible rendering; there is no difference in semantics between the two symbols. However, rendering is such a problematic issue for tie bars in general that many linguists simply do not use them. Just looking at a few different fonts already indicates that actually no font designer really gets the placement right in combination with superscripts and subscripts. If really necessary, we propose to

standardize on the tie bar above the base characters, using a COMBINING DOUBLE INVERTED BREVE at U+0361.⁸

Times new Roman: $\text{t}^{\text{h}}\text{s} \text{t}^{\text{h}}\text{s}$

CharisSIL: $\text{t}^{\text{h}}\text{s} \text{t}^{\text{h}}\text{s}$

Monaco: $\text{t}^{\text{h}}\text{s} \text{t}^{\text{h}}\text{s}$

DoulosSIL: $\text{t}^{\text{h}}\text{s} \text{t}^{\text{h}}\text{s}$

Linux Libertine: $\text{t}^{\text{h}}\text{s} \text{t}^{\text{h}}\text{s}$

Tie bars are a special type of character in the sense that they do not belong to a segment, but bind two graphemes together. This actually turns out to be rather different from Unicode conceptions. The Unicode encoding of this character belong to the Combining Diacritical Marks, namely either COMBINING DOUBLE INVERTED BREVE at U+0361 or COMBINING DOUBLE BREVE BELOW at U+035C. Such a combining mark is by definition tied to the character in front, but not the character following it. The Unicode treatment of this character thus only partly corresponds to the IPA conception, which ideally would have the tie bar linked both to the character in front and to the character following.

Further, according to the spirit of the IPA, it would also be possible to combine more than two base characters into one tie bar, but this is not possible with Unicode (i.e. there is no possibility to draw a tie bar over three or four characters). It is possible to indicate such larger groups by repeating the tie bar, like for a triphthong <æ̯œ̯> in the English word *hour*. If really necessary, we consider this possible, even though the rendering will never look good.

Most importantly though, in comparison to normal Unicode processing, the tie-bar actually takes a reversed approach to complex graphemes. Basically, the Unicode principle (see Section 2.3) is that fixed sequences in a writing system have to be specified as “tailored” grapheme clusters. Only in case the sequence is not a cluster, then this has to be explicitly indicated. IPA takes a different approach. In IPA by default different base letters are not connected into larger clusters; only when it is specified in the string itself (using the tie bar).

⁸ Also note that the UNDERTIE at U+203F looks like the tie bar below and is easily confused with it. However, it is a different character and has a different function in IPA. The undertie is used as a linking symbol to indicate the lack of a boundary, e.g. French *petit ami* [pətit_ami] ‘boyfriend’.

4.7 Pitfall: Ligatures and digraphs

One important distinction to acknowledge is the difference between multigraphs and ligatures. Multigraphs are groups of characters (in the context of IPA e.g. <tf> or <ou>) while ligatures are single characters (e.g. <ɸ> LATIN SMALL LETTER TESH DIGRAPH at U+02A7). Ligatures arose in the context of printing easier-to-read texts, and are included in the Unicode Standard for reasons of legacy encoding. However, their usage is discouraged by the Unicode core specification. Specifically related to IPA, various phonetic combinations of characters (typically affricates) are available as single code-points in the Unicode Standard, but are designated DIGRAPHS. Such glyphs might be used by software to produce a pleasing display, but they should not be hard-coded into the text itself. In the context of IPA, characters like the following ligatures should thus *not* be used. Instead a combination of two characters is preferred:

<ɖz> LATIN SMALL LETTER DZ DIGRAPH at U+02A3 (use <ɖz>)
 <ɖʒ> LATIN SMALL LETTER DEZH DIGRAPH at U+02A4 (use <ɖʒ>)
 <ɖ͡z> LATIN SMALL LETTER DZ DIGRAPH WITH CURL at U+02A5 (use <ɖ͡z>)
 <ts> LATIN SMALL LETTER TS DIGRAPH at U+02A6 (use <ts>)
 <ɸ> LATIN SMALL LETTER TESH DIGRAPH at U+02A7 (use <ɸ>)
 <ɸ͡> LATIN SMALL LETTER TC DIGRAPH WITH CURL at U+02A8 (use <ɸ͡>)
 <ɸ͡ɲ> LATIN SMALL LETTER FENG DIGRAPH at U+02A9 (use <ɸ͡ɲ>)

However, there are a few Unicode characters that are historically ligatures, but which are today considered as simple characters in the Unicode Standard and thus should be used when writing IPA, namely:

<ɰ> LATIN SMALL LETTER LEZH at U+026E
 <œ> LATIN SMALL LIGATURE OE at U+0153
 <Ɔ> LATIN LETTER SMALL CAPITAL OE at U+0276
 <æ> LATIN SMALL LETTER AE at U+00E6

4.8 Pitfall: Missing decomposition

Although many combinations of base character with diacritic are treated as canonical equivalent with precomposed characters, there are a few combinations in IPA that allow for multiple, apparently identical, encodings that are not canonical equivalent (see Section 2.9 on the principle of canonical equivalence). For that reason, the following elements should not be treated as diacritics when encoding IPA in Unicode:

- <◌̺> COMBINING PALATALIZED HOOK BELOW at U+0321
- <◌̻> COMBINING RETROFLEX HOOK BELOW at U+0322
- <◌̽> COMBINING SHORT STROKE OVERLAY at U+0335
- <◌̾> COMBINING SHORT SOLIDUS OVERLAY at U+0337

There turn out to be a lot of characters in the IPA that could be conceived as using any of these elements, like <ɲ>, <ɳ>, <ɿ> or <ø>. However, all such characters exist as well as precomposed combination in Unicode, and these precomposed characters should preferably be used. When instead combinations of a base character with diacritic are used, then these combinations are not canonical equivalent to the precomposed combinations. This means that any search will not find both at the same time.

A similar problem arises with the rhotic hook. There are two precomposed characters in Unicode with a rhotic hook, which are not canonical equivalent with a combination of the vowel with a separately encoded hook:

- <ø̥> LATIN SMALL LETTER SCHWA WITH HOOK at U+025A
- <ɜ̥> LATIN SMALL LETTER REVERSED OPEN E WITH HOOK at U+025D

All other combinations of vowels with rhotic hooks will have to be made by using <◌̥> MODIFIER LETTER RHOTIC HOOK at U+02DE, because there is no complete set of precomposed characters with rhotic hooks in Unicode. For that reason we propose to not use the two precomposed characters with hooks mentioned above, but always use the separate rhotic hook at U+02DE in IPA.

A similar situation arises with <◌̽> COMBINING TILDE OVERLAY at U+0334. The main reason some phoneticians like to use this in IPA is to mark the “dark” <ɫ> in English codas, using the character <ɫ̽> LATIN SMALL LETTER L WITH MIDDLE TILDE at U+026B. This character is not canonically equivalent to the combination <ɫ> + <◌̽>, so one of the two possible encodings has to be chosen. Because the tilde overlay is described as a general mechanism by the IPA, we propose to use the separated <◌̽> COMBINING TILDE OVERLAY at U+0334. However, note that phonetically this seems to be (almost) superfluous (see Section 4.5) and the typical usage in the form of <ɫ̽> is (almost) a homoglyph with <ɫ̥> (see Section 4.4). For these reasons we also suggest to try and avoid the tilde overlay completely.

Reversely, note that the <ç> LATIN SMALL LETTER C WITH CEDILLA at U+00E7 is canonically equivalent with <c> with <◌̸> COMBINING CEDILLA at U+0327, so it will be separated into two characters by Unicode canonical decomposition, also if such a decomposition is not intended in the IPA. However, because of the nature of canonical equivalence (see Section 2.9), these two encodings are completely

identical in any computational treatment, so this decomposition does not have any practical consequences.

4.9 Pitfall: Different notions of diacritics

Another pitfall relates to the question of what are diacritics. The problem is that the meaning of the term diacritics as used by the IPA is not the same as is used in the Unicode Standard. Specifically, diacritics in the IPA-sense are either so-called COMBINING DIACRITICAL MARKS or SPACING MODIFIER LETTERS in the Unicode Standard. Crucially, Combining Diacritical Marks are by definition combined with the character before them (to form so-called default grapheme clusters, see Section 1.4). In contrast, Spacing Modifier Letters are by definition *not* combined into grapheme clusters with the preceding character, but simply treated as separate letters. In the context of the IPA, the following IPA-diacritics are actually Spacing Modifier Letters in the Unicode Standard:

Length marks, namely:

<◌:̤> MODIFIER LETTER TRIANGULAR COLON at U+02D0

<◌:̥> MODIFIER LETTER HALF TRIANGULAR COLON at U+02D1

Tone letters, like:

<◌̥̥̥> MODIFIER LETTER EXTRA-HIGH TONE BAR at U+02E5

<◌̬̬̬> MODIFIER LETTER LOW TONE BAR at U+02E8

and others like this

Superscript letters, like:

<◌^h> MODIFIER LETTER SMALL H at U+02B0

<◌^ʔ> MODIFIER LETTER SMALL REVERSED GLOTTAL STOP at U+02E4

<◌ⁿ> SUPERSCRIPT LATIN SMALL LETTER N at U+207F

and many more like this

The rhotic hook:⁹

<◌̤̥̥̥> MODIFIER LETTER RHOTIC HOOK at U+02DE

⁹ It is really unfortunate that the rhotic hook in Unicode is classified as a Spacing Modifier, and not as a Combining Diacritical Mark. Although the rhotic hook is placed to the right of its base character (and not above or below), it still is always connected to the character in front, even physically connected to it. We cannot find any reason for this treatment, and consider it an error in Unicode. We hope it will be possible to change this classification in the future.

Although linguists might expect these characters to belong together with the character in front of them, at least for tone letters, stress symbols and <sup>h> MODIFIER LETTER SMALL H at U+02B0 the Unicode Consortium's decision to treat it as a separate character is also linguistically correct.

- According to the IPA, <sup>h> can be used both as <sup>h> for post-aspiration (following the base character) and as <sup>h> for pre-aspiration (preceding the base character), so there is no consistent direction in which this diacritic should bind. Note that there is yet another homoglyph, namely the COMBINING LATIN SMALL LETTER H at U+036A. We propose not to use this combining diacritical mark, but to standardize on Unicode modifier letters for all superscript letters in IPA.
- Tone letters are normally written at the end of the syllable, possibly occurring immediately adjacent to a consonant in the coda of the syllable. Such tone markers should of course not be treated as belonging to this consonant, so we propose to treat tone letters as separate segments.
- Stress markers <sup>ˈ> at U+02C8 and <sup>ˌ> at U+02CC have a very similar distribution in that they normally are written at the start of the stressed syllable. In a sense, they thus belong to the characters following the stress marker, but it would be wrong to cluster them together with whatever segment is at the start of the syllable. So, like tone letters, we propose to treat stress markers as separate segments.

If intended, then any default combination of Spacing Modifiers with the preceding character can be specified in orthography profiles (see Chapter 5).

4.10 Pitfall: No unique diacritic ordering

Also related to diacritics is the question of ordering. To our knowledge, the International Phonetic Association does not specify an ordering for diacritics that combine with phonetic base symbols; this exercise is left to the reasoning of the transcriber. However, such marks have to be explicitly ordered if sequences of them are to be interoperable and compatible computationally. An example is a labialized aspirated alveolar plosive: $\text{<t}^{\text{wh}}\text{>}$. There is nothing holding linguists back from using $\text{<t}^{\text{hw}}\text{>}$ instead (with exactly the same intended meaning). However, from a technical standpoint, these two sequences are different, e.g. if both sequences are used in a document, searching for $\text{<t}^{\text{wh}}\text{>}$ will not find any instances of $\text{<t}^{\text{hw}}\text{>}$, and vice versa. Likewise, a creaky voiced syllabic dental nasal can be encoded in various orders, e.g. $\text{<n}^{\text{̤}}\text{>}$, $\text{<n}^{\text{̥}}\text{>}$ or $\text{<n}^{\text{̚}}\text{>}$.

Canonical combining classes

In accordance with the absence of any specification of ordering in the IPA, the Unicode Standard likewise does not propose any standardized orders. Both leave it to the user to be consistent; this approach naturally invites inconsistency across different authored resources.

There is one (minor) aspect of ordering for which the Unicode Standard does present a canonical solution. Fortunately, this is uncontroversial from a linguistic perspective. Diacritics in the Unicode Standard (i.e. Combining Diacritical Marks, see Section 4.9) are classified in so-called CANONICAL COMBINING CLASSES. In practice, the diacritics are distinguished by their position relative to the base character.¹⁰ When applying a Unicode normalization (NFC or NFD, see Section 2.9), the diacritics in different positions are put in a specified order. This process therefore harmonizes the difference between different encodings in some situations, for example in the case of an extra-short creaky voice vowel <ë>. This grapheme cluster can be encoded either as <e>+<◌◌̥>+<◌◌̚> or as <e>+<◌◌̚>+<◌◌̥>. To prevent this twofold encoding, the Unicode Standard specifies the second ordering as canonical (namely, diacritics below are put before diacritics above).

When encoding a string according to the Unicode Standard, it is possible to do this either using the NFC (composition) or NFD (decomposition) normalization (see Section 2.9). Decomposition implies that precomposed characters (like <á> LATIN SMALL LETTER A WITH ACUTE at U+00E1) will be split into its parts. This might sound preferable for a linguistic analysis, as the different diacritics are separated from the base characters. However, note that most attached elements like strokes (e.g. in the <i>), retroflex hooks (e.g. in <z̥>) or rhotic hooks (e.g. in <ʒ̥>) will not be decomposed, but strangely enough a cedilla (like in <ç>) will be decomposed (see Section 4.8). In general, Unicode decomposition does not behave like a feature decomposition as expected from a linguistic perspective. It is thus important to consider Unicode decomposition only as a technical procedure, and not assume that it is linguistically sensible.

Proposal for diacritic ordering

Facing the problem of specifying a consistent ordering of diacritics while developing a large database of phonological inventories from the world's languages, Moran (2012: 540) defined a set of diacritic ordering conventions. The conventions are influenced by the linguistic literature, though some ad-hoc decisions

¹⁰ See http://unicode.org/reports/tr44/#Canonical_Combining_Class_Values for a detailed description.

had to be taken given the vast variability of phonological segments described by linguists.

By Unicode Canonical Combining Classes, the diacritics on top of a character, like <◌̂> (Combining Class number 1), always come before diacritics below (Combining Class number 220), which in turn always come before diacritics above (Combining Class number 230), which in turn come before diacritics over multiple characters like the tie bar <◌̳> (Combining Class number 233). We follow this order, but add the other IPA diacritics (which are not diacritics in the Unicode sense) between diacritics below and the tie bar. Further, *within* all these classes of diacritics there is no canonical ordering specified by Unicode, so we propose an explicit ordering here.

Starting with the diacritics below: if a character sequence contains more than one diacritic below the base character, then the place features are applied first (linguolabial, dental, apical, laminal, advanced, retracted), followed by the manner features (raised, lowered, advanced and retracted tongue root), then secondary articulations (more round, less round), laryngeal settings (creaky, breathy, voiced, devoiced), and finally the syllabic or non-syllabic marker. So, the order that is proposed is the following, where <|> indicates *or* and <→> indicates *precedes*. Note that the groups of alternatives (as marked by <|>) are supposed never to occur together with the same base character. In effect, this represents yet another restriction on possible diacritic sequences.

COMBINING DIACRITICAL MARKS (BELOW) ORDERING:

- linguolabial <◌̦> | dental <◌̤> | apical <◌̥> | laminal <◌̪>
- advanced <◌̦̰> | retracted <◌̦̱>
- raised <◌̥̥̥> | lowered <◌̥̥̥̥>
- advanced tongue root <◌̠> | retracted tongue root <◌̡>
- more rounded <◌̙> | less rounded <◌̘>
- creaky voiced <◌̰̰̰> | breathy voiced <◌̰̰̰̰> | voiced <◌̰̰̰̰̰> | voiceless <◌̰̰̰̰̰̰>
- syllabic <◌̰̰̰̰̰̰̰> | non-syllabic <◌̰̰̰̰̰̰̰̰>

Next, if a character sequence contains more than one diacritic above the base character, we propose the following order:

COMBINING DIACRITICAL MARKS (ABOVE) ORDERING:

- nasalized <◌̃>
- centralized <◌̞̞̞> | mid-centralized <◌̞̞̞̞>
- extra short <◌̞̞̞̞̞>

→ no audible release <◌̚>

Then, when a character sequence contains more than one character of the Spacing Modifier Letters, these will be placed after all combining diacritic marks in the following order:

SPACING MODIFIER LETTERS ORDERING:

- rhotic hook <◌̺>
- lateral release <◌̜> | nasal release <◌̝>
- labialized <◌̥>
- palatalized <◌̦>
- velarized <◌̧>
- pharyngealized <◌̨>
- aspirated <◌̩> | ejective <◌̪>
- long <◌̫> | half-long <◌̬>

Finally, the tie bar follows at the very end of any such sequence:

TIE BAR:

- tie bar <◌̯>

4.11 Pitfall: Revisions to the IPA

With each revision of the IPA, many decisions need to be made by the Association as to which symbols should be added, removed or changed. For example, in the 1989 revision of the IPA at the Kiel Convention, changes to specific symbols (in previous charts) were debated and the Association's members certain decisions. The prevailing mood at the convention was not to change specific symbols unless a strong case was made (Ladefoged 1990). For example, two such decisions included:

- symbols for clicks were changed from <ɀ Ɂ ɂ> to <ǀ ǁ ǃ> because the latter were the symbols used by nearly all Khoisanists and Bantuists.
- the Americanist tradition of using <ǃ>, a COMBINING CARON at U+030C for all post-alveolar sounds, like in <š ž č j>, was not adopted because the Association members at the convention “were not sufficiently impressed by arguments ... to the effect that these sounds formed a natural class, and thus it would be appropriate to recognize this by maintaining a common aspect to their symbolism” (Ladefoged 1990: 62).

These decisions have practical consequences for transcribers of IPA, particularly those who wish to follow recommended practices of encoding electronic text in the Unicode Standard. For example, the Unicode Standard contains removed symbols as labels them as clicks, e.g. <ɿ> LATIN SMALL LETTER TURNED T at U+0287 in the IPA EXTENSIONS block has the comment ‘dental click (sound of “tsk tsk”)’. In this case, the IPA transcriber must know the status of click symbols in the current version of the IPA and then identify those characters within the Unicode Standard.

The most controversial issue regarding symbols debated at the convention was the representation for voiceless implosives (Ladefoged 1990: 62). In accordance with the principles of the IPA, as outlined in Section 3, distinct symbols are favoured for cases of phonological contrast. Further, also convenience of display in the chart must be taken into account when arguing for or against the inclusion or deletion of IPA symbols in the IPA chart. Finally, the inclusion or deletion of symbols should consider the current state of phonetic knowledge of the world’s languages.

Ladefoged (1990) argued against the inclusion of the symbols <ɓ, ɗ, ɠ, ɡ> for voiceless implosives, noting (i) that they are not contrastive (e.g. in Mayan languages); (ii) that there is no instrumental evidence supporting voiceless implosives in Africa; and (iii) that the sounds are sufficiently rare as to not need a whole new row of symbols in the chart. Ladefoged favored symbolizing the sounds using a voiceless diacritic ring below voiced implosives, e.g. <ɓ̥>. Nevertheless, in the 1989 IPA chart there is indeed a row for implosives containing voiceless and voiced pairs.¹¹ But already in the next revision, in 1993 (with an update in 1996), the voiceless implosives were dropped. The implosives row from the IPA consonantal chart disappeared and voiced implosives were given a column in the non-pulmonic consonants table (which is still reflected in the latest revision to date, IPA 2005).

The *Journal of the International Phonetic Association* follows its own published standard for the IPA at the time of publication, even when it may conflict with the Association’s principle of using different symbols for contrastive sounds and diacritics for phonetic variation. For example, in the case of voiceless implosives, Mc Laughlin (2005) shows that Seereer-Siin (Niger-Congo, Atlantic; ISO 639-3 srr) has a phonologically contrastive set of voiced and voiceless implosive stops at the labial, coronal and palatal places of articulation. These symbols are transcribed in an *Illustrations of the IPA* article in the IPA journal as <ɓ̥, ɗ̥, ɠ̥>.

¹¹ https://en.wikipedia.org/wiki/File:IPA_as_of_1989.png

The point of this pitfall is to highlight that revisions to the IPA will continue into the future, albeit they are infrequent. Nevertheless, given the Unicode Standard's principle of maintaining backwards compatibility (at all costs), transcribers and consumers of IPA cannot rely solely on remarks in the Unicode Standard to reflect current standard IPA usage. There is the possibility that at a later revision of the IPA, symbols that are not currently encoded in the Unicode Standard are added to the IPA – although we think this is unlikely.

4.12 Additions to the IPA

In the course of collecting a large sample of phoneme systems across the world's languages (Moran 2012), we found that in order to preserve distinctions both within and across language descriptions, additions to the approved IPA glyph set were needed. Wherever possible these additions were drawn from the extIPA symbols for disordered speech.¹² This section describes our proposed additions to the IPA glyph set. These additions are not part of the official IPA recommendations, so they should be used with care.

- RETROFLEX CLICK

Retroflex clicks can be represented by <!!> DOUBLE EXCLAMATION MARK at U+203C. Note that the (post-)alveolar click <!> at U+01C3 is confusingly referred to as LATIN LETTER RETROFLEX CLICK in the Unicode standard, which is probably best seen as an error.

- VOICED RETROFLEX IMPLOSIVE

Although the IPA includes a series of voiced implosives (marked with a hook on top, see Section 4.8), there is no voiced retroflex implosive. Following the spirit of the IPA, we propose to use <ɖ̥> LATIN SMALL LETTER D WITH HOOK AND TAIL at U+1D91 for this sound.

- FORTIS/LENIS

Languages described as having a fortis/plain/lenis distinction that corresponds poorly with the traditional voiced/voiceless-unaspirated/voiceless-aspirated continuum can be marked using the voiceless glyph for the plain phoneme, and then <◌͡> COMBINING DOUBLE VERTICAL LINE BELOW at U+0348 to mark the fortis articulation, and/or <◌͜> COMBINING LEFT ANGLE BELOW at U+0349 for the lenis articulation.

¹² **Add better reference here:** <https://www.internationalphoneticassociation.org/sites/default/files/extIPAChart2008.pdf>

- **FRICTIONALIZATION**

The diacritic <◌̤> COMBINING X BELOW at U+0353 can be used to represent three types of frictionalized sounds: First, click consonants where the release of the anterior closure involves an ingressive “sucking” sound similar to a fricative, for example <k^{!h}_{̤̤}>.

- **DERHOTICIZATION**

For derhoticization we propose to use <◌̣> COMBINING BREVE BELOW at U+032E. **WHAT EXACTLY IS THIS?**

- **CORONAL NON-SIBILANT**

Languages described as having a sibilant/non-sibilant distinction among coronal fricatives and affricates can be handled using the subscript <◌̱> COMBINING EQUALS SIGN BELOW at U+0347 to mark the non-sibilant phoneme.

- **GLOTTALIZATION**

Glottalized sounds can be indicated using <◌̚> MODIFIER LETTER GLOTTAL STOP at U+02C0, unless it is clear that either “ejective” or “creaky voicing” are the intended sounds (in which cases the standard IPA diacritics should be used). Pre-glottalized sounds can be marked with <^ʔ◌̚> to the left of the base glyph, for example <^ʔt̚>.

- **VOICED PRE-ASPIRATION**

Voiced sounds having pre-aspiration can be marked with <^h◌̚> MODIFIER LETTER SMALL H WITH HOOK at U+02B1 to the left of the base glyph, for example <^hd̚>.

- **EPILARYNGEAL PHONATION**

There are some rare articulations that make use of an epilaryngeal phonation mechanism (e.g., the “sphincteric vowels” of !Xóḏ). To represent these vowels, we propose to use the modifier <◌̥> MODIFIER LETTER CAPITAL E at U+1D31 to denote such sphincteric phonation.

4.13 Recommendations

Summarizing the pitfalls as discussed in this chapter, we propose to define three different IPA encodings: strict-IPA, valid-IPA and widened-IPA. Informally speaking, valid-IPA represents the current state of the IPA (International Phonetic As-

sociation 2005). Strict-IPA represents a more constrained version of IPA, while widened-IPA is a slightly extended version of IPA, allowing a few more symbols.

Strict-IPA encoding is supposed to be used when interoperability of phonetic resources is intended. It is a strongly constrained subset of IPA geared towards uniqueness of encoding. Ideally, for each transcription there should be exactly one possible strict-IPA encoding. For each phonetic feature there is only one possibility (see Section 4.5) and the IPA diacritics are forced into a canonical ordering (see Section 4.10).

Valid-IPA does allow alternative symbols with the same phonetic meaning, as specified in the official IPA specifications. Also, valid-IPA does not enforce a specific ordering of diacritics, because the IPA does not propose any such ordering. This means that in valid-IPA the same phonetic intention can be encoded in multiple ways. This is sufficient for phonetically trained human eyes, but it is not sufficient for automatic interoperability.

Finally, widened-IPA includes a few more symbols which seem to be useful for various special cases (see Section 4.12).

At the end of this chapter we have added a few longish tables summarizing all 159 different unicode codepoints that form the basis of strict-IPA encoding (107 letters, 36 diacritics and 16 remaining symbols). Each of these tables shows a typical glyph, and then lists the Unicode Codepoint, Unicode Name and IPA description for each symbol. Further, there is a table with the additional options for valid-IPA and a table with the additional options for widened-IPA.

- **STRICT-IPA LETTERS**

The 107 different IPA letters as allowed in strict-IPA encoding are listed in Table 4.1 starting on page 70.

- **STRICT-IPA DIACRITICS**

The 36 different IPA diacritics and tone markers (both Unicode Modifier Letters and Combining Diacritical Marks) as allowed in strict-IPA encoding are listed in Table 4.2 starting on page 73.

- **STRICT-IPA REMAINDERS**

The 16 remaining IPA symbols (boundary, stress, tone letters and intonation markers) as allowed in strict-IPA encoding are listed in Table 4.3 on page 74.

- **VALID-IPA ADDITIONS**

The 16 additional symbols as allowed in valid-IPA encoding are listed in Table 4.4 on page 75.

- **WIDENED-IPA ADDITIONS**

The 10 proposed additions to the IPA are listed in Table 4.5 on page 76.

Table 4.1: Strict-IPA letters with Unicode encodings

	Code	Unicode name	IPA name
a	U+0061	LATIN SMALL LETTER A	open front unrounded
æ	U+00E6	LATIN SMALL LETTER AE	raised open front unrounded
ɐ	U+0250	LATIN SMALL LETTER TURNED A	lowered schwa
ɑ	U+0251	LATIN SMALL LETTER ALPHA	open back unrounded
ɒ	U+0252	LATIN SMALL LETTER TURNED ALPHA	open back rounded
b	U+0062	LATIN SMALL LETTER B	voiced bilabial plosive
B	U+0299	LATIN LETTER SMALL CAPITAL B	voiced bilabial trill
ɸ	U+0253	LATIN SMALL LETTER B WITH HOOK	voiced bilabial implosive
c	U+0063	LATIN SMALL LETTER C	voiceless palatal plosive
ç	U+00E7	LATIN SMALL LETTER C WITH CEDILLA	voiceless palatal fricative
ɕ	U+0255	LATIN SMALL LETTER C WITH CURL	voiceless alveolo-palatal fricative
d	U+0064	LATIN SMALL LETTER D	voiced alveolar plosive
ð	U+00F0	LATIN SMALL LETTER ETH	voiced dental fricative
ɖ	U+0256	LATIN SMALL LETTER D WITH TAIL	voiced retroflex plosive
ɗ	U+0257	LATIN SMALL LETTER D WITH HOOK	voiced dental/alveolar implosive
e	U+0065	LATIN SMALL LETTER E	close-mid front unrounded
ə	U+0259	LATIN SMALL LETTER SCHWA	mid-central schwa
ɛ	U+025B	LATIN SMALL LETTER OPEN E	open-mid front unrounded
ɜ	U+0258	LATIN SMALL LETTER REVERSED E	close-mid central unrounded
ɞ	U+025C	LATIN SMALL LETTER REVERSED OPEN E	open-mid central unrounded
ɐ̃	U+025E	LATIN SMALL LETTER CLOSED REVERSED OPEN E	open-mid central rounded
f	U+0066	LATIN SMALL LETTER F	voiceless labiodental fricative
g	U+0261	LATIN SMALL LETTER SCRIPT G	voiced velar plosive
ɢ	U+0262	LATIN LETTER SMALL CAPITAL G	voiced uvular plosive
ɣ	U+0260	LATIN SMALL LETTER G WITH HOOK	voiced velar implosive
ʒ	U+029B	LATIN LETTER SMALL CAPITAL G WITH HOOK	voiced uvular implosive
ɣ̥	U+0264	LATIN SMALL LETTER RAMS HORN	close-mid back unrounded
ɣ̥	U+0263	LATIN SMALL LETTER GAMMA	voiced velar fricative
h	U+0068	LATIN SMALL LETTER H	voiceless glottal fricative
ħ	U+0127	LATIN SMALL LETTER H WITH STROKE	voiceless pharyngeal fricative
H	U+029C	LATIN LETTER SMALL CAPITAL H	voiceless epiglottal fricative

continued on next page

Table 4.1 Strict-IPA letters with Unicode encodings — *continued*

	Code	Unicode name	IPA name
ɦ	U+0266	LATIN SMALL LETTER H WITH HOOK	voiced glottal fricative
ɸ	U+0267	LATIN SMALL LETTER HENG WITH HOOK	simultaneous voiceless postalveolar+velar fricative
ɥ	U+0265	LATIN SMALL LETTER TURNED H	voiced labial-palatal approximant
ɪ	U+0069	LATIN SMALL LETTER I	close front unrounded
ɪ	U+026A	LATIN LETTER SMALL CAPITAL I	lax close front unrounded
ɨ	U+0268	LATIN SMALL LETTER I WITH STROKE	close central unrounded
ɹ	U+006A	LATIN SMALL LETTER J	voiced palatal approximant
ɹ	U+029D	LATIN SMALL LETTER J WITH CROSSED TAIL	voiced palatal fricative
ɹ̥	U+025F	LATIN SMALL LETTER DOTLESS J WITH STROKE	voiced palatal plosive
ɹ̥	U+0284	LATIN SMALL LETTER DOTLESS J WITH STROKE AND HOOK	voiced palatal implosive
ɰ	U+006B	LATIN SMALL LETTER K	voiceless velar plosive
ɭ	U+006C	LATIN SMALL LETTER L	voiced alveolar lateral approximant
ɮ	U+029F	LATIN LETTER SMALL CAPITAL L	voiced velar lateral approximant
ɮ̥	U+026C	LATIN SMALL LETTER L WITH BELT	voiceless alveolar lateral fricative
ɮ̥	U+026D	LATIN SMALL LETTER L WITH RETROFLEX HOOK	voiced retroflex lateral approximant
ɮ̥	U+026E	LATIN SMALL LETTER LEZH	voiced alveolar lateral fricative
ɮ̥	U+028E	LATIN SMALL LETTER TURNED Y	voiced palatal lateral approximant
ɱ	U+006D	LATIN SMALL LETTER M	voiced bilabial nasal
ɱ̥	U+0271	LATIN SMALL LETTER M WITH HOOK	voiced labiodental nasal
ɳ	U+006E	LATIN SMALL LETTER N	voiced alveolar nasal
ɳ	U+0274	LATIN LETTER SMALL CAPITAL N	voiced uvular nasal
ɳ̥	U+0272	LATIN SMALL LETTER N WITH LEFT HOOK	voiced palatal nasal
ɳ̥	U+0273	LATIN SMALL LETTER N WITH RETROFLEX HOOK	voiced retroflex nasal
ɳ̥	U+014B	LATIN SMALL LETTER ENG	voiced velar nasal
ɔ	U+006F	LATIN SMALL LETTER O	close-mid back rounded
ø	U+00F8	LATIN SMALL LETTER O WITH STROKE	close-mid front rounded
œ	U+0153	LATIN SMALL LIGATURE OE	open-mid front rounded
œ	U+0276	LATIN LETTER SMALL CAPITAL OE	open front rounded
ɔ̥	U+0254	LATIN SMALL LETTER OPEN O	open-mid back rounded
ɵ	U+0275	LATIN SMALL LETTER BARRED O	close-mid central rounded

continued on next page

Table 4.1 Strict-IPA letters with Unicode encodings — *continued*

	Code	Unicode name	IPA name
p	U+0070	LATIN SMALL LETTER P	voiceless bilabial plosive
ɸ	U+0278	LATIN SMALL LETTER PHI	voiceless bilabial fricative
q	U+0071	LATIN SMALL LETTER Q	voiceless uvular plosive
r	U+0072	LATIN SMALL LETTER R	voiced alveolar trill
ʀ	U+0280	LATIN LETTER SMALL CAPITAL R	voiced uvular trill
ɹ	U+0279	LATIN SMALL LETTER TURNED R	voiced alveolar approximant
ɽ	U+027A	LATIN SMALL LETTER TURNED R WITH LONG LEG	voiced alveolar lateral flap
ɻ	U+027B	LATIN SMALL LETTER TURNED R WITH HOOK	voiced retroflex approximant
ɽ	U+027D	LATIN SMALL LETTER R WITH TAIL	voiced retroflex tap
ɾ	U+027E	LATIN SMALL LETTER R WITH FISHHOOK	voiced alveolar tap
ʁ	U+0281	LATIN LETTER SMALL CAPITAL INVERTED R	voiced uvular fricative
s	U+0073	LATIN SMALL LETTER S	voiceless alveolar fricative
ʂ	U+0282	LATIN SMALL LETTER S WITH HOOK	voiceless retroflex fricative
ʃ	U+0283	LATIN SMALL LETTER ESH	voiceless postalveolar fricative
t	U+0074	LATIN SMALL LETTER T	voiceless alveolar plosive
ʈ	U+0288	LATIN SMALL LETTER T WITH RETROFLEX HOOK	voiceless retroflex plosive
u	U+0075	LATIN SMALL LETTER U	close back rounded
ʊ	U+0289	LATIN SMALL LETTER U BAR	close central rounded
ʉ	U+026F	LATIN SMALL LETTER TURNED M	close back unrounded
ɯ	U+0270	LATIN SMALL LETTER TURNED M WITH LONG LEG	voiced velar approximant
ʊ	U+028A	LATIN SMALL LETTER UPSILON	lax close back rounded
v	U+0076	LATIN SMALL LETTER V	voiced labiodental fricative
ʋ	U+028B	LATIN SMALL LETTER V WITH HOOK	voiced labiodental approximant
ʌ	U+2C71	LATIN SMALL LETTER V WITH RIGHT HOOK	voiced labiodental tap
ʌ	U+028C	LATIN SMALL LETTER TURNED V	open-mid back unrounded
w	U+0077	LATIN SMALL LETTER W	voiced labial-velar approximant
ʍ	U+028D	LATIN SMALL LETTER TURNED W	voiceless labial-velar fricative
x	U+0078	LATIN SMALL LETTER X	voiceless velar fricative
y	U+0079	LATIN SMALL LETTER Y	close front rounded
ʏ	U+028F	LATIN LETTER SMALL CAPITAL Y	lax close front rounded
z	U+007A	LATIN SMALL LETTER Z	voiced alveolar fricative
ʐ	U+0290	LATIN SMALL LETTER Z WITH RETROFLEX HOOK	voiced retroflex fricative

continued on next page

Table 4.1 Strict-IPA letters with Unicode encodings — *continued*

	Code	Unicode name	IPA name
ɹ	U+0291	LATIN SMALL LETTER Z WITH CURL	voiced alveolo-palatal fricative
ʒ	U+0292	LATIN SMALL LETTER EZH	voiced postalveolar fricative
ʔ	U+0294	LATIN LETTER GLOTTAL STOP	voiceless glottal plosive
ʕ	U+0295	LATIN LETTER PHARYNGEAL VOICED FRICATIVE	voiced pharyngeal fricative
ʔ̥	U+02A1	LATIN LETTER GLOTTAL STOP WITH STROKE	epiglottal plosive
ʕ̥	U+02A2	LATIN LETTER REVERSED GLOTTAL STOP WITH STROKE	voiced epiglottal fricative
ǀ	U+01C0	LATIN LETTER DENTAL CLICK	voiceless dental click
ǁ	U+01C1	LATIN LETTER LATERAL CLICK	voiceless alveolar lateral click
ǃ	U+01C2	LATIN LETTER ALVEOLAR CLICK	voiceless palatoalveolar click
ǂ	U+01C3	LATIN LETTER RETROFLEX CLICK	voiceless (post)alveolar click
ǃ̥	U+0298	LATIN LETTER BILABIAL CLICK	voiceless bilabial click
β	U+03B2	GREEK SMALL LETTER BETA	voiced bilabial fricative
θ	U+03B8	GREEK SMALL LETTER THETA	voiceless dental fricative
χ	U+03C7	GREEK SMALL LETTER CHI	voiceless uvular fricative

Table 4.2: Strict-IPA diacritics with Unicode encodings

	Code	Unicode name	IPA name
̠	U+0334	COMBINING TILDE OVERLAY	velarized or pharyngealized
̡	U+033C	COMBINING SEAGULL BELOW	linguolabial
̢	U+032A	COMBINING BRIDGE BELOW	dental
̣	U+033B	COMBINING SQUARE BELOW	laminal
̤	U+033A	COMBINING INVERTED BRIDGE BELOW	apical
̥	U+031F	COMBINING PLUS SIGN BELOW	advanced
̦	U+0320	COMBINING MINUS SIGN BELOW	retracted
̧	U+031D	COMBINING UP TACK BELOW	raised
̨	U+031E	COMBINING DOWN TACK BELOW	lowered
̩	U+0318	COMBINING LEFT TACK BELOW	advanced tongue root
̪	U+0319	COMBINING RIGHT TACK BELOW	retracted tongue root

continued on next page

Table 4.2 Strict-IPA diacritics with Unicode encodings — *continued*

	Code	Unicode name	IPA name
◌̥	U+031C	COMBINING LEFT HALF RING BELOW	less rounded
◌̦	U+0339	COMBINING RIGHT HALF RING BELOW	more rounded
◌̧	U+032C	COMBINING CARON BELOW	voiced
◌̨	U+0325	COMBINING RING BELOW	voiceless
◌̩	U+0330	COMBINING TILDE BELOW	creaky voiced
◌̪	U+0324	COMBINING DIAERESIS BELOW	breathy voiced
◌̫	U+0329	COMBINING VERTICAL LINE BELOW	syllabic
◌̬	U+032F	COMBINING INVERTED BREVE BELOW	non-syllabic
◌̭	U+0303	COMBINING TILDE	nasalized
◌̮	U+0308	COMBINING DIAERESIS	centralized
◌̯	U+033D	COMBINING X ABOVE	mid-centralized
◌̰	U+0306	COMBINING BREVE	extra-short
◌̱	U+031A	COMBINING LEFT ANGLE ABOVE	no audible release
◌̲	U+02DE	MODIFIER LETTER RHOTIC HOOK	rhotacized
◌̳	U+02E1	MODIFIER LETTER SMALL L	lateral release
◌̴	U+020F	SUPERSCRIPIT LATIN SMALL LETTER N	nasal release
◌̵	U+02B7	MODIFIER LETTER SMALL W	labialized
◌̶	U+02B2	MODIFIER LETTER SMALL J	palatalized
◌̷	U+02E0	MODIFIER LETTER SMALL GAMMA	velarized
◌̸	U+02E4	MODIFIER LETTER SMALL REVERSED GLOTTAL STOP	pharyngealized
◌̹	U+02B0	MODIFIER LETTER SMALL H	aspirated
◌̺	U+02BC	MODIFIER LETTER APOSTROPHE	ejective
◌̻	U+02D0	MODIFIER LETTER TRIANGULAR COLON	long
◌̼	U+02D1	MODIFIER LETTER HALF TRIANGULAR COLON	half-long
◌̽	U+0361	COMBINING DOUBLE INVERTED BREVE	tie bar

Table 4.3: Other Strict-IPA symbols with Unicode encodings

Code	Unicode name	IPA name
------	--------------	----------

continued on next page

Table 4.3 Other Strict-IPA symbols with Unicode encodings – *continued*

	Code	Unicode name	IPA name
'	U+02C8	MODIFIER LETTER VERTICAL LINE	primary stress
,	U+02CC	MODIFIER LETTER LOW VERTICAL LINE	secondary stress
Ɑ	U+02E5	MODIFIER LETTER EXTRA-HIGH TONE BAR	extra high tone
Ɱ	U+02E6	MODIFIER LETTER HIGH TONE BAR	high tone
Ɐ	U+02E7	MODIFIER LETTER MID TONE BAR	mid tone
Ɒ	U+02E8	MODIFIER LETTER LOW TONE BAR	low tone
ⱱ	U+02E9	MODIFIER LETTER EXTRA-LOW TONE BAR	extra low tone
↗	U+2191	UPWARDS ARROW	global rise
↘	U+2193	DOWNWARDS ARROW	global fall
↗↘	U+2197	NORTH EAST ARROW	global rise
↘↗	U+2198	SOUTH EAST ARROW	global fall
	U+0020	SPACE	word break
.	U+002E	FULL STOP	syllable break
	U+007C	VERTICAL LINE	minor group break (foot)
	U+2016	DOUBLE VERTICAL LINE	major group break (intonation)
˘	U+203F	UNDERTIE	linking (absence of a break)

Table 4.4: Additional characters for valid-IPA with Unicode encodings

	Code	Unicode name	Phonetic description
◌̥	U+030A	COMBINING RING ABOVE	voiceless (above)
◌̑	U+0067	LATIN SMALL LETTER G	voiced velar plosive
◌̚	U+030B	COMBINING DOUBLE ACUTE ACCENT	extra high tone
◌̊	U+0301	COMBINING ACUTE ACCENT	high tone
◌̋	U+0304	COMBINING MACRON	mid tone
◌̌	U+0300	COMBINING GRAVE ACCENT	low tone
◌̍	U+030F	COMBINING DOUBLE GRAVE ACCENT	extra low tone
◌̎	U+0302	COMBINING CIRCUMFLEX ACCENT	falling
◌̏	U+030C	COMBINING CARON	rising
◌̐	U+1DC4	COMBINING MACRON-ACUTE	high rising

continued on next page

Table 4.4 Additional characters for valid-IPA with Unicode encodings — *continued*

	Code	Unicode name	Phonetic description
◌̎	U+1DC5	COMBINING GRAVE-MACRON	low rising
◌̍	U+1DC6	COMBINING MACRON-GRAVE	low falling
◌̏	U+1DC7	COMBINING ACUTE-MACRON	high falling
◌̐	U+1DC8	COMBINING GRAVE-ACUTE-GRAVE	rising-falling
◌̑	U+1DC9	COMBINING ACUTE-GRAVE-ACUTE	falling-rising
◌̽	U+035C	COMBINING DOUBLE BREVE BELOW	tie bar (below)

Table 4.5: Additions to widened-IPA with Unicode encodings

	Code	Unicode name	Phonetic description
!!	U+203C	DOUBLE EXCLAMATION MARK	retroflex click
ɖ	U+1D91	LATIN SMALL LETTER D WITH HOOK AND TAIL	voiced retroflex implosive
◌̹	U+0348	COMBINING DOUBLE VERTICAL LINE BELOW	fortis
◌̸	U+0349	COMBINING LEFT ANGLE BELOW	lenis
◌̥	U+0353	COMBINING X BELOW	frictionalized
◌̦	U+032E	COMBINING BREVE BELOW	derhoticized
◌̨	U+0347	COMBINING EQUALS SIGN BELOW	non-sibilant
◌̚	U+02C0	MODIFIER LETTER GLOTTAL STOP	glottalized
ɦ̥	U+02B1	MODIFIER LETTER SMALL H WITH HOOK	voiced pre-aspirated
◌̜	U+1D31	MODIFIER LETTER CAPITAL E	epilaryngeal phonation

5 Orthography profiles

5.1 Characterizing writing systems

At this point in the course of rapid ongoing developments, we are left with a situation in which the Unicode Standard offers a highly detailed and flexible approach to deal computationally with writing systems, but it has unfortunately not influenced the linguistic practice very much. In many practical situations, the Unicode Standard is far too complex for the day-to-day practice in linguistics because it does not offer practical solutions for the down-to-earth problems of many linguists. In this section, we propose some simple practical guidelines and methods to improve on this situation.

Our central aims for linguistics, to be approached with a Unicode-based solution, are: (i) to improve the consistency of the encoding of sources, (ii) to transparently document knowledge about the writing system (including transliteration), and (iii) to do all of that in a way that is easy and quick to manage for many different sources with many different writing systems. The central concept in our proposal is the **ORTHOGRAPHY PROFILE**, a simple tab-separated text file, that characterizes and documents a writing system. We also offer basic implementations in Python and R to assist with the production of such files, and to apply orthography profiles for consistency testing, grapheme tokenization and transliteration. Not only can orthography profiles be helpful in the daily practice of linguistics, they also succinctly document the orthographic details of a specific source, and, as such, might fruitfully be published alongside sources (e.g. in digital archives). Also, in high-level linguistic analyzes in which the graphemic detail is of central importance (e.g. phonotactic or comparative-historical studies), orthography profiles can transparently document the decisions that have been taken in the interpretation of the orthography in the sources used.

Given these goals, Unicode locale descriptions (see Section 1.4) might seem like the ideal orthography profiles. However, there are various practical obstacles preventing the use of such locale descriptions in the daily linguistic practice, namely: (i) the XML-structure is too verbose to easily and quickly produce or correct manually, (ii) locale descriptions are designed for a wide scope on infor-

mation (like date formats or names of weekdays) most of which is not applicable for documenting writing systems, and (iii) most crucially, even if someone made the effort to produce a technically correct locale description for a specific source at hand, then it is nigh impossible to deploy the description. This is because a locale description has to be submitted to and accepted by the Unicode Common Locale Data Repository. The repository is (rightly so) not interested in descriptions that only apply to a limited set of sources (e.g. descriptions for only a single dictionary).

The major challenge then is developing an infrastructure to identify the elements that are individual graphemes in a source, specifically for the enormous variety of sources using some kind of alphabetic writing system. Authors of source documents (e.g. dictionaries, wordlists, corpora) use a variety of writing systems that range from their own idiosyncratic transcriptions to already well-established practical or longstanding orthographies. Although the IPA is one practical choice as a sound-based normalization for writing systems (which can act as an interlingual pivot to attain interoperability across writing systems), graphemes in each writing system must also be identified and standardized if interoperability across different sources is to be achieved. In most cases, this amounts to more than simply mapping a grapheme to an IPA segment because graphemes must first be identified in context (e.g. is the sequence one sound or two sounds or both?) and strings must be tokenized, which may include taking orthographic rules into account (e.g. between vowels is /n/ and after a vowel but before a consonant is a nasalized vowel / \tilde{V} /).

In our experience, data from each source must be individually tokenized into graphemes so that its orthographic structure can be identified and its contents can be extracted. To extract data for analysis, a source-by-source approach is required before an orthography profile can be created. For example, almost each available lexicon on the world's languages is idiosyncratic in its orthography and thus requires lexicon-specific approaches to identify graphemes in the writing system and to map graphemes to phonemes, if desired.

Our key proposal for the characterization of a writing system is to use a grapheme tokenization as an inter-orthographic pivot. Basically, any source document is tokenized by graphemes, and only then a mapping to IPA (or any other orthographic transliteration) is performed. An ORTHOGRAPHY PROFILE then is a description of the units and rules that are needed to adequately model a graphemic tokenization for a language variety as described in a particular source document. An orthography profile summarizes the Unicode (tailored) graphemes and ortho-

graphic rules used to write a language (the details of the structure and assumptions of such a profile will be presented in the next section).

As an example of graphemic tokenization, note the three different levels of technological and linguistic elements that interact in the hypothetical lexical form <ts^hǫshi>:

1. code points (10 text elements): t s ^h o ǫ ɔ s h i
2. grapheme clusters (7 text elements): t s ^h ǫ s h i
3. tailored grapheme clusters (4 text elements): ts^h ǫ sh i

In (1), the string <ts^hǫshi> has been tokenized into ten Unicode code points (using NFD normalization), delimited here by space. Unicode tokenization is required because sequences of code points can differ in their visual and logical orders. For example, <ǫ> is ambiguous to whether it is the sequence of <o> + <̥> + <̥> or <o> + <̥> + <̥>. Although these two variants are visually homographs, computationally they are different. Unicode normalization should be applied to this string to reorder the code points into a canonical order, allowing the data to be treated canonically equivalently for search and comparison.

In (2), the Unicode code points have been logically normalized and visually organized into grapheme clusters, as specified by the Unicode Standard. The combining character sequence <̥> is normalized and visually grouped together. Note that, the MODIFIER LETTER SMALL H at U+02B0, is not grouped with any other character. This is because it belongs to Spacing Modifier Letters category in the Unicode Standard. These characters are underspecified for the direction in which they modify a host character. For example, can indicate either pre- or post-aspiration (whereas the nasalization or creaky diacritic is defined in the Unicode Standard to apply to a specified base character).

Finally, to arrive at the graphemic tokenization in (3), tailored grapheme clusters are needed (possibly specified in an orthography profile). For example, an orthography profile might specify that the sequence of characters <ts^h> form a single grapheme. The orthography profile could also specify orthographic rules, e.g. when tokenizing graphemes, in say English, the sequences <sh> in the forms <mishap> and <mishmash> should be treated as distinct sequences depending on their contexts.

5.2 Informal description

An orthography profile describes the Unicode code points, characters, graphemes and orthographic rules in a writing system. An orthography profile is

5 Orthography profiles

a language-specific (and often even resource-specific) description of the units and rules that are needed to adequately model a writing system. An important assumption is that we assume a resource is encoded in Unicode (or has been converted to Unicode). Any data source that the Unicode Standard is unable to capture, will also not be captured by an orthography profile.

Informally, an orthography profile specifies the graphemes (or, in Unicode parlance, `TAILORED GRAPHEME CLUSTERS`) that are expected to occur in any data to be analyzed or checked for consistency. These graphemes are first identified throughout the whole data (a step which we call `TOKENIZATION`), and possibly simply returned as such, possibly including error messages about any parts of the data that are not specified by the orthography profile. Once the graphemes are identified, they might also be changed into other graphemes (a step which we call `TRANSLITERATION`). When a grapheme has different possible transliterations, then these differences should be separated by contextual specification, possibly down to listing individual exceptional cases.

The crucial difference between our current proposal and traditional computational approaches to transliteration is the strict separation between tokenization and transliteration. Most computational approaches to transliteration are based on finite-state transducers (including the transliteration as described in the Unicode Locale Data Markup Language.¹ Finite-state transducers attempt to describe the mapping from input to output string directly as a set of rewrite rules. Although such systems are computationally well understood, we feel that they are not well-suited for day-to-day linguistic practice. First, by forcing a first step of grapheme tokenization, our system tries to keep close to the logic of the writing system. Second, by separating tokenization from transliteration there is no problem with ‘feeding’ and ‘bleeding’ of rules, common with transducers (cf. Section 6.5).

Note that to deal with ambiguous parsing cases, it is still possible to use the Unicode approach including the `ZERO-WIDTH NON-JOINER` character at U+200C into the text. The idea is to add this character into the text to identify cases in which a sequence of characters is *not* supposed to be a complex grapheme cluster (even though the sequence is in the orthography profile).

In practice, we foresee a workflow in which orthography profiles are iteratively refined, while at the same time inconsistencies and errors in the data to be tokenized are corrected. In some more complex use-cases there might even be a need for multiple different orthography profiles to be applied in sequence (see Section 7 on various exemplary use-cases). The result of any such workflow

¹ See <http://www.unicode.org/reports/tr35/>

will normally be a cleaned dataset and an explicit description of the orthographic structure in the form of an orthography profile. Subsequently, the orthography profiles can be easily distributed in scholarly channels alongside the cleaned data, for example in supplementary material added to journal papers or in electronic archives.

5.3 Formal specification

File Format

The formal specifications of an orthography profile (or simply `PROFILE` for short) are the following:

- A1. A PROFILE IS A UNICODE UTF-8 ENCODED TEXT FILE, ideally using NFC, no-BOM, and LF (see Section 2.11), that includes the information pertinent to the orthography.
- A2. A PROFILE IS A TAB-SEPARATED CSV FILE WITH AN OBLIGATORY HEADER LINE. A minimal profile can have just a single column, in which case there will of course be no tabs, but the first line will still be the header. For all columns we assume the name in the header of the CSV file to be crucial. The actual ordering of the columns is unimportant. Because tabs and newlines are field separators of the profile, these characters will lead to problems when they occur in the original data.²
- A3. METADATA ARE ADDED IN A SEPARATE UTF-8 TEXT FILE using a basic TAG: VALUE format. Metadata about the orthographic description given in the orthography profile includes, minimally, (i) author, (ii) date, (iii) title of the profile, (iv) a stable language identifier encoded in BCP 47/ISO 639-3 of the target language of the profile, and (v) bibliographic data for resource(s) that illustrate the orthography described in the profile. Further, (vi) the tokenization method and (vii) the unicode normalisation used should be documented here (see below).
- A4. SEPARATE LINES WITH COMMENTS ARE NOT ALLOWED. Comments that belong to specific lines will have to be put in a separate column of the CSV file, e.g. add a column called `COMMENTS`.

² The main reason to choose for tab-delimited, and not for quoted comma-separated files is that we want the profiles to be easily manipulated by hand. Tab-separation is strongly preferred for this reason.

5 Orthography profiles

The content of a profile consists of lines, each describing a grapheme of the orthography, using the following columns:

- A5. A MINIMAL PROFILE CONSISTS OF A SINGLE COLUMN with a header called `Grapheme`, listing each of the different graphemes in a separate line. The name of this column is crucial for automatic processing.
- A6. OPTIONAL COLUMNS CAN BE USED TO SPECIFY THE LEFT AND RIGHT CONTEXT OF THE GRAPHEME, to be designated with the headers `Left` and `Right`, respectively. The same grapheme can occur multiple times with different contextual specifications, for example to distinguish different pronunciations depending on the context.
- A7. THE COLUMNS `Grapheme`, `Left` AND `Right` CAN USE REGULAR EXPRESSION METACHARACTERS. If regular expressions are used, then all literal usage of the special symbols, like full stops `<.>` or dollar signs `<$>` (so-called METACHARACTERS) have to be explicitly escaped by adding a backslash before them (i.e. use `<\.>` or `<\$>`). Note that any specification of context automatically expects regular expressions, so it is probably better to always escape all regular expression metacharacters when used literally in the orthography. The following symbols will need to be preceded by a backslash: `[] () { } | + * . - ! ? ^ $` and the backslash `\` itself.
- A8. AN OPTIONAL COLUMN CAN BE USED TO SPECIFY CLASSES OF GRAPHEMES, to be identified by the header `Class`. For example, this column can be used to define a class of vowels. Users can simply add ad-hoc identifiers in this column to indicate a group of graphemes, which can then be used in the description of the graphemes or the context. The identifiers should of course be chosen such that they do not conflate with any symbols used in the orthography themselves. Note that such classes only refer to the graphemes, not to the context.
- A9. COLUMNS DESCRIBING TRANSLITERATIONS FOR EACH GRAPHEMES CAN BE ADDED AND NAMED AT WILL. Often more than a single possible transliteration will be of interest. Any software application using these profiles should prompt the user to name any of these columns to select a specific transliteration.
- A10. ANY OTHER COLUMNS CAN BE ADDED FREELY, BUT WILL MOSTLY BE IGNORED BY ANY SOFTWARE APPLICATION USING THE PROFILES. As orthography profiles are also intended to be read and interpreted by humans, it is often highly useful to add extra information on the graphemes in further columns, like for example Unicode codepoints, Unicode names, frequency of occur-

rence, examples of occurrence, explanation of the contextual restrictions, or comments.

Implementation

For the automatic processing of the profiles, the following technical standards will be expected:

- B1. EACH LINE OF A PROFILE WILL BE INTERPRETED AS A REGULAR EXPRESSION. Software applications using profiles can also offer to interpret a profile in the literal sense to avoid the necessity for the user to escape regular expressions metacharacters in the profile. However, this only is possible when no contexts or classes are described, so this seems only useful in the most basic orthographies.
- B2. THE CLASS COLUMN WILL BE USED TO PRODUCE EXPLICIT OR CHAINS OF REGULAR EXPRESSIONS, which will then be inserted in the `Grapheme`, `Left` and `Right` columns at the position indicated by the class-identifiers. For example, a class called `V` as a context specification might be replaced by a regular expression like: `(a|e|i|o|u|ei|au)`. Only the graphemes themselves are included here, not any contexts specified for the elements of the class. Note that in some cases the ordering inside this regular expression might be crucial.
- B3. THE LEFT AND RIGHT CONTEXTS WILL BE INCLUDED INTO THE REGULAR EXPRESSIONS BY USING LOOKBEHIND AND LOOKAHEAD. Basically, the actual regular expression syntax of lookbehind and lookahead is simply hidden to the users by allowing them to only specify the contexts themselves. Internally, the contexts in the columns `Left` and `Right` are combined with the column `Grapheme` to form a complex regular expression like: `(?<=Left)Grapheme(?=Right)`.
- B4. THE REGULAR EXPRESSIONS WILL BE APPLIED IN THE ORDER AS SPECIFIED IN THE PROFILE, FROM TOP TO BOTTOM. A software implementation can offer help in figuring out the optimal ordering of the regular expressions, but should then explicitly report on the order used.

The actual implementation of the profile on some text-string will function as follows:

- B5. ALL GRAPHEMES ARE MATCHED IN THE TEXT BEFORE THEY ARE TOKENIZED OR TRANSLITERATED. In this way, there is no necessity for the user to consider

‘feeding’ and ‘bleeding’ situations, in which the application of a rule either changes the text so another rule suddenly applies (feeding) or prevents another rule to apply (‘bleeding’).

- B6. THE MATCHING OF THE GRAPHEMES CAN OCCUR EITHER GLOBALLY OR LINEARLY. From a computer science perspective, the most natural way to match graphemes from a profile in some text is by walking linearly through the text-string from left to right, and at each position go through all graphemes in the profile to see which one matches, then go to the position at the end of the matched grapheme and start over. This is basically how a finite state transducer works, which is a well-established technique in computer science. However, from a linguistic point of view, our experience is that most linguists find it more natural to think from a global perspective. In this approach, the first grapheme in the profile is matched everywhere in the text-string first, before moving to the next grapheme in the profile. Theoretically, these approaches will lead to different results, though in practice of actual natural language orthographies they almost always lead to the same result. Still, we suggest that any software application using orthography profiles should offer both approaches (i.e. GLOBAL or LINEAR) to the user. The approach used should be documented in the metadata as `TOKENIZATION METHOD`.
- B7. THE MATCHING OF THE GRAPHEMES CAN OCCUR EITHER IN NFC OR NFD. By default, both the profile and the text-string to be tokenized should be treated as NFC (see Section 2.9). However, in some use-cases it turns out to be practical to treat both text and profile as NFD. This typically happens when very many different combinations of diacritics occur in the data. An NFD-profile can then be used to first check which individual diacritics are used, before turning to the more cumbersome inspection of all combinations. We suggest that any software application using orthography profiles should offer both approaches (i.e. NFC or NFD) to the user. The approach used should be documented in the metadata as `UNICODE NORMALIZATION`.
- B8. THE TEXT-STRING IS ALWAYS RETURNED IN TOKENIZED FORM by separating the matched graphemes by a user-specified symbols-string. Any transliteration will be returned on top of the tokenization.
- B9. LEFTOVER CHARACTERS, i.e. CHARACTERS THAT ARE NOT MATCHED BY THE PROFILE, SHOULD BE REPORTED TO THE USER AS ERRORS. Typically, the unmatched character are replaced in the tokenization by a user-specified symbol-string.

Software applications

Any software application offering to use orthography profile:

1. SHOULD OFFER USER-OPTIONS to specify:
 - C1. THE NAME OF THE COLUMN TO BE USED FOR TRANSLITERATION (if any).
 - C2. THE SYMBOL-STRING TO BE INSERTED BETWEEN GRAPHEMES. Optionally, a warning might be given if the chosen string includes characters from the orthography itself.
 - C3. THE SYMBOL-STRING TO BE INSERTED FOR UNMATCHED STRINGS in the tokenized and transliterated output.
 - C4. THE TOKENIZATION METHOD, i.e. whether the tokenization should proceed GLOBAL or LINEAR (see B6 above).
 - C5. UNICODE NORMALIZATION, i.e. whether the text-string and profile should use NFC or NFD.
2. MIGHT OFFER USER-OPTIONS:
 - C6. TO ASSIST IN THE ORDERING OF THE GRAPHEMES. In our experience, it makes sense to apply larger graphemes before shorter graphemes, and to apply graphemes with context before graphemes without context. Further, frequently relevant rules might be applied after rarely relevant rules (though frequency is difficult to establish in practice, as it depends on the available data). Also, if this all fails to give any decisive ordering between rules, it seems useful to offer linguists the option to reverse the ordering from any manual specified ordering, because linguists tend to write the more general rule first, before turning to exceptions or special cases.
 - C7. TO ASSIST IN DEALING WITH UPPER AND LOWER CASE CHARACTERS. It seems practical to offer some basic case matching, so characters like <a> and <A> are treated equally. This will be useful in many concrete cases, although the user should be warned that case matching does not function universally in the same way across orthographies. Ideally, users should prepare orthography profiles with all lowercase and uppercase variants explicitly mentioned, so by default no case matching should be performed.
 - C8. TO TREAT THE PROFILE LITERAL, i.e. to not interpret regular expression metacharacters. Matching graphemes literally often leads to strong speed increase, and would allow users to not needing to worry about escaping metacharacters. However, in our experience all actually

5 *Orthography profiles*

interesting use-cases of orthography profiles include some contexts, which automatically prevents any literal interpretation.

3. SHOULD RETURN THE FOLLOWING INFORMATION to the user:

- C9. THE ORIGINAL TEXT-STRINGS to be processed in the specified unicode normalization, i.e. in either NFC or NFD as specified by the user.
- C10. THE TOKENIZED STRINGS, with additionally any transliterated strings, if transliteration is requested.
- C11. A SURVEY OF ALL ERRORS ENCOUNTERED, ideally both in which text-strings any errors occurred and which characters in the text-strings lead to errors.
- C12. A REORDERED PROFILE, when any automatic reordering is offered

6 Implementation

6.1 Available implementations

To illustrate the practical application of orthography profiles as defined in the previous Chapter 5, we have implemented two different versions of the specifications presented there: one in Python and one in R. These two implementations have rather different implementation histories, and we are not 100% sure that they will in all situations give the same results. Also the performance with larger datasets differs, and the code is not always as clean as we would like it to be. In sum, the two implementations should be considered as ‘proof of concept’ and not as the final word on the practical application of the specifications above. In our own experience, the current implementations are sufficiently fast and stable to be useful for academic practice (e.g. checking data consistency, or analyzing and transliterating small to medium sized data sets), but they should probably better not be used for full-scale industry applications.

This chapter will introduce the implementations, and give practical step-by-step guidelines for installing them and using them. Various simple and sometimes somewhat abstract examples will be discussed to show the different options available, and to illustrate the intended usage of orthography profiles. In the next chapter we will then discuss a few full-scale examples to illustrate the application and the usefulness of orthography profiles.

Installing the R implementation

The R implementation is available in the package `qlcData`, which is directly available from the central R repository CRAN (Comprehensive R Archive Network). The R software environment itself has to be downloaded from its website.¹ After starting the included R program, the `qlcData` package for dealing with orthography profiles can be simply installed as follows:

```
# download and install the qlcData software
install.packages('qlcData')
```

¹ <https://www.r-project.org>

6 Implementation

```
# load the software, so it can be used
library(qlcData)
```

The version available through CRAN is the latest stable version. To obtain the most recent bug-fixes and experimental additions, please use the development version, which is available on GitHub.² This development version can be easily installed using the github-install helper software from the devtools package.

```
# download and install helper software
install.packages('devtools')
# install the qlcData package from GitHub
devtools::install_github('cysouw/qlcData', build_vignettes = TRUE)
# load the software, so it can be used
library(qlcData)
```

Inside the qlcData package, there are two functions for orthography processing, `write.profile` and `tokenize`. R includes help files with illustrative examples, and also a so-called ‘vignette’ with explanations and examples.

```
# view help files
help(write.profile)
help(tokenize)
# view vignette with explanation and examples
vignette('orthography_processing')
```

Basically, the idea is to use `write.profile` to produce a basic orthography profile from some data and then `tokenize` to apply the (possibly edited) profile on some data, as exemplified in the next section. This can of course be performed though R, but additionally there are two more interfaces to the R code supplied in the qlcData package: (i) bash executables and (ii) ‘shiny’ webapps.

The bash executables are little files providing an interface to the R code that can be used in a shell on a UNIX-alike machine. The location of these executables is a bit hidden away by the install procedure of R packages. The location can be found by the following command in R. These executables can be used from here, or they can be linked and/or copied to any location as wanted.

```
# show the path to the bash executables
file.path(find.package('qlcData'), 'exec')
```

For example, a good way to use the executables in a terminal is to make soft-links (using `ln`) from the executables to a directory in your PATH, e.g. to `/usr/local/bin/`. The two executables are named `tokenize` and `writeprofile`, and the links can be made directly by using Rscript to get the paths to the executables within the terminal.

² <http://github.com/cysouw/qlcData>

6.1 Available implementations

```
# get the paths to the R executables in bash
pathT=`Rscript -e 'cat(file.path(find.package("qlcData"), "exec", "tokenize"))'`
pathW=`Rscript -e 'cat(file.path(find.package("qlcData"), "exec", "writeprofile"))'`

# make softlinks to the R executables in /usr/local/bin
# you will have to enter your user's password!
sudo ln -is $pathT $pathW /usr/local/bin
```

After inserting this softlink it should be possible to access the `tokenize` function from the shell. Try `tokenize --help` to test the functionality.

```
tokenize --help
```

```
## USAGE:
##  tokenize [-hlrd] [-s SEP -p REPL] [-t TRANS -m MISSING] [<STRINGS> <PROFILE>]
##
## DESCRIPTION:
##  Using the function tokenize() from the R-package qlcData for tokenization and
##  transliteration of character strings based on an orthography profile. Details
##  http://www.rdocumentation.org/packages/qlcData/functions/tokenize.html
##
##  STRINGS can be piped through. When no PROFILE is specified, default Unicode
##  tokenization is performed. Note that the R-code allows for even more options,
##  not all options are made available here for ease of use.
##  Errors are likewise (not yet) returned
##
## OPTIONS:
##  -h, --help      Showing this help text
##  -l, --linear     Use linear method instead of default global method
##  -r, --regex      Use regex matching, including contexts in matching graphemes
##  -d, --NFD        Use NFD normalization instead of default NFC
##  -s SEP           Separator to be inserted after tokenization
##                   defaults to space [default: SPACE]
##  -p REPL          Replacement symbol in case separator occurs in the STRINGS
##                   defaults to nothing [default: NULL]
##  -t TRANS         Column in profile to use for transliteration [default: NULL]
##  -m MISSING       Character to be inserted for characters not in the PROFILE
##                   defaults to DOUBLE QUESTION MARK at U+2047 [default: ??]
##
## EXAMPLES:
##  Try the difference between the following:
##
##      tokenize áüî
##      tokenize --NFD -s '___' áüî
##
##  Piping strings should work:
##
##      ls | tokenize
```

To make the functionality even more accessible, we have prepared webapps with the shiny framework for the R functions. These webapps are available online at [TODO](#). The webapps are also included inside the `qlcData` package and can be started with the following helper function:

6 Implementation

```
launch_shiny('tokenize')
```

TO DO

Installing the Python implementation

TO DO

6.2 Profiles and error reporting

The first example of how to use these functions concerns finding errors in the encoding of texts. In the following example, it looks as if we have two identical strings, AABB. However, this is just an surface-impression delivered by the current font, which renders Latin and Cyrillic capitals identically. We can identify this problem when we produce an orthography profile from the strings. Using here the R implementation of orthography profiles, we first assign the two strings to a variable `test`, and then produce an orthography profile with the function `write.profile`. As it turns out, some of the letters are Cyrillic.

```
( test <- c('AABB', '\u0041\u0410\u0042\u0412') )

## [1] "AABB" "AABB"

write.profile(test)

##      Grapheme Frequency Codepoint      UnicodeName
## 1         A           3    U+0041    LATIN CAPITAL LETTER A
## 2         B           3    U+0042    LATIN CAPITAL LETTER B
## 3         A           1    U+0410    CYRILLIC CAPITAL LETTER A
## 4         B           1    U+0412    CYRILLIC CAPITAL LETTER VE
```

The function of error-message reporting can also nicely be illustrated with this example. Suppose we made an orthography profile with just the two Latin letters <A> and as possible graphemes, then this profile would not be sufficient to tokenize the strings. There are graphemes in the data that are not in the profile, so the tokenization produces an error, which can be used to fix the encoding (or the profile). In the example below, we can see that the Cyrillic encoding is found in the second string of the `test` input.

```
test <- c('AABB', '\u0041\u0410\u0042\u0412')
tokenize( test, profile = c('A', 'B') )
```

```
## Warning in tokenize(test, profile = c("A", "B")):
## There were unknown characters found in the input data.
## Check output$errors for a table with all problematic strings.

## $strings
##   originals tokenized
## 1      AABB   A A B B
## 2      AABB   A ?? B ??
##
## $profile
##   matched_rules Grapheme
## 1              3       B
## 2              3       A
##
## $errors
##   originals  errors
## 2      AABB A ?? B ??
##
## $missing
##   Grapheme Frequency Codepoint      UnicodeName
## 1      A           1    U+0410  CYRILLIC CAPITAL LETTER A
## 2      B           1    U+0412  CYRILLIC CAPITAL LETTER VE
```

6.3 Different ways to write a profile

The function `write.profile` can be used to prepare a skeleton for an orthography profile from some data. The preparation of an orthography profile from some data might sound like a trivial problem, but actually there are various different ways in which strings can be separated into graphemes by `write.profile`. Consider the following string of characters called `example` below. The default settings of `write.profile` separates the string into Unicode graphemes according to grapheme clusters ('user-perceived characters', see Section 1.4). As it turns out, some of these graphemes are single codepoints, others are combinations of two codepoints (see Section 2.2).

```
example <- '\u00d9\u00da\u00db\u0055\u0300\u0055\u0301\u0055\u0302'
profile_1 <- write.profile(example)
```

Table 6.1: Profile 1 (default settings, splitting grapheme clusters)

Gr.	Freq.	Codepoint	Unicode Name
Ú	1	U+00DA	LATIN CAPITAL LETTER U WITH ACUTE
Û	1	U+0055, U+0301	LATIN CAPITAL LETTER U, COMBINING ACUTE ACCENT
Û	1	U+00D9	LATIN CAPITAL LETTER U WITH GRAVE
Û	1	U+0055, U+0300	LATIN CAPITAL LETTER U, COMBINING GRAVE ACCENT
Û	1	U+00DB	LATIN CAPITAL LETTER U WITH CIRCUMFLEX
Û	1	U+0055, U+0302	LATIN CAPITAL LETTER U, COMBINING CIRCUMFLEX ACCENT

By specifying the splitting separator as the empty string `sep = ""`, it is possible to split the string into Unicode codepoints, thus separating the combining diacritics. The idea behind this option `sep` is that separating by a character allows for user-determined separation. The most extreme choice here is the empty string `sep = ""`, which is interpreted as separation everywhere. The other extreme is the default setting `sep = NULL`, which means that the separation is not user-defined, but relegated to the Unicode grapheme definitions.

```
profile_2 <- write.profile(example, sep = "")
```

Table 6.2: Profile 2 (splitting by codepoints)

Grapheme	Frequency	Codepoint	Unicode Name
´	1	U+0301	COMBINING ACUTE ACCENT
`	1	U+0300	COMBINING GRAVE ACCENT
ˆ	1	U+0302	COMBINING CIRCUMFLEX ACCENT
U	3	U+0055	LATIN CAPITAL LETTER U
Ú	1	U+00DA	LATIN CAPITAL LETTER U WITH ACUTE
Û	1	U+00D9	LATIN CAPITAL LETTER U WITH GRAVE
Û	1	U+00DB	LATIN CAPITAL LETTER U WITH CIRCUMFLEX

Some characters look identical, although they are encoded differently. Unicode offers different ways of normalization (see Section 2.9), which can be invoked here as well using the option `normalize`. NFC-normalization turns everything into the precomposed characters, while NFD-normalization separates everything into base characters with combining diacritics. Splitting by codepoints (i.e. `sep = ""`) shows the results of these two normalizations.

```
# after NFC normalization unicode codepoints have changed
profile_3 <- write.profile(example, normalize = "NFC", sep = "")
# NFD normalization gives yet another structure of the codepoints
profile_4 <- write.profile(example, normalize = "NFD", sep = "")
```

Table 6.3: Profile 3 (splitting by NFC codepoints)

Grapheme	Frequency	Codepoint	Unicode Name
Ú	2	U+00DA	LATIN CAPITAL LETTER U WITH ACUTE
Û	2	U+00D9	LATIN CAPITAL LETTER U WITH GRAVE
Û	2	U+00DB	LATIN CAPITAL LETTER U WITH CIRCUMFLEX

Table 6.4: Profile 4 (splitting by NFD codepoints)

Grapheme	Frequency	Codepoint	Unicode Name
´	2	U+0301	COMBINING ACUTE ACCENT
`	2	U+0300	COMBINING GRAVE ACCENT
ˆ	2	U+0302	COMBINING CIRCUMFLEX ACCENT
U	6	U+0055	LATIN CAPITAL LETTER U

It is important to realize that for Unicode grapheme definitions, NFC and NFD normalization are equivalent. This can be shown by normalizing the example in either NFC or NFD, but using default separation in `write.profile`. To be precise, default separation means setting `sep = NULL`, but that has not been added explicitly.

```
# note that NFC and NFD normalization are identical
# for unicode grapheme definitions
profile_5 <- write.profile(example, normalize = "NFD")
profile_6 <- write.profile(example, normalize = "NFC")
```

Table 6.5: Profile 5 (splitting by graphemes after NFD)

Gr.	Freq.	Codepoint	Unicode Name
Ú	2	U+0055, U+0301	LATIN CAPITAL LETTER U, COMBINING ACUTE ACCENT
Û	2	U+0055, U+0300	LATIN CAPITAL LETTER U, COMBINING GRAVE ACCENT
Û	2	U+0055, U+0302	LATIN CAPITAL LETTER U, COMBINING CIRCUMFLEX ACCENT

Table 6.6: Profile 6 (splitting by graphemes after NFC)

Gr.	Freq.	Codepoint	Unicode Name
Ú	2	U+00DA	LATIN CAPITAL LETTER U WITH ACUTE
Û	2	U+00D9	LATIN CAPITAL LETTER U WITH GRAVE
Û	2	U+00DB	LATIN CAPITAL LETTER U WITH CIRCUMFLEX

6 Implementation

Note that these different profiles can also be produced using the bash executable `writeprofile` (see 6.1 for notes on how to install the bash executable). Exactly this example is also included in the help file of the executable.

```
writeprofile --help
```

```
## USAGE:
##   writeprofile [-hiec] [-s SEP] [-n NORMALIZE] [<STRINGS>]
##
## DESCRIPTION:
##   Using the function write.profile() from the R-package qlcData to prepare
##   a skeleton for an orthography profile to be used with tokenize. Details:
##   http://www.rdocumentation.org/packages/qlcData/functions/readwrite_orthography
##
##   STRINGS can be piped through. The results is a tab-delimited file for manual
##   editing. Note that the R-code allows for even more options, not all options
##   are made available here for ease of use.
##
## OPTIONS:
##   -h, --help           Showing this help text
##   -i, --info           Add columns with Unicode information to profile
##   -e, --editing        Add empty columns for easy editing of profile
##   -c, --codepoints     Separate by Unicode codepoints (overrides any -s setting)
##   -s SEP               Separator used to separate the strings
##                       defaults to unicode character definition [default: NULL]
##   -n NORMALIZE         Normalization used to make profile. Possible options: NFD, NFC
##                       defaults to no normalization [default: NULL]
##
## EXAMPLES:
##   Note the differences between the following variants:
##
##   example='ÛÛÛÛÛ'
##   echo -e $example
##   echo -e $example | writeprofile -i
##   echo -e $example | writeprofile -ic
##   echo -e $example | writeprofile -ic -n NFC
##   echo -e $example | writeprofile -ic -n NFD
```

6.4 Using an orthography profile skeleton

A common workflow to use these functions is to first make a skeleton for an orthography profile and then edit this profile by hand. For example, Table 6.7 shows the profile skeleton after a few graphemes have been added to the file. Note that in this example, the profile is written to the desktop, and this file has to be edited manually. We simply add a few multigraphs to the column `Grapheme` and leave the other columns empty. these new graphemes are then included in the graphemic parsing.

6.4 Using an orthography profile skeleton

```
# a few words to be graphemically parsed
example <- c("mishmash", "mishap", "mischief", "scheme")
# write a profile skeleton to a file
write.profile(example, file = "~/Desktop/profile_skeleton.txt")
# edit the profile, and then use the edited profile to tokenize
tokenize(example, profile = "~/Desktop/profile_skeleton.txt")$strings

##  originals      tokenized
## 1  shampoo  sh a m p oo
## 2  mishap   m i sh a p
## 3  mischief m i sch ie f
## 4  scheme   sch e m e
```

Table 6.7: Manually edited profile skeleton

Grapheme	Frequency	Codepoint	UnicodeName
sh			
ch			
sch			
ie			
oo			
a	2	U+0061	LATIN SMALL LETTER A
c	2	U+0063	LATIN SMALL LETTER C
e	3	U+0065	LATIN SMALL LETTER E
f	1	U+0066	LATIN SMALL LETTER F
h	4	U+0068	LATIN SMALL LETTER H
i	3	U+0069	LATIN SMALL LETTER I
m	4	U+006D	LATIN SMALL LETTER M
o	2	U+006F	LATIN SMALL LETTER O
p	2	U+0070	LATIN SMALL LETTER P
s	4	U+0073	LATIN SMALL LETTER S

It is also possible to leave out the Unicode information in the profile skeleton by using the option `info = FALSE`. It is also possible not to use a separate file at all, but process everything within R. In simple situations this is often useful (see below), but in general we prefer to handle everything through a separately saved orthography profile. This profile often contains highly useful information that is nicely coded and saved inside this one file, and can thus be easily distributed and shared. Doing the same as above completely within R might look as follows:

```
# make a profile, just select the column 'Grapheme'
profile <- write.profile(example)[, "Grapheme"]
# extend the profile with multigraphs
profile <- c("sh", "ch", "sch", "ie", "oo", profile)
# use the profile to tokenize
tokenize(example, profile)$strings
```

6 Implementation

```
##      originals      tokenized
## 1    shampoo    sh a m p oo
## 2     mishap    m i sh a p
## 3  mischief m i sch ie f
## 4     scheme    sch e m e
```

6.5 Rule ordering

Not yet everything is correct in this graphemic parsing of the example discussed previously. The sequence <sh> in ‘mishap’ should not be a digraph, and conversely the sequence <sch> in ‘mischief’ should of course be separated into <s> and <ch>. One of the important issues to get the graphemic parsing right is the order in which graphemes are parsed. For example, currently the grapheme <sch> is parsed before the grapheme <ch>, leading to <m i sch ie f> instead of the intended <m i s ch ie f>. The reason that <sch> is parsed before <ch> is that by default longer graphemes are parsed before shorter ones. Our experience is that in most cases this is expected behaviour. You can change the ordering by specifying the option `ordering`. Setting this option to `NULL` results in no preferential ordering, i.e. the graphemes are parsed in the order of the profile, from top to bottom. Now ‘mischief’ is parsed correctly, but ‘scheme’ is wrong. So this ordering is not the solution in this case.

```
# do not reorder the profile
# just apply the graphemes from top to bottom
tokenize( example
  , profile = "~/Desktop/profile_skeleton.txt"
  , ordering = NULL
) $strings
```

```
##      originals      tokenized
## 1    shampoo    sh a m p oo
## 2     mishap    m i sh a p
## 3  mischief m i s ch ie f
## 4     scheme    s ch e m e
```

There are various additional options for rule ordering implemented. Please check the help description in `R help(tokenize)` for more details on the possible rule ordering specifications. In summary, there are four different ordering options, that can also be combined:

- `SIZE`

This option orders the lines in the profile by the size of the grapheme, largest first. Size is measured by number of Unicode characters after normalization as specified in the option `normalize`. For example, `<é>` has a size of 1 with `normalize = "NFC"`, but a size of 2 with `normalize = "NFD"`.

- **CONTEXT**

This option orders the lines by whether they have any context specified (see next section). Lines with context will then be used first. Note that this only works when the option `regex = TRUE` is also chosen (otherwise context specifications are not used).

- **REVERSE**

This option orders the lines from bottom to top. Reversing order can be useful because hand-written profiles tend to put general rules before specific rules, which mostly should be applied in reverse order.

- **FREQUENCY**

This option orders the lines by the frequency with which they match in the specified strings before tokenization, least frequent coming first. This frequency of course depends crucially on the available strings, so it will lead to different orderings when applied to different data. Also note that this frequency is (necessarily) measured before graphemes are identified, so these ordering frequencies are not the same as the final frequencies shown in the output. Frequency of course also strongly differs on whether context is used for the matching through `regex = TRUE`.

By specifying more than one ordering, these orderings are used to break ties, e.g. the default setting `ordering = c("size", "context", "reverse")` will first order by size, and for those with the same size, it will order by whether there is any context specified or not. For lines that are still tied (i.e. they have the same size and both/neither have context) the order will be reversed compared to the order as attested in the profile, because most hand-written specifications of graphemes will first write the general rule, followed by more specific regularities. To get the right tokenization, these rules should in most cases be applied in reverse order.

It is important to realize that different ordering of the rules does not have so-called ‘feeding’ and ‘bleeding’ effects as known from finite-state rewrite rules.³

³ ‘Bleeding’ is the effect that the application of a rule changes the string, so as to prevent a next rule to apply. ‘Feeding’ is the opposite: a specific rule will only be applied because a previous rule changed the string already. The interaction of rules with such ‘feeding’ and ‘bleeding’ effects is extremely difficult to predict.

6 Implementation

The graphemic parsing advocated here works crucially different from rewrite rules in that there is nothing being rewritten: each line in an orthography profile specifies a grapheme to be ‘captured’ in the string. All lines in the profile are processed in a specified order (as determined by the option `ordering`). At the processing of a specific line, all matching graphemes in the data are marked as ‘captured’, but not changed. Captured parts cannot be captured again, but they can still be used to match contexts of other lines in the profile. Only when all lines are processed the captured graphemes are separated (and possibly transliterated). In this way the result of the application of the ‘rules’ is rather easy to predict.

To document a specific case of graphemic parsing, it is highly useful to save all results of the tokenization to file by using the option `file.out`, for example as follows:

```
# save the results to various files
tokenize( example
  , profile = "~/Desktop/profile_skeleton.txt"
  , file.out = "~/Desktop/result"
)
```

This will lead to the following four files being written. Crucially, a new profile is produced with the re-ordered orthography profile. To reproduce the tokenization, this re-ordered profile can be used with the option `ordering = NULL`.

- **RESULT_STRINGS.TSV:**
A tab-separated file with the original and the tokenized/transliterated strings
- **RESULT_PROFILE.TSV:**
A tab-separated file with the graphemes with added frequencies of occurrence in the data. The lines in the file is re-ordered according to the order that resulted from the ordering specifications (see Section 6.5).
- **RESULT_ERRORS.TSV:**
A tab-separated file with all original strings that contain unmatched parts. Unmatched parts are indicated with the character as specified with the option `missing`. By default the character DOUBLE QUESTION MARK `<??>` at U+2047 is used. When there are no errors, then this file is absent.
- **RESULT_MISSING.TSV:**
A tab-separated file with the graphemes that are missing from the original orthography profile, as indicated in the errors. When there are no errors, then this file is absent.

6.6 Contextually specified graphemes

To refine a profile, it is also possible to add graphemes with contextual specifications. An orthography profile can have columns called `Left` and `Right` to specify the context in which the grapheme is to be separated.⁴ For example, we are adding an extra line to the profile from above, resulting in the profile shown in Table 6.8. The extra line specifies that `<s>` is a grapheme when it occurs after `<mi>`. Such contextually-specified graphemes are based on regular expressions so you can also use regular expressions in the description of the context. For such contextually specified graphemes to be included in the graphemic parsing we have to specify the option `regex = TRUE`. This contextually specified grapheme should actually be handled first, so we could try `ordering = NULL`. However, we can also explicitly specify that rules with contextual information should be applied first by using `ordering = "context"`. That gives the right results for this toy example.

```
# add a contextual grapheme, and then use the edited profile to tokenize
tokenize( example
  , profile = "~/Desktop/profile_skeleton.txt"
  , regex = TRUE
  , ordering = "context"
) $strings

##      originals      tokenized
## 1  shampoo  sh a m p oo
## 2   mishap  m i s h a p
## 3 mischief m i s ch ie f
## 4   scheme  s ch e m e
```

It is important to realize that with the option `regex = TRUE` all content in the profile is treated as regular expressions, so the characters with special meaning in regular expressions should be either omitted or escaped (by putting a `< \ >` REVERSE SOLIDUS at U+005C before the character). Specifically, this concerns the following characters:

```
<-> HYPHEN-MINUS at U+002D
<!> EXCLAMATION MARK at U+0021
<?> QUESTION MARK at U+003F
<.> FULL STOP at U+002E
```

⁴ The column names `Left`, `Right` and `Grapheme` are currently hard-coded, so exactly these column names should be used for these effects to take place. The position of the columns in the profile is unimportant. So the column `Left` can occur anywhere.

Table 6.8: Orthography profile with contextual specification for <s>

Left	Grapheme	Frequency	Codepoint	UnicodeName
mi	s			
	sh			
	ch			
	sch			
	ie			
	oo			
	a	2	U+0061	LATIN SMALL LETTER A
	c	2	U+0063	LATIN SMALL LETTER C
	e	3	U+0065	LATIN SMALL LETTER E
	f	1	U+0066	LATIN SMALL LETTER F
	h	4	U+0068	LATIN SMALL LETTER H
	i	3	U+0069	LATIN SMALL LETTER I
	m	4	U+006D	LATIN SMALL LETTER M
	o	2	U+006F	LATIN SMALL LETTER O
	p	2	U+0070	LATIN SMALL LETTER P
	s	4	U+0073	LATIN SMALL LETTER S

<(> LEFT PARENTHESIS at U+0028
 <)> RIGHT PARENTHESIS at U+0029
 <[> LEFT SQUARE BRACKET at U+005B
 <]> RIGHT SQUARE BRACKET at U+005D
 <{> LEFT CURLY BRACKET at U+007B
 <}> RIGHT CURLY BRACKET at U+007D
 <|> VERTICAL LINE at U+007C
 <*> ASTERISK at U+002A
 <\\> REVERSE SOLIDUS at U+005C
 <^> CIRCUMFLEX ACCENT at U+005E
 <+> PLUS SIGN at U+002B
 <\$> DOLLAR SIGN at U+0024

6.7 Profile skeleton with columns for editing

When it is expected that context might be important for a profile, then the profile skeleton can be created with columns prepared for the contextual specifications. This is done by using the option `editing = TRUE` (cf. Table 6.9 for a toy profile of some Italian words).

```
example <- c('cane', 'cena', 'cine')
write.profile(example
              , file = "~/Desktop/profile_skeleton.txt")
```

```
, editing = TRUE
, info = FALSE
)
```

Table 6.9: Orthography profile with empty columns for editing contexts

Left	Grapheme	Right	Class	Replacement
	a			a
	c			c
	e			e
	i			i
	n			n

Besides the columns `Left`, `Grapheme`, and `Right` as discussed in the previous sections, there are also columns `Class` and `Replacement`. The column `Class` can be used to specify classes of graphemes that can then be used in the contextual specification. The column `Replacement` is just a copy of the column `Grapheme`, providing a skeleton to specify transliteration. The name of the column `Replacement` is not fixed – there can actually be multiple columns with different kinds of transliterations in a single profile.

To achieve contextually determined replacements it is possible to use a regular expression in the contextual column. For example, consider the edited toy profile for Italian in Table 6.10 (where `<c>` becomes `/k/` except before `<i,e>`, then it becomes `/tʃ/`).

Table 6.10: Orthography profile with regex as context

Left	Grapheme	Right	Class	IPA
	c	[ie]		tʃ
	a			a
	n			n
	c			k
	e			e
	i			i

To use this profile, you have to add the option `regex = TRUE`. Also note that we have changed the name of the transliteration-column, so we have to tell the tokenization process to use this column to transliterate. This is done by adding the option `transliterate = "IPA"`.

```
# add a contextual grapheme, and then use the edited profile to tokenize
tokenize( example
```

6 Implementation

```
, profile = "~/Desktop/profile_skeleton.txt"
, regex = TRUE
, transliterate = "IPA"
)$strings

## originals tokenized transliterated
## 1      cane  c a n e      k a n e
## 2      cena  c e n a      tʃ e n a
## 3      cine  c i n e      tʃ i n e
```

Another equivalent possibility is to use a column `Class` to specify a class of graphemes, and then use this class in the specification of context. This is useful to keep track of recurrent classes in larger profiles. You are free to use any class-name you like, as long as it doesn't clash with the rest of the profile. The example shown in Table 6.11 should give the same result as obtained previously by using a regular expression.

Table 6.11: Orthography profile with `Class` as context

Left	Grapheme	Right	Class	IPA
	c	Vfront		tʃ
	a			a
	n			n
	c			k
	e		Vfront	e
	i		Vfront	i

```
# add a class, and then use the edited profile to tokenize
tokenize( example
, profile = "~/Desktop/profile_skeleton.txt"
, regex = TRUE
, transliterate = "IPA"
)$strings

## originals tokenized transliterated
## 1      cane  c a n e      k a n e
## 2      cena  c e n a      tʃ e n a
## 3      cine  c i n e      tʃ i n e
```

6.8 Formatting grapheme separation

In all examples above we have used the default formatting for grapheme separation using space as a separator, which is obtained by the default setting `sep = " "`.

It is possible to specify any other separator here, including the empty string, i.e. `sep = ""`. This will not show the graphemic tokenization anymore (although it has of course been used in the background).

```
# Use the empty string as separator
tokenize( example
  , profile = "~/Desktop/profile_skeleton.txt"
  , regex = TRUE
  , transliterate = "IPA"
  , sep = ""
) $strings

##   originals tokenized transliterated
## 1      cane      cane          kane
## 2      cena      cena      tʃena
## 3      cine      cine      tʃine
```

Normally, the separator specified should not occur in the data. If it does, unexpected things might happen, so consider removing the chosen separator from your strings first. However, there is also an option `sep.replace` to replace the separator with something else. When `sep.replace` is specified, this mark is inserted in the string at those places where the separator occurs. Typical usage in linguistics would be `sep = " "`, `sep.replace = "#"` adding spaces between graphemes and replacing spaces in the input string by hashes in the output string.

```
# Replace separator in string to be tokenized
tokenize( "test test test"
  , sep = " "
  , sep.replace = "#"
) $strings$tokenized

## [1] "t e s t # t e s t # t e s t"
```

6.9 Remaining issues

Given a set of graphemes, there are at least two different methods to tokenize strings. The first is called `method = "global"`. This approach takes the first grapheme in the profile, then matches this grapheme globally at all places in the string, and then turns to the next string in the profile. To us, this seems to be the most intuitive way for linguists to deal with graphemic parsing. The other approach is called `method = "linear"`. This approach walks through the

6 Implementation

string from left to right. At the first character it looks through all graphemes whether there is any match, and then walks further to the end of the match and starts again. This approach is more alike to finite-state rewrite rules (though note that it still works differently from such rewrite rules, see Section 6.5). The global method is used by default.

In some special cases these two tokenization methods can lead to different results, but these special situations are very unlikely to happen in natural language. The example below shows that a string 'abc' can be parsed differently in case of a very special profile with a very special ordering of the graphemes.

```
# different parsing methods can lead to different results
# the global method first catches 'bc'
tokenize( "abc"
  , profile = c("bc", "ab", "a", "c")
  , order = NULL
  , method = "global"
)$strings

##   originals tokenized
## 1      abc      a bc

# the linear method catches the first grapheme, which is 'ab'
tokenize( "abc"
  , profile = c("bc", "ab", "a", "c")
  , order = NULL
  , method = "linear"
)$strings

##   originals tokenized
## 1      abc      ab c
```

Further, the current R-implementation has a limitation when regular expressions are used. The problem is that overlapping matches are not captured when using regular expressions.⁵ Everything works as expected without regular expressions, but there might be warnings/errors in case of `regex = TRUE`. However, just as in the previous issue, this problem should only very rarely (when at all) happen in natural language data.

The problem can be exemplified by a sequence <bbbb> in which a grapheme <bb> should be matched. With the default `regex = FALSE` there are three possible matches, but with `regex = TRUE` only the first two s or the last two

⁵ This restriction is an effect of the underlyingly used ICU implementation of the Unicode standard as implemented in R through the package `stringi`.

s are matched. The middle two s are not matched because they overlap with the other matches. In the example below this leads to an error, because the second <bb> is not matched. However, we have not been able to produce a real example in any natural language in which this limitation might lead to an error.

```
# Everything perfekt without regular expressions
tokenize( "abbb"
, profile = c("ab", "bb")
, order = NULL
, regex = FALSE
)$strings

##    originals tokenized
## 1      abbb      ab bb

# Matching with regular expressions does not catch overlap
tokenize( "abbb"
, profile = c("ab", "bb")
, order = NULL
, regex = TRUE
)$strings

## Warning in tokenize("abbb", profile = c("ab", "bb"), order =
NULL, regex = TRUE):
## There were unknown characters found in the input data.
## Check output$errors for a table with all problematic strings.

##    originals tokenized
## 1      abbb      ab ?? ??
```


7 Use cases

7.1 Introduction

We now present several use cases that have motivated the development of orthography profiles. These include:

- tokenization and error checking
- normalization of orthographic systems to attain interoperability across different source documents,
- cognate identification for detecting the similarity between words from different languages
- ...

An important assumption of our work is that input sources are encoded in Unicode (UTF-8 to be precise). Anything that the Unicode Standard is unable to capture, cannot be captured by an orthography profile. We also use Unicode normalization form NFD to decompose all incoming text input into normalized and (Unicode) logically ordered strings. This process is of fundamental importance when working with text data in the Unicode Standard.

7.2 Tokenization and error checking

The most basic use case is tokenization of language data. Tokenization comes in three types:

- Unicode character tokenization
- grapheme tokenization
- tailored grapheme tokenization

The simplest form of tokenization provided by any Unicode Standard-compliant software will split an input text stream on Unicode character code points. The input text is split on the character bit sequences as they have been encoded by the user. This may be in various Unicode Normalization Forms (see Section ??)

and tokenization split function returns a byte stream sequence given a particular encoding (see Section ??).

There is increasing support for the regular expression match, commonly “X”, to identify Unicode graphemes (see Section ??) and perform tokenization on these sets of Unicode characters. Examples are given in Section ??.

The orthography profile provides the mechanism for tokenizing an input stream on tailored grapheme clusters. The orthography’s formal specification is given in Section ?. To sum, an orthography profile is a list of tailored grapheme clusters and / or orthographic rules that specify how a specific resource of text input should be tokenized.

Once text has been normalized, an additional and straightforward process is Unicode grapheme tokenization. More recently this regular expression, often available as short cut “X”, identifies all sequences of “base” characters followed by 0 or more combining diacritics. For example, the sequence, a :>. For a first pass at identifying grapheme clusters, this is a straightforward tokenization. However, the linguist will note that both the tie bar (COMBINING DOUBLE INVERTED BREVE, a Unicode Combining Modifier) and the length marker (MODIFIER LETTER TRIANGULAR COLON, a Unicode Letter Modifier) do not appear “correctly”, i.e. the tie bar is grouped with the first character in the sequence and the length marker appears singly. This is necessary, as defined by the Unicode Standard, because of these character’s semantic properties – with Unicode we cannot know if, for example, the MODIFIER LETTER TRIANGULAR COLON should appear before or after the base character that it modifies, cf. the IPA aspiration marker, , the MODIFIER LETTER SMALL H, which is also a Unicode Letter Modifier and for linguists a diacritic that be be used for pre- or post-aspiration. These ambiguities of position in tokenization are not uncommon in IPA, thus orthographic (or source) normalization and tokenization is needed.

Bibliography

- Anderson, Lloyd B. 1984. Multilingual text processing in a two-byte code. In *Proceedings of the 10th international conference on computational linguistics* (COLING '84), 1–4. Stanford, California: Association for Computational Linguistics. <http://dx.doi.org/10.3115/980431.980492>.
- Association, International Phonetic et al. 1989. Report on the 1989 kiel convention. *Journal of the International Phonetic Association* 19(2). 67–80.
- Becker, Joseph D. 1984. Multilingual word processing. *Scientific American* 251(1). 96–107.
- Beider, Alexander & Stephen P. Morse. 2008. Beider-morse phonetic matching: an alternative to soundex with fewer false hits. *Avotaynu: the International Review of Jewish Genealogy* 24(2). 12. <https://web.archive.org/web/20150513075152/http://stevemorse.org/phonetics/bmpm.htm>.
- Belongie, Serge, Jitendra Malik & Jan Puzicha. 2002. Shape matching and object recognition using shape contexts. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24(4). 509–522.
- Bird, Steven & Gary F. Simons. 2003. Seven Dimensions of Portability for Language Documentation and Description. *Language* 79(3). 557–582. <http://www.language-archives.org/documents/portability.pdf>.
- Brown, Cecil H., Eric W. Holman & Søren Wichmann. 2013. Sound correspondences in the world's languages. *Language* 89(1). 4–29.
- Computer, Apple. 1985. The font manager. *Inside Machintosh* 1. 215–240.
- Computer, Apple. 1986. The font manager. *Inside Machintosh* 4. 27–46.
- Computer, Apple. 1988. The script manager. *Inside Machintosh* 5. 293–322.
- Daniels, Peter T. 1990. Fundamentals of Grammatology. *Journal of the American Oriental Society* 110(4). 727–731.
- Daniels, Peter T. 1996. The Study of Writing Systems. In P. T. Daniels & W. Bright (eds.), *The world's writing systems*. New York, NY: Oxford University Press.
- Daniels, Peter T. & William Bright. 1996. *The World's Writing Systems*. New York, NY: Oxford University Press.
- Dolgopolsky, A. B. 1986. A probabilistic hypothesis concerning the oldest relationships among the language families of northern eurasia. In Vitalij V. Shevoroshkin

Bibliography

- (ed. and trans.). Trans. from the Russian by Vitalij V. Shevoroshkin, *Typology, relationship and time: A collection of papers on language change and relationship by soviet linguists*, 27–50. Ann Arbor: Karoma Publisher.
- Esling, John. 1990. Computer coding of the ipa: supplementary report. *Journal of the International Phonetic Association* 20(01). 22–26.
- Esling, John H. & Harry Gaylord. 1993. Computer codes for phonetic symbols. *Journal of the International Phonetic Association* 23(2). 83–97.
- Evans, Nicholas & Stephen C Levinson. 2009. The myth of language universals: language diversity and its importance for cognitive science. *Behavioral and brain sciences* 32(05). 429–448.
- Gaultney, J. Victor. 2002. *Problems of diacritic design for latin script text faces*. SIL International.
- Hieronymus, James L. 1993. Ascii phonetic symbols for the world's languages: worldbet. *Journal of the International Phonetic Association* 23.
- Huurdeman, Anton A. 2003. *The worldwide history of telecommunications*. New York, NY: John Wiley & Sons.
- International Phonetic Association. 2005. *International Phonetic Alphabet*. Tech. rep. International Phonetic Association. <http://www.arts.gla.ac.uk/IPA/>.
- Kemp, Alan. 2006. Phonetic transcription: history. *The Encyclopedia of Language and Linguistics* 6. 396–410.
- Knuth, Donald Ervin. 1973. *The art of programming: sorting and searching*. Vol. 3. Reading, Mass.: Addison-Wesley.
- Kohrt, Manfred. 1986. The Term 'Grapheme' in the History and Theory of Linguistics. In Gerhard Augst (ed.), *New trends in graphemics and orthography*, 80–96. Berlin: de Gruyter.
- Ladefoged, Peter. 1990. Some Reflections on the IPA. *Journal of Phonetics* 18. 335–346.
- List, Johann-Mattis. 2012. SCA: phonetic alignment based on sound classes. In Daniel Lassiter & Marija Slavkovik (eds.), *New directions in logic, language, and computation*, 32–51. Berlin & Heidelberg: Springer.
- Maddieson, Ian. 1984. *Pattern of Sounds*. Cambridge, UK: Cambridge University Press.
- Mania, Hubert. 2008. *Gauss: eine biographie*. Reinbek bei Hamburg: Rowohlt.
- Mc Laughlin, Fiona. 2005. Voiceless Implosives in Seereer-Siin. *Journal of the International Phonetic Association* 35(2). 201–214.
- Meinhof, Carl & Daniel Jones. 1928. Principles of Practical Orthography for African Languages. *Africa: Journal of the International African Institute* 1(2). 228–239.

- Meyer, Julien. 2015. *Whistled languages: a worldwide inquiry on human whistled speech*. Berlin: Springer. <http://search.ebscohost.com/login.aspx?direct=true&scope=site&db=nlebk&db=nlabk&AN=957780>.
- Meyer, Julien, Laure Dentel & Frank Seifart. 2012. A methodology for the study of rhythm in drummed forms of languages: application to bora manguaré of amazon. In *Proceedings of interspeech*, 686–690.
- Mielke, Jeff. 2009. Segment Inventories. *Language and Linguistics Compass* 3/2. 700–718.
- Moran, Steven. 2006. *A grammatical sketch of western sisaala*. Ypsilanti, MI: Eastern Michigan University MA thesis.
- Moran, Steven. 2012. *Phonetics Information Base and Lexicon*. University of Washington PhD thesis.
- Olúmúyìw, Tèmitópé. 2013. Yoruba writing: standards and trends. *Journal of Arts and Humanities* 2(1). 40–51. <http://www.theartsjournal.org/index.php/site/article/view/50>.
- Postel, Hans Joachim. 1969. Die kölnen phonetik: ein verfahren zur identifizierung von personennamen auf der grundlage der gestaltanalyse. *IBM-Nachrichten* 19. 925–931.
- Powell, Barry B. 2012. *Writing theory and history of the technology of civilization*. New York, NY: John Wiley & Sons. <http://nbn-resolving.de/urn:nbn:de:101:1-201502112957>.
- Pullum, Geoffrey K.. & William Allen Ladusaw. 1986. *Phonetic symbol guide*. University of Chicago Press.
- Robinson, Andrew. 1995. *The story of writing*. London: Thames & Hudson.
- Simons, Gary F. 1989. Working With Special Characters. In *Occasional publications in academic computing* (Occasional Publications in Academic Computing). Dallas, TX: Summer Institute of Linguistics.
- Simons, Gary F. 1996. The nature of linguistic data and the requirements of a computing environment for linguistic research. *Dutch Studies on Near Eastern Languages and Literature* 2(1). 111–128.
- Singh, Simon. 1999. *The code book*. New York: Doubleday.
- Sproat, Richard. 2000. *A Computational Theory of Writing Systems*. Cambridge, UK: Cambridge University Press.
- The International Phonetic Association. 1999. *Handbook of the International Phonetic Association: A Guide to the Use of the International Phonetic Alphabet*. Cambridge University Press.

Subject index

- The International Phonetic Association. 2007. *Handbook of the International Phonetic Association: A Guide to the Use of the International Phonetic Alphabet*. Vol. Eighth printing. Cambridge University Press.
- Unicode Consortium, The. 2014. *The Unicode Standard, Version 7.0.0*. Tech. rep. Mountain View, CA: The Unicode Consortium.
- Wells, John, William Barry, Martine Grice, Adrian Fourcin & Dafydd Gibbon. 1992. Standard computer-compatible transcription. *Esprit project 2589 (SAM)*, Doc. no. SAM-UCL 37.
- Wells, John C. 1987. Computer-coded phonetic transcription. *Journal of the International Phonetic Association* 17(02). 94–114.
- Wells, John C. 1995. Computer-coding the ipa: a proposed extension of sampa. *Revised draft 4(28)*. 1995.