*Productive*

recognition this is exactly what we do. But for many NLP applications this isn't possible because *-ing* is a **productive** suffix; by this we mean that it applies to every verb. Similarly *-s* applies to almost every noun. Productive suffixes even apply to new words; thus the new word *fax* can automatically be used in the *-ing* form: *faxing*. Since new words (particularly acronyms and proper nouns) are created every day, the class of nouns in English increases constantly, and we need to be able to add the plural morpheme *-s* to each of these. Additionally, the plural form of these new nouns depends on the spelling/pronunciation of the singular form; for example if the noun ends in *-z* then the plural form is *-es* rather than *-s*. We'll need to encode these rules somewhere.

Finally, we certainly cannot list all the morphological variants of every word in morphologically complex languages like Turkish, which has words like:

(3.1) uygarlaştıramadıklarımızdanmışsınızcasına
*uygar*   *+laş* *+tır*   *+ama*   *+dık*   *+lar* *+ımız* *+dan* *+mış*   *+sınız* *+casına*
civilized +BEC +CAUS +NABL +PART +PL  +P1PL +ABL +PAST +2PL   +AsIf

"(behaving) as if you are among those whom we could not civilize"

The various pieces of this word (the **morphemes**) have these meanings:

|  |  |
|---|---|
| +BEC | "become" |
| +CAUS | the causative verb marker ('cause to X') |
| +NABL | "not able" |
| +PART | past participle form |
| +P1PL | 1st person pl possessive agreement |
| +2PL | 2nd person pl |
| +ABL | ablative (from/among) case marker |
| +AsIf | derivationally forms an adverb from a finite verb |

Not all Turkish words look like this; the average Turkish word has about three morphemes. But such long words do exist; indeed Kemal Oflazer, who came up with this example, notes (p.c.) that verbs in Turkish have 40,000 possible forms not counting derivational suffixes. Adding derivational suffixes, such as causatives, allows a theoretically infinite number of words, since causativization can be repeated in a single word (*You cause X to cause Y to . . . do W*). Thus we cannot store all possible Turkish words in advance, and must do morphological parsing dynamically.

In the next section we survey morphological knowledge for English and some other languages. We then introduce the key algorithm for morphological parsing, the **finite-state transducer**. Finite-state transducers are a crucial technology throughout speech and language processing, so we will return to them again in later chapters.

After describing morphological parsing, we will introduce some related algorithms in this chapter. In some applications we don't need to parse a word, but we do need to map from the word to its root or stem. For example in information retrieval and web search (IR), we might want to map from *foxes* to *fox*; but might not need to also know

*Stemming*

that *foxes* is plural. Just stripping off such word endings is called **stemming** in IR. We will describe a simple stemming algorithm called the **Porter stemmer**.

For other speech and language processing tasks, we need to know that two words have a similar root, despite their surface differences. For example the words *sang*, *sung*, and *sings* are all forms of the verb *sing*. The word *sing* is sometimes called the common

*Lemmatization*

*lemma* of these words, and mapping from all of these to *sing* is called **lemmatization**.[2]

<span style="float:left">*Tokenization*</span>

Next, we will introduce another task related to morphological parsing. **Tokenization** or **word segmentation** is the task of separating out (tokenizing) words from running text. In English, words are often separated from each other by blanks (whitespace), but whitespace is not always sufficient; we'll need to notice that *New York* and *rock 'n' roll* are individual words despite the fact that they contain spaces, but for many applications we'll need to separate *I'm* into the two words *I* and *am*.

Finally, for many applications we need to know how similar two words are orthographically. Morphological parsing is one method for computing this similarity, but another is to just compare the strings of letters to see how similar they are. A common way of doing this is with the **minimum edit distance** algorithm, which is important throughout NLP. We'll introduce this algorithm and also show how it can be used in spell-checking.

# 3.1    Survey of (Mostly) English Morphology

<span style="float:left">*Morpheme*</span>

Morphology is the study of the way words are built up from smaller meaning-bearing units, **morphemes**. A morpheme is often defined as the minimal meaning-bearing unit in a language. So for example the word *fox* consists of a single morpheme (the morpheme *fox*) while the word *cats* consists of two: the morpheme *cat* and the morpheme *-s*.

<span style="float:left">*Stem*<br>*Affix*</span>

As this example suggests, it is often useful to distinguish two broad classes of morphemes: **stems** and **affixes**. The exact details of the distinction vary from language to language, but intuitively, the stem is the "main" morpheme of the word, supplying the main meaning, while the affixes add "additional" meanings of various kinds.

Affixes are further divided into **prefixes**, **suffixes**, **infixes**, and **circumfixes**. Prefixes precede the stem, suffixes follow the stem, circumfixes do both, and infixes are inserted inside the stem. For example, the word *eats* is composed of a stem *eat* and the suffix *-s*. The word *unbuckle* is composed of a stem *buckle* and the prefix *un-*. English doesn't have any good examples of circumfixes, but many other languages do. In German, for example, the past participle of some verbs is formed by adding *ge-* to the beginning of the stem and *-t* to the end; so the past participle of the verb *sagen* (to say) is *gesagt* (said). Infixes, in which a morpheme is inserted in the middle of a word, occur very commonly for example in the Philipine language Tagalog. For example the affix *um*, which marks the agent of an action, is infixed to the Tagalog stem *hingi* "borrow" to produce *humingi*. There is one infix that occurs in some dialects of English in which the taboo morphemes "f**king" or "bl**dy" or others like them are inserted in the middle of other words ("Man-f**king-hattan", "abso-bl**dy-lutely"[3]) (McCawley, 1978).

A word can have more than one affix. For example, the word *rewrites* has the prefix

---

[2]    Lemmatization is actually more complex, since it sometimes involves deciding on which sense of a word is present. We return to this issue in Ch. 20.

[3]    Alan Jay Lerner, the lyricist of My Fair Lady, bowdlerized the latter to *abso-bloomin'lutely* in the lyric to "Wouldn't It Be Loverly?" (Lerner, 1978, p. 60).

*re-*, the stem *write*, and the suffix *-s*. The word *unbelievably* has a stem (*believe*) plus three affixes (*un-*, *-able*, and *-ly*). While English doesn't tend to stack more than four or five affixes, languages like Turkish can have words with nine or ten affixes, as we saw above. Languages that tend to string affixes together like Turkish does are called **agglutinative** languages.

There are many ways to combine morphemes to create words. Four of these methods are common and play important roles in speech and language processing: **inflec-tion**, **derivation**, **compounding**, and **cliticization**.

*Inflection*

*Derivation*

*Compounding*

*Cliticization*

**Inflection** is the combination of a word stem with a grammatical morpheme, usu-ally resulting in a word of the same class as the original stem, and usually filling some syntactic function like agreement. For example, English has the inflectional morpheme *-s* for marking the **plural** on nouns, and the inflectional morpheme *-ed* for marking the past tense on verbs. **Derivation** is the combination of a word stem with a grammatical morpheme, usually resulting in a word of a *different* class, often with a meaning hard to predict exactly. For example the verb *computerize* can take the derivational suffix *-ation* to produce the noun *computerization*. **Compounding** is the combination of mul-tiple word stems together. For example the noun *doghouse* is the concatenation of the morpheme *dog* with the morpheme *house*. Finally, **cliticization** is the combination of

*Clitic*

a word stem with a **clitic**. A clitic is a morpheme that acts syntactically like a word, but is reduced in form and attached (phonologically and sometimes orthographically) to another word. For example the English morpheme *'ve* in the word *I've* is a clitic, as is the French definite article *l'* in the word *l'opera*. In the following sections we give more details on these processes.

### 3.1.1   Inflectional Morphology

English has a relatively simple inflectional system; only nouns, verbs, and sometimes adjectives can be inflected, and the number of possible inflectional affixes is quite small.

*Plural*

English nouns have only two kinds of inflection: an affix that marks **plural** and an affix that marks **possessive**. For example, many (but not all) English nouns can either

*Singular*

appear in the bare stem or **singular** form, or take a plural suffix. Here are examples of the regular plural suffix *-s* (also spelled *-es*), and irregular plurals:

|              | Regular Nouns |          | Irregular Nouns |      |
| ------------ | ------------- | -------- | --------------- | ---- |
| **Singular** | cat           | thrush   | mouse           | ox   |
| **Plural**   | cats          | thrushes | mice            | oxen |

While the regular plural is spelled *-s* after most nouns, it is spelled *-es* after words ending in *-s* (*ibis/ibises*), *-z* (*waltz/waltzes*), *-sh* (*thrush/thrushes*), *-ch* (*finch/finches*), and sometimes *-x* (*box/boxes*). Nouns ending in *-y* preceded by a consonant change the *-y* to *-i* (*butterfly/butterflies*).

The possessive suffix is realized by apostrophe + *-s* for regular singular nouns (*llama's*) and plural nouns not ending in *-s* (*children's*) and often by a lone apostro-phe after regular plural nouns (*llamas'*) and some names ending in *-s* or *-z* (*Euripides' comedies*).

English verbal inflection is more complicated than nominal inflection. First, English has three kinds of verbs; **main verbs**, (*eat, sleep, impeach*), **modal verbs** (*can, will, should*), and **primary verbs** (*be, have, do*) (using the terms of Quirk et al., 1985). In this chapter we will mostly be concerned with the main and primary verbs, because

*Regular verb*

it is these that have inflectional endings. Of these verbs a large class are **regular**, that is to say all verbs of this class have the same endings marking the same functions. These regular verbs (e.g. *walk*, or *inspect*) have four morphological forms, as follow:

| Morphological Class | Regularly Inflected Verbs | | | |
|---|---|---|---|---|
| stem | walk | merge | try | map |
| *-s* form | walks | merges | tries | maps |
| *-ing* participle | walking | merging | trying | mapping |
| Past form or *-ed* participle | walked | merged | tried | mapped |

These verbs are called regular because just by knowing the stem we can predict the other forms by adding one of three predictable endings and making some regular spelling changes (and as we will see in Ch. 7, regular pronunciation changes). These regular verbs and forms are significant in the morphology of English first because they cover a majority of the verbs, and second because the regular class is **productive**. As discussed earlier, a productive class is one that automatically includes any new words that enter the language. For example the recently-created verb *fax* (*My mom faxed me the note from cousin Everett*) takes the regular endings *-ed*, *-ing*, *-es*. (Note that the *-s* form is spelled *faxes* rather than *faxs*; we will discuss spelling rules below).

*Irregular verb*

The **irregular verbs** are those that have some more or less idiosyncratic forms of inflection. Irregular verbs in English often have five different forms, but can have as many as eight (e.g., the verb *be*) or as few as three (e.g. *cut* or *hit*). While constituting a much smaller class of verbs (Quirk et al. (1985) estimate there are only about 250 irregular verbs, not counting auxiliaries), this class includes most of the very frequent verbs of the language.[4] The table below shows some sample irregular forms. Note that

*Preterite*

an irregular verb can inflect in the past form (also called the **preterite**) by changing its vowel (*eat/ate*), or its vowel and some consonants (*catch/caught*), or with no change at all (*cut/cut*).

| Morphological Class | Irregularly Inflected Verbs | | |
|---|---|---|---|
| stem | eat | catch | cut |
| *-s* form | eats | catches | cuts |
| *-ing* participle | eating | catching | cutting |
| preterite | ate | caught | cut |
| past participle | eaten | caught | cut |

The way these forms are used in a sentence will be discussed in the syntax and semantics chapters but is worth a brief mention here. The *-s* form is used in the "habitual present" form to distinguish the third-person singular ending (*She jogs every Tuesday*)

---

[4]  In general, the more frequent a word form, the more likely it is to have idiosyncratic properties; this is due to a fact about language change; very frequent words tend to preserve their form even if other words around them are changing so as to become more regular.

from the other choices of person and number (*I/you/we/they jog every Tuesday*). The stem form is used in the infinitive form, and also after certain other verbs (*I'd rather walk home*, *I want to walk home*). The *-ing* participle is used in the **progressive** construction to mark present or ongoing activity (*It is raining*), or when the verb is treated as a noun; this particular kind of nominal use of a verb is called a **gerund** use: *Fishing is fine if you live near water.* The *-ed/-en* participle is used in the **perfect** construction (*He's eaten lunch already*) or the passive construction (*The verdict was overturned yesterday*).

*Progressive*

*Gerund*

*Perfect*

In addition to noting which suffixes can be attached to which stems, we need to capture the fact that a number of regular spelling changes occur at these morpheme boundaries. For example, a single consonant letter is doubled before adding the *-ing* and *-ed* suffixes (*beg/begging/begged*). If the final letter is "c", the doubling is spelled "ck" (*picnic/picnicking/picnicked*). If the base ends in a silent *-e*, it is deleted before adding *-ing* and *-ed* (*merge/merging/merged*). Just as for nouns, the *-s* ending is spelled *-es* after verb stems ending in *-s* (*toss/tosses*) , *-z*, (*waltz/waltzes*) *-sh*, (*wash/washes*) *-ch*, (*catch/catches*) and sometimes *-x* (*tax/taxes*). Also like nouns, verbs ending in *-y* preceded by a consonant change the *-y* to *-i* (*try/tries*).

The English verbal system is much simpler than for example the European Spanish system, which has as many as fifty distinct verb forms for each regular verb. Fig. 3.1 shows just a few of the examples for the verb *amar*, 'to love'. Other languages can have even more forms than this Spanish example.

|      | Present Indicative | Imperfect Indicative | Future | Preterite | Present Subjunctive | Conditional | Imperfect Subjunctive | Future Subjunctive |
|------|-------------------|---------------------|--------|-----------|--------------------|-------------|----------------------|-------------------|
| 1SG  | amo    | amaba    | amaré    | amé      | ame   | amaría   | amara    | amare    |
| 2SG  | amas   | amabas   | amarás   | amaste   | ames  | amarías  | amaras   | amares   |
| 3SG  | ama    | amaba    | amará    | amó      | ame   | amaría   | amara    | amáreme  |
| 1PL  | amamos | amábamos | amaremos | amamos   | amemos | amaríamos | amáramos | amáremos |
| 2PL  | amáis  | amabais  | amaréis  | amasteis | améis | amaríais | amarais  | amareis  |
| 3PL  | aman   | amaban   | amarán   | amaron   | amen  | amarían  | amaran   | amaren   |

**Figure 3.1**    To love in Spanish. Some of the inflected forms of the verb *amar* in European Spanish. *1SG* stands for "first person singular", 3PL for "third person plural", and so on.

## 3.1.2   Derivational Morphology

While English inflection is relatively simple compared to other languages, derivation in English is quite complex. Recall that derivation is the combination of a word stem with a grammatical morpheme, usually resulting in a word of a *different* class, often with a meaning hard to predict exactly.

*nominalization*

A very common kind of derivation in English is the formation of new nouns, often from verbs or adjectives. This process is called **nominalization**. For example, the suffix *-ation* produces nouns from verbs ending often in the suffix *-ize* (*computerize → computerization*). Here are examples of some particularly productive English nominalizing suffixes.

| Suffix | Base Verb/Adjective | Derived Noun |
|--------|---------------------|--------------|
| -ation | computerize (V) | computerization |
| -ee | appoint (V) | appointee |
| -er | kill (V) | killer |
| -ness | fuzzy (A) | fuzziness |

Adjectives can also be derived from nouns and verbs. Here are examples of a few suffixes deriving adjectives from nouns or verbs.

| Suffix | Base Noun/Verb | Derived Adjective |
|--------|----------------|-------------------|
| -al | computation (N) | computational |
| -able | embrace (V) | embraceable |
| -less | clue (N) | clueless |

Derivation in English is more complex than inflection for a number of reasons. One is that it is generally less productive; even a nominalizing suffix like *-ation*, which can be added to almost any verb ending in *-ize*, cannot be added to absolutely every verb. Thus we can't say *\*eatation* or *\*spellation* (we use an asterisk (\*) to mark "non-examples" of English). Another is that there are subtle and complex meaning differences among nominalizing suffixes. For example *sincerity* has a subtle difference in meaning from *sincereness*.

### 3.1.3    Cliticization

Recall that a clitic is a unit whose status lies in between that of an affix and a word. The phonological behavior of clitics is like affixes; they tend to be short and unaccented (we will talk more about phonology in Ch. 8). Their syntactic behavior is more like words, often acting as pronouns, articles, conjunctions, or verbs. Clitics preceding a word are called **proclitics**, while those following are **enclitics**.

*Proclitic*
*Enclitic*

English clitics include these auxiliary verbal forms:

| Full Form | Clitic | Full Form | Clitic |
|-----------|--------|-----------|--------|
| am | 'm | have | 've |
| are | 're | has | 's |
| is | 's | had | 'd |
| will | 'll | would | 'd |

Note that the clitics in English are ambiguous; Thus *she's* can mean *she is* or *she has*. Except for a few such ambiguities, however, correctly segmenting off clitics in English is simplified by the presence of the apostrophe. Clitics can be harder to parse in other languages. In Arabic and Hebrew, for example, the definite article (*the*; *Al* in Arabic, *ha* in Hebrew) is cliticized on to the front of nouns. It must be segmented off in order to do part-of-speech tagging, parsing, or other tasks. Other Arabic proclitics include prepositions like *b* 'by/with', and conjunctions like *w* 'and'. Arabic also has *enclitics* marking certain pronouns. For example the word *and by their virtues* has clitics meaning *and*, *by*, and *their*, a stem *virtue*, and a plural affix. Note that since

Arabic is read right to left, these would actually appear ordered from right to left in an Arabic word.

|           | **proclitic** | **proclitic** | **stem** | **affix** | **enclitic** |
|-----------|---------------|---------------|----------|-----------|--------------|
| **Arabic** | w            | b             | Hsn      | At        | hm           |
| **Gloss**  | and          | by            | virtue   | s         | their        |

### 3.1.4   Non-concatenative Morphology

*Concatenative morphology*

The kind of morphology we have discussed so far, in which a word is composed of a string of morphemes concatenated together is often called **concatenative morphology**. A number of languages have extensive **non-concatenative morphology**, in which morphemes are combined in more complex ways. The Tagalog infixation example above is one example of non-concatenative morphology, since two morphemes (*hingi* and *um*) are intermingled.

Another kind of non-concatenative morphology is called **templatic morphology** or **root-and-pattern** morphology. This is very common in Arabic, Hebrew, and other Semitic languages. In Hebrew, for example, a verb (as well as other parts-of-speech) is constructed using two components: a root, consisting usually of three consonants (CCC) and carrying the basic meaning, and a template, which gives the ordering of consonants and vowels and specifies more semantic information about the resulting verb, such as the semantic voice (e.g., active, passive, middle). For example the Hebrew tri-consonantal root *lmd*, meaning 'learn' or 'study', can be combined with the active voice CaCaC template to produce the word *lamad*, 'he studied', or the intensive CiCeC template to produce the word *limed*, 'he taught', or the intensive passive template CuCaC to produce the word *lumad*, 'he was taught'. Arabic and Hebrew combine this templatic morphology with concatenative morphology (like the cliticization example shown in the previous section).

### 3.1.5   Agreement

*Agreement*

*Gender*

*Noun class*

We introduced the plural morpheme above, and noted that plural is marked on both nouns and verbs in English. We say that the subject noun and the main verb in English have to **agree** in number, meaning that the two must either be both singular or both plural. There are other kinds of agreement processes. For example nouns, adjectives, and sometimes verbs in many languages are marked for **gender**. A gender is a kind of equivalence class that is used by the language to categorize the nouns; each noun falls into one class. Many languages (for example Romance languages like French, Spanish, or Italian) have 2 genders, which are referred to as masculine and feminine. Other languages (like most Germanic and Slavic languages) have three (masculine, feminine, neuter). Some languages, for example the Bantu languages of Africa, have as many as 20 genders. When the number of classes is very large, we often refer to them as **noun classes** instead of genders.

Gender is sometimes marked explicitly on a noun; for example Spanish masculine words often end in *-o* and feminine words in *-a*. But in many cases the gender is not marked in the letters or phones of the noun itself. Instead, it is a property of the word