

CHAPTER

12

Machine Translation

“I want to talk the dialect of your people. It’s no use of talking unless people understand what you say.”

Zora Neale Hurston, *Moses, Man of the Mountain* 1939, p. 121

machine
translation
MT

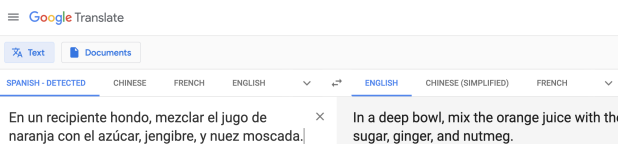
This chapter introduces **machine translation (MT)**, the use of computers to translate from one language to another.

Of course translation, in its full generality, such as the translation of literature, or poetry, is a difficult, fascinating, and intensely human endeavor, as rich as any other area of human creativity.

information
access

Machine translation in its present form therefore focuses on a number of very practical tasks. Perhaps the most common current use of machine translation is for **information access**. We might want to translate some instructions on the web, perhaps the recipe for a favorite dish, or the steps for putting together some furniture. Or we might want to read an article in a newspaper, or get information from an online resource like Wikipedia or a government webpage in some other language.

MT for information access is probably one of the most common uses of NLP technology, and Google



digital divide

Translate alone (shown above) translates hundreds of billions of words a day between over 100 languages. Improvements in machine translation can thus help reduce what is often called the **digital divide** in information access: the fact that much more information is available in English and other languages spoken in wealthy countries. Web searches in English return much more information than searches in other languages, and online resources like Wikipedia are much larger in English and other higher-resourced languages. High-quality translation can help provide information to speakers of lower-resourced languages.

post-editing

Another common use of machine translation is to aid human translators. MT systems are routinely used to produce a draft translation that is fixed up in a **post-editing** phase by a human translator. This task is often called **computer-aided translation** or **CAT**. CAT is commonly used as part of **localization**: the task of adapting content or a product to a particular language community.

CAT
localization

Finally, a more recent application of MT is to in-the-moment human communication needs. This includes incremental translation, translating speech on-the-fly before the entire sentence is complete, as is commonly used in simultaneous interpretation. Image-centric translation can be used for example to use OCR of the text on a phone camera image as input to an MT system to translate menus or street signs.

encoder-
decoder

The standard algorithm for MT is the **encoder-decoder** network/ We briefly mentioned in Chapter 7 that encoder-decoder or sequence-to-sequence models are used for tasks in which we need to map an input sequence to an output sequence that is a complex function of the entire input sequence, like machine translation or

speech recognition. Indeed, in machine translation, the words of the target language don't necessarily agree with the words of the source language in number or order. Consider translating the following made-up English sentence into Japanese.

- (12.1) English: *He wrote a letter to a friend*
 Japanese: *tomodachi ni tegami-o kaita*
 friend to letter wrote

Note that the elements of the sentences are in very different places in the different languages. In English, the verb is in the middle of the sentence, while in Japanese, the verb *kaita* comes at the end. The Japanese sentence doesn't require the pronoun *he*, while English does.

Such differences between languages can be quite complex. In the following actual sentence from the United Nations, notice the many changes between the Chinese sentence (we've given in red a word-by-word gloss of the Chinese characters) and its English equivalent produced by human translators.

- (12.2) 大会/General Assembly 在/on 1982年/1982 12月/December 10日/10 通过
 了/adopted 第37号/37th 决议/resolution , 核准了/approved 第二
 次/second 探索/exploration 及/and 和平/peaceful 利用/using 外层空
 间/outer space 会议/conference 的/of 各项/various 建议/suggestions 。

On 10 December 1982 , the General Assembly adopted resolution 37 in which it endorsed the recommendations of the Second United Nations Conference on the Exploration and Peaceful Uses of Outer Space .

Note the many ways the English and Chinese differ. For example the ordering differs in major ways; the Chinese order of the noun phrase is “peaceful using outer space conference of suggestions” while the English has “suggestions of the ... conference on peaceful use of outer space”). And the order differs in minor ways (the date is ordered differently). English requires *the* in many places that Chinese doesn't, and adds some details (like “in which” and “it”) that aren't necessary in Chinese. Chinese doesn't grammatically mark plurality on nouns (unlike English, which has the “-s” in “recommendations”), and so the Chinese must use the modifier 各项/*various* to make it clear that there is not just one recommendation. English capitalizes some words but not others. Encoder-decoder networks are very successful at handling these sorts of complicated cases of sequence mappings.

We'll begin in the next section by considering the linguistic background about how languages vary, and the implications this variance has for the task of MT. Then we'll sketch out the standard algorithm, give details about things like input tokenization and creating training corpora of parallel sentences, give some more low-level details about the encoder-decoder network, and finally discuss how MT is evaluated, introducing the simple chrF metric.

12.1 Language Divergences and Typology

universal

There are about 7,000 languages in the world. Some aspects of human language seem to be **universal**, holding true for every one of these languages, or are statistical universals, holding true for most of these languages. Many universals arise from the functional role of language as a communicative system by humans. Every language, for example, seems to have words for referring to people, for talking about eating and drinking, for being polite or not. There are also structural linguistic universals; for example, every language seems to have nouns and verbs (Chapter 17), has

ways to ask questions, or issue commands, has linguistic mechanisms for indicating agreement or disagreement.

translation
divergence

Yet languages also **differ** in many ways (as has been pointed out since ancient times; see Fig. 12.1). Understanding what causes such **translation divergences** (Dorr, 1994) can help us build better MT models. We often distinguish the **idiosyncratic** and lexical differences that must be dealt with one by one (the word for “dog” differs wildly from language to language), from **systematic** differences that we can model in a general way (many languages put the verb before the grammatical object; others put the verb after the grammatical object). The study of these systematic cross-linguistic similarities and differences is called **linguistic typology**. This section sketches some typological facts that impact machine translation; the interested reader should also look into WALS, the World Atlas of Language Structures, which gives many typological facts about languages (Dryer and Haspelmath, 2013).

typology



Figure 12.1 The Tower of Babel, Pieter Bruegel 1563. Wikimedia Commons, from the Kunsthistorisches Museum, Vienna.

12.1.1 Word Order Typology

SVO

SOV

VSO

As we hinted at in our example above comparing English and Japanese, languages differ in the basic word order of verbs, subjects, and objects in simple declarative clauses. German, French, English, and Mandarin, for example, are all **SVO** (**Subject-Verb-Object**) languages, meaning that the verb tends to come between the subject and object. Hindi and Japanese, by contrast, are **SOV** languages, meaning that the verb tends to come at the end of basic clauses, and Irish and Arabic are **VSO** languages. Two languages that share their basic word order type often have other similarities. For example, **VO** languages generally have **prepositions**, whereas **OV** languages generally have **postpositions**.

Let’s look in more detail at the example we saw above. In this SVO English sentence, the verb *wrote* is followed by its object *a letter* and the prepositional phrase

to a friend, in which the preposition *to* is followed by its argument *a friend*. Arabic, with a VSO order, also has the verb before the object and prepositions. By contrast, in the Japanese example that follows, each of these orderings is reversed; the verb is *preceded* by its arguments, and the postposition follows its argument.

- (12.3) English: *He wrote a letter to a friend*
 Japanese: *tomodachi ni tegami-o kaita*
 friend to letter wrote
 Arabic: *katabt risāla li šadq*
 wrote letter to friend

Other kinds of ordering preferences vary idiosyncratically from language to language. In some SVO languages (like English and Mandarin) adjectives tend to appear before nouns, while in others languages like Spanish and Modern Hebrew, adjectives appear after the noun:

- (12.4) Spanish *bruja verde* English *green witch*

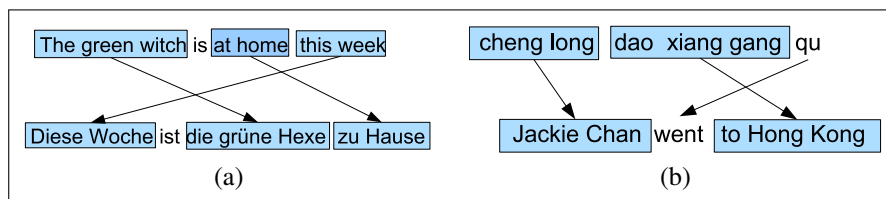


Figure 12.2 Examples of other word order differences: (a) In German, adverbs occur in initial position that in English are more natural later, and tensed verbs occur in second position. (b) In Mandarin, preposition phrases expressing goals often occur pre-verbally, unlike in English.

Fig. 12.2 shows examples of other word order differences. All of these word order differences between languages can cause problems for translation, requiring the system to do huge structural reorderings as it generates the output.

12.1.2 Lexical Divergences

Of course we also need to translate the individual words from one language to another. For any translation, the appropriate word can vary depending on the context. The English source-language word *bass*, for example, can appear in Spanish as the fish *lubina* or the musical instrument *bajo*. German uses two distinct words for what in English would be called a *wall*: *Wand* for walls inside a building, and *Mauer* for walls outside a building. Where English uses the word *brother* for any male sibling, Chinese and many other languages have distinct words for *older brother* and *younger brother* (Mandarin *gege* and *didi*, respectively). In all these cases, translating *bass*, *wall*, or *brother* from English would require a kind of specialization, disambiguating the different uses of a word. For this reason the fields of MT and Word Sense Disambiguation (Appendix G) are closely linked.

Sometimes one language places more grammatical constraints on word choice than another. We saw above that English marks nouns for whether they are singular or plural. Mandarin doesn't. Or French and Spanish, for example, mark grammatical gender on adjectives, so an English translation into French requires specifying adjective gender.

The way that languages differ in lexically dividing up conceptual space may be more complex than this one-to-many translation problem, leading to many-to-many

mappings. For example, Fig. 12.3 summarizes some of the complexities discussed by Hutchins and Somers (1992) in translating English *leg*, *foot*, and *paw*, to French. For example, when *leg* is used about an animal it's translated as French *patte*; but about the leg of a journey, as French *etape*; if the leg is of a chair, we use French *pied*.

lexical gap

Further, one language may have a **lexical gap**, where no word or phrase, short of an explanatory footnote, can express the exact meaning of a word in the other language. For example, English does not have a word that corresponds neatly to Mandarin *xiào* or Japanese *oyakōkō* (in English one has to make do with awkward phrases like *filial piety* or *loving child*, or *good son/daughter* for both).

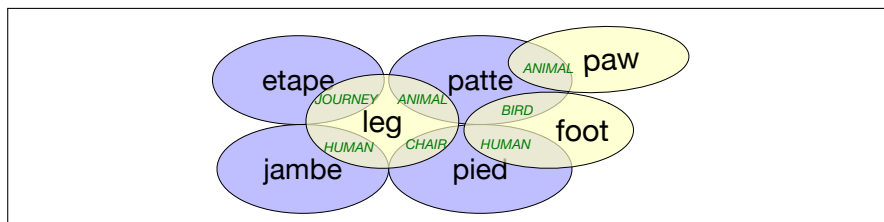


Figure 12.3 The complex overlap between English *leg*, *foot*, etc., and various French translations as discussed by Hutchins and Somers (1992).

Finally, languages differ systematically in how the conceptual properties of an event are mapped onto specific words. Talmy (1985, 1991) noted that languages can be characterized by whether direction of motion and manner of motion are marked on the verb or on the “satellites”: particles, prepositional phrases, or adverbial phrases. For example, a bottle floating out of a cave would be described in English with the direction marked on the particle *out*, while in Spanish the direction would be marked on the verb:

(12.5) English: *The bottle floated out.*

Spanish: *La botella salió flotando.*

The bottle exited floating.

verb-framed

satellite-framed

Verb-framed languages mark the direction of motion on the verb (leaving the satellites to mark the manner of motion), like Spanish *acercarse* ‘approach’, *alcanzar* ‘reach’, *entrar* ‘enter’, *salir* ‘exit’. **Satellite-framed** languages mark the direction of motion on the satellite (leaving the verb to mark the manner of motion), like English *crawl out*, *float off*, *jump down*, *run after*. Languages like Japanese, Tamil, and the many languages in the Romance, Semitic, and Mayan languages families, are verb-framed; Chinese as well as non-Romance Indo-European languages like English, Swedish, Russian, Hindi, and Farsi are satellite framed (Talmy 1991, Slobin 1996).

12.1.3 Morphological Typology

isolating

polysynthetic

agglutinative

fusion

Morphologically, languages are often characterized along two dimensions of variation. The first is the number of morphemes per word, ranging from **isolating** languages like Vietnamese and Cantonese, in which each word generally has one morpheme, to **polysynthetic** languages like Siberian Yupik (“Eskimo”), in which a single word may have very many morphemes, corresponding to a whole sentence in English. The second dimension is the degree to which morphemes are segmentable, ranging from **agglutinative** languages like Turkish, in which morphemes have relatively clean boundaries, to **fusion** languages like Russian, in which a single affix

may conflate multiple morphemes, like *-om* in the word *stolom* (table-SG-INSTR-DECL1), which fuses the distinct morphological categories instrumental, singular, and first declension.

Translating between languages with rich morphology requires dealing with structure below the word level, and for this reason modern systems generally use subword models like the wordpiece or BPE models of Section 12.2.1.

12.1.4 Referential density

Finally, languages vary along a typological dimension related to the things they tend to omit. Some languages, like English, require that we use an explicit pronoun when talking about a referent that is given in the discourse. In other languages, however, we can sometimes omit pronouns altogether, as the following example from Spanish shows¹:

(12.6) [El jefe]_i dio con un libro. \emptyset_i Mostró su hallazgo a un descifrador ambulante.
[The boss] came upon a book. [He] showed his find to a wandering decoder.

pro-drop

referential
density

cold language

hot language

Languages that can omit pronouns are called **pro-drop** languages. Even among the pro-drop languages, there are marked differences in frequencies of omission. Japanese and Chinese, for example, tend to omit far more than does Spanish. This dimension of variation across languages is called the dimension of **referential density**. We say that languages that tend to use more pronouns are more **referentially dense** than those that use more zeros. Referentially sparse languages, like Chinese or Japanese, that require the hearer to do more inferential work to recover antecedents are also called **cold** languages. Languages that are more explicit and make it easier for the hearer are called **hot** languages. The terms *hot* and *cold* are borrowed from Marshall McLuhan's 1964 distinction between hot media like movies, which fill in many details for the viewer, versus cold media like comics, which require the reader to do more inferential work to fill out the representation (Bickel, 2003).

Translating from languages with extensive pro-drop, like Chinese or Japanese, to non-pro-drop languages like English can be difficult since the model must somehow identify each zero and recover who or what is being talked about in order to insert the proper pronoun.

12.2 Machine Translation using Encoder-Decoder

The standard architecture for MT is the **encoder-decoder transformer** or **sequence-to-sequence** model, an architecture we saw for RNNs in Chapter 13. We'll see the details of how to apply this architecture to transformers in Section 12.3, but first let's talk about the overall task.

Most machine translation tasks make the simplification that we can translate each sentence independently, so we'll just consider individual sentences for now. Given a sentence in a **source** language, the MT task is then to generate a corresponding sentence in a **target** language. For example, an MT system is given an English sentence like

The green witch arrived

and must translate it into the Spanish sentence:

¹ Here we use the \emptyset -notation; we'll introduce this and discuss this issue further in Chapter 23

Llegó la bruja verde

MT uses supervised machine learning: at training time the system is given a large set of **parallel** sentences (each sentence in a source language matched with a sentence in the target language), and learns to map source sentences into target sentences. In practice, rather than using words (as in the example above), we split the sentences into a sequence of subword tokens (tokens can be words, or subwords, or individual characters). The systems are then trained to maximize the probability of the sequence of tokens in the target language y_1, \dots, y_m given the sequence of tokens in the source language x_1, \dots, x_n :

$$P(y_1, \dots, y_m | x_1, \dots, x_n) \quad (12.7)$$

Rather than use the input tokens directly, the encoder-decoder architecture consists of two components, an **encoder** and a **decoder**. The encoder takes the input words $x = [x_1, \dots, x_n]$ and produces an intermediate context \mathbf{h} . At decoding time, the system takes \mathbf{h} and, word by word, generates the output y :

$$\mathbf{h} = \text{encoder}(x) \quad (12.8)$$

$$y_{t+1} = \text{decoder}(\mathbf{h}, y_1, \dots, y_t) \quad \forall t \in [1, \dots, m] \quad (12.9)$$

In the next two sections we'll talk about subword tokenization, and then how to get parallel corpora for training, and then we'll introduce the details of the encoder-decoder architecture.

12.2.1 Tokenization

Machine translation systems use a vocabulary that is fixed in advance, and rather than using space-separated words, this vocabulary is generated with subword tokenization algorithms, like the **BPE** algorithm sketched in Chapter 2. A shared vocabulary is used for the source and target languages, which makes it easy to copy tokens (like names) from source to target. Using subword tokenization with tokens shared between languages makes it natural to translate between languages like English or Hindi that use spaces to separate words, and languages like Chinese or Thai that don't.

We build the vocabulary by running a subword tokenization algorithm on a corpus that contains both source and target language data.

Rather than the simple BPE algorithm from Fig. ??, modern systems often use more powerful tokenization algorithms. Some systems (like BERT) use a variant of BPE called the **wordpiece** algorithm, which instead of choosing the most frequent set of tokens to merge, chooses merges based on which one most increases the language model probability of the tokenization. Wordpieces use a special symbol at the beginning of each token; here's a resulting tokenization from the Google MT system (Wu et al., 2016):

words: Jet makers feud over seat width with big orders at stake
wordpieces: _J et _makers _fe ud _over _seat _width _with _big _orders _at _stake

The wordpiece algorithm is given a training corpus and a desired vocabulary size V , and proceeds as follows:

1. Initialize the wordpiece lexicon with characters (for example a subset of Unicode characters, collapsing all the remaining characters to a special unknown character token).

2. Repeat until there are V wordpieces:

- (a) Train an n -gram language model on the training corpus, using the current set of wordpieces.
- (b) Consider the set of possible new wordpieces made by concatenating two wordpieces from the current lexicon. Choose the one new wordpiece that most increases the language model probability of the training corpus.

Recall that with BPE we had to specify the number of merges to perform; in wordpiece, by contrast, we specify the total vocabulary, which is a more intuitive parameter. A vocabulary of 8K to 32K word pieces is commonly used.

unigram
SentencePiece

An even more commonly used tokenization algorithm is (somewhat ambiguously) called the **unigram** algorithm (Kudo, 2018) or sometimes the **SentencePiece** algorithm, and is used in systems like ALBERT (Lan et al., 2020) and T5 (Raffel et al., 2020). (Because unigram is the default tokenization algorithm used in a library called SentencePiece that adds a useful wrapper around tokenization algorithms (Kudo and Richardson, 2018), authors often say they are using SentencePiece tokenization but really mean they are using the **unigram** algorithm).

In unigram tokenization, instead of building up a vocabulary by merging tokens, we start with a huge vocabulary of every individual unicode character plus all frequent sequences of characters (including all space-separated words, for languages with spaces), and iteratively remove some tokens to get to a desired final vocabulary size. The algorithm is complex (involving suffix-trees for efficiently storing many tokens, and the EM algorithm for iteratively assigning probabilities to tokens), so we don't give it here, but see Kudo (2018) and Kudo and Richardson (2018). Roughly speaking the algorithm proceeds iteratively by estimating the probability of each token, tokenizing the input data using various tokenizations, then removing a percentage of tokens that don't occur in high-probability tokenization, and then iterates until the vocabulary has been reduced down to the desired number of tokens.

Why does unigram tokenization work better than BPE? BPE tends to create lots of very small non-meaningful tokens (because BPE can only create larger words or morphemes by merging characters one at a time), and it also tends to merge very common tokens, like the suffix *ed*, onto their neighbors. We can see from these examples from Bostrom and Durrett (2020) that unigram tends to produce tokens that are more semantically meaningful:

Original:	corrupted	Original:	Completely preposterous suggestions
BPE:	corrupted	BPE:	Completely preposterous suggestions
Unigram:	corrupt ed	Unigram:	Complete ly pre post er ous suggestion s

12.2.2 Creating the Training data

parallel corpus

Europarl

Machine translation models are trained on a **parallel corpus**, sometimes called a **bitext**, a text that appears in two (or more) languages. Large numbers of parallel corpora are available. Some are governmental; the **Europarl** corpus (Koehn, 2005), extracted from the proceedings of the European Parliament, contains between 400,000 and 2 million sentences each from 21 European languages. The United Nations Parallel Corpus contains on the order of 10 million sentences in the six official languages of the United Nations (Arabic, Chinese, English, French, Russian, Spanish) Ziemski et al. (2016). Other parallel corpora have been made from movie and TV subtitles, like the **OpenSubtitles** corpus (Lison and Tiedemann, 2016), or from general web text, like the **ParaCrawl** corpus of 223 million sentence pairs between 23 EU languages and English extracted from the CommonCrawl Bañón et al. (2020).

Sentence alignment

Standard training corpora for MT come as aligned pairs of sentences. When creating new corpora, for example for underresourced languages or new domains, these sentence alignments must be created. Fig. 12.4 gives a sample hypothetical sentence alignment.

E1: "Good morning," said the little prince.	F1: -Bonjour, dit le petit prince.
E2: "Good morning," said the merchant.	F2: -Bonjour, dit le marchand de pilules perfectionnées qui apaisent la soif.
E3: This was a merchant who sold pills that had been perfected to quench thirst.	F3: On en avale une par semaine et l'on n'éprouve plus le besoin de boire.
E4: You just swallow one pill a week and you won't feel the need for anything to drink.	F4: -C'est une grosse économie de temps, dit le marchand.
E5: "They save a huge amount of time," said the merchant.	F5: Les experts ont fait des calculs.
E6: "Fifty-three minutes a week."	F6: On épargne cinquante-trois minutes par semaine.
E7: "If I had fifty-three minutes to spend?" said the little prince to himself.	F7: "Moi, se dit le petit prince, si j'avais cinquante-trois minutes à dépenser, je marcherais tout doucement vers une fontaine..."
E8: "I would take a stroll to a spring of fresh water"	

Figure 12.4 A sample alignment between sentences in English and French, with sentences extracted from Antoine de Saint-Exupéry's *Le Petit Prince* and a hypothetical translation. Sentence alignment takes sentences e_1, \dots, e_n , and f_1, \dots, f_m and finds minimal sets of sentences that are translations of each other, including single sentence mappings like (e_1, f_1) , (e_4, f_3) , (e_5, f_4) , (e_6, f_6) as well as 2-1 alignments $(e_2/e_3, f_2)$, $(e_7/e_8, f_7)$, and null alignments (f_5) .

Given two documents that are translations of each other, we generally need two steps to produce sentence alignments:

- a cost function that takes a span of source sentences and a span of target sentences and returns a score measuring how likely these spans are to be translations.
- an alignment algorithm that takes these scores to find a good alignment between the documents.

To score the similarity of sentences across languages, we need to make use of a **multilingual embedding space**, in which sentences from different languages are in the same embedding space (Artetxe and Schwenk, 2019). Given such a space, cosine similarity of such embeddings provides a natural scoring function (Schwenk, 2018). Thompson and Koehn (2019) give the following cost function between two sentences or spans x, y from the source and target documents respectively:

$$c(x, y) = \frac{(1 - \cos(x, y))nSents(x) nSents(y)}{\sum_{s=1}^S 1 - \cos(x, y_s) + \sum_{s=1}^S 1 - \cos(x_s, y)} \quad (12.10)$$

where $nSents()$ gives the number of sentences (this biases the metric toward many alignments of single sentences instead of aligning very large spans). The denominator helps to normalize the similarities, and so $x_1, \dots, x_S, y_1, \dots, y_S$, are randomly selected sentences sampled from the respective documents.

Usually dynamic programming is used as the alignment algorithm (Gale and Church, 1993), in a simple extension of the minimum edit distance algorithm we introduced in Chapter 2.

Finally, it's helpful to do some corpus cleanup by removing noisy sentence pairs. This can involve handwritten rules to remove low-precision pairs (for example removing sentences that are too long, too short, have different URLs, or even pairs

that are too similar, suggesting that they were copies rather than translations). Or pairs can be ranked by their multilingual embedding cosine score and low-scoring pairs discarded.

12.3 Details of the Encoder-Decoder Model

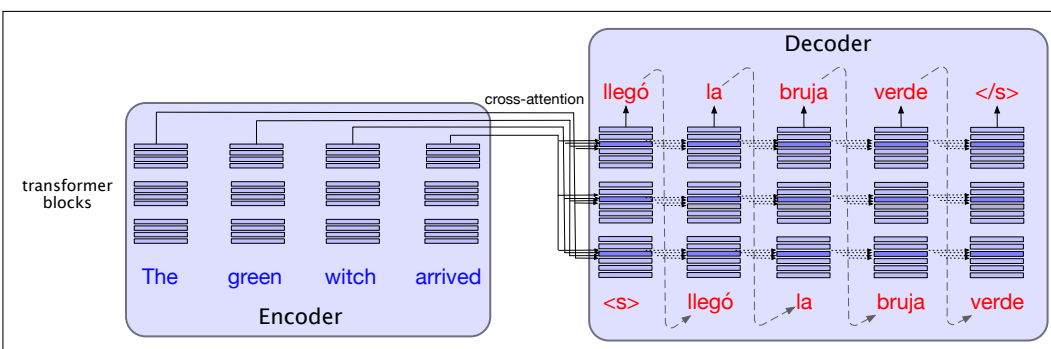


Figure 12.5 The encoder-decoder transformer architecture for machine translation. The encoder uses the transformer blocks we saw in Chapter 8, while the decoder uses a more powerful block with an extra **cross-attention** layer that can attend to all the encoder words. We’ll see this in more detail in the next section.

The standard architecture for MT is the encoder-decoder transformer. (For those of you who studied RNNs, the encoder-decoder architecture was introduced already for RNNs in Chapter 13.) Fig. 12.5 shows the intuition of the architecture at a high level. You’ll see that the encoder-decoder architecture is made up of two transformers: an **encoder**, which is the same as the basic transformers from Chapter 8, and a **decoder**, which is augmented with a special new layer called the **cross-attention** layer. The encoder takes the source language input word tokens $\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_n$ and maps them to an output representation $\mathbf{H}^{enc} = \mathbf{h}_1, \dots, \mathbf{h}_n$ via a stack of encoder blocks.

The decoder is essentially a conditional language model that attends to the encoder representation and generates the target words one by one, at each timestep conditioning on the source sentence and the previously generated target language words to generate a token. Decoding can use any of the decoding methods discussed in Chapter 8 like greedy, or temperature or nucleus sampling. But the most common decoding algorithm for MT is the beam search algorithm that we’ll introduce in Section 12.4.

But the components of the architecture differ somewhat from the transformer block we’ve seen. First, in order to attend to the source language, the transformer blocks in the decoder have an extra **cross-attention** layer. Recall that the transformer block of Chapter 8 consists of a self-attention layer that attends to the input from the previous layer, preceded by layer norm, and followed by another layer norm and the feed forward layer. The decoder transformer block includes an extra layer with a special kind of attention, **cross-attention** (also sometimes called **encoder-decoder attention** or **source attention**). Cross-attention has the same form as the multi-head attention in a normal transformer block, except that while the queries as usual come from the previous layer of the decoder, the keys and values come from the output of the *encoder*.

cross-attention

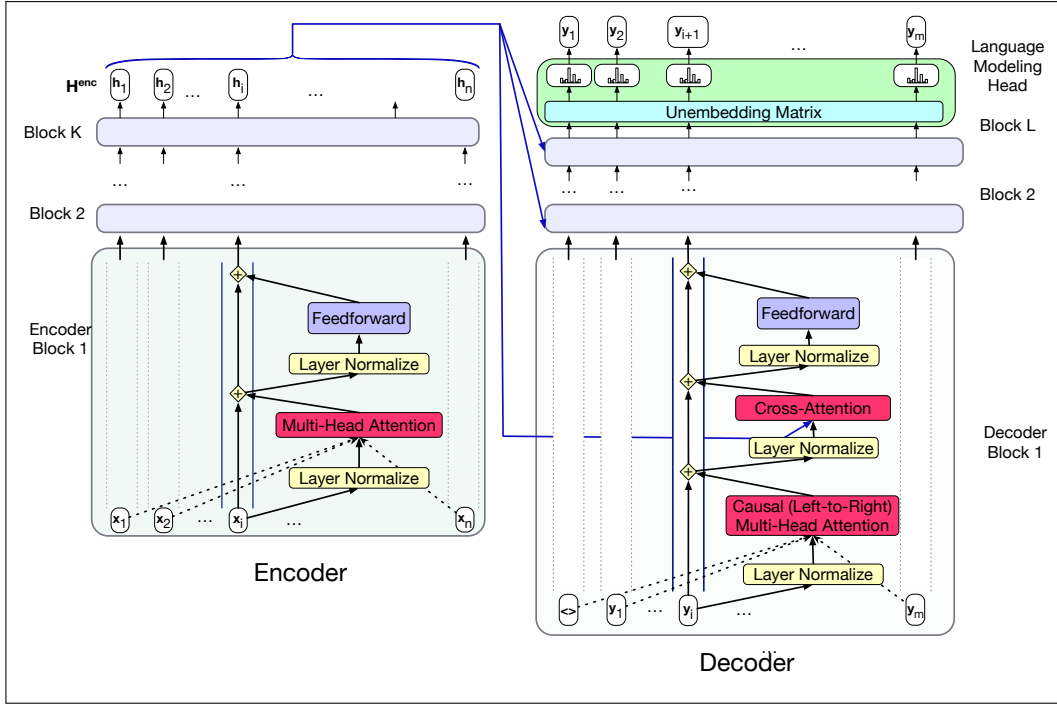


Figure 12.6 The transformer block for the encoder and the decoder, showing the residual stream view. The final output of the encoder $\mathbf{H}^{enc} = \mathbf{h}_1, \dots, \mathbf{h}_n$ is the context used in the decoder. The decoder is a standard transformer except with one extra layer, the **cross-attention** layer, which takes that encoder output \mathbf{H}^{enc} and uses it to form its **K** and **V** inputs.

That is, where in standard multi-head attention the input to each attention layer is \mathbf{X} , in cross attention the input is the the final output of the encoder $\mathbf{H}^{enc} = \mathbf{h}_1, \dots, \mathbf{h}_n$. \mathbf{H}^{enc} is of shape $[n \times d]$, each row representing one input token. To link the keys and values from the encoder with the query from the prior layer of the decoder, we multiply the encoder output \mathbf{H}^{enc} by the cross-attention layer's key weights \mathbf{W}^K and value weights \mathbf{W}^V . The query comes from the output from the prior decoder layer $\mathbf{H}^{dec[\ell-1]}$, which is multiplied by the cross-attention layer's query weights \mathbf{W}^Q :

$$\mathbf{Q} = \mathbf{H}^{dec[\ell-1]} \mathbf{W}^Q; \mathbf{K} = \mathbf{H}^{enc} \mathbf{W}^K; \mathbf{V} = \mathbf{H}^{enc} \mathbf{W}^V \quad (12.11)$$

$$\text{CrossAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q} \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V} \quad (12.12)$$

The cross attention thus allows the decoder to attend to each of the source language words as projected into the entire encoder final output representations. The other attention layer in each decoder block, the multi-head attention layer, is the same causal (left-to-right) attention that we saw in Chapter 8. The multi-head attention in the encoder, however, is allowed to look ahead at the entire source language text, so it is not masked.

To train an encoder-decoder model, we use the same self-supervision model we used for training encoder-decoders RNNs in Chapter 13. The network is given the source text and then starting with the separator token is trained autoregressively to predict the next token using cross-entropy loss. Recall that cross-entropy loss for language modeling is determined by the probability the model assigns to the correct

next word. So at time t the CE loss is the negative log probability the model assigns to the next word in the training sequence:

$$L_{CE}(\hat{\mathbf{y}}_t, \mathbf{y}_t) = -\log \hat{\mathbf{y}}_t[w_{t+1}] \quad (12.13)$$

teacher forcing As in that case, we use **teacher forcing** in the decoder. Recall that in teacher forcing, at each time step in decoding we force the system to use the gold target token from training as the next input x_{t+1} , rather than allowing it to rely on the (possibly erroneous) decoder output \hat{y}_t .

12.4 Decoding in MT: Beam Search

Recall the **greedy decoding** algorithm from Chapter 8: at each time step t in generation, the output y_t is chosen by computing the probability for each word in the vocabulary and then choosing the highest probability word (the argmax):

$$\hat{w}_t = \operatorname{argmax}_{w \in V} P(w | \mathbf{w}_{<t}) \quad (12.14)$$

A problem with greedy decoding is that what looks high probability at word t might turn out to have been the wrong choice once we get to word $t + 1$. The **beam search** algorithm maintains multiple choices until later when we can see which one is best.

In beam search we model decoding as searching the space of possible generations, represented as a **search tree** whose **branches** represent actions (generating a token), and **nodes** represent states (having generated a particular prefix). We search for the best action sequence, i.e., the string with the highest probability.

An illustration of the problem

Fig. 12.7 shows a made-up example. The most probable sequence is *ok ok EOS* (its probability is $.4 \times .7 \times 1.0$). But greedy search doesn't find it, incorrectly choosing *yes* as the first word since it has the highest local probability (0.5).

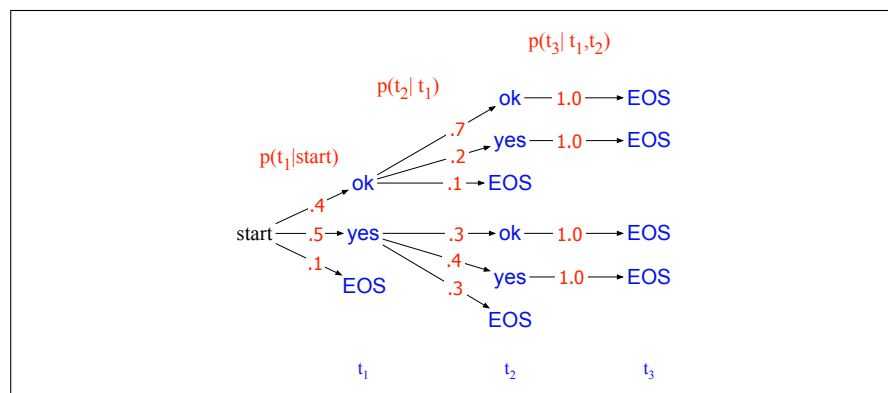


Figure 12.7 A search tree for generating the target string $T = t_1, t_2, \dots$ from vocabulary $V = \{\text{yes}, \text{ok}, \langle s \rangle\}$, showing the probability of generating each token from that state. Greedy search chooses *yes* followed by *yes*, instead of the globally most probable sequence *ok ok*.

For some problems, like part-of-speech tagging or parsing as we will see in Chapter 17 or Chapter 18, we can use dynamic programming search (the Viterbi

algorithm) to address this problem. Unfortunately, dynamic programming is not applicable to generation problems with long-distance dependencies between the output decisions. The only method guaranteed to find the best solution is exhaustive search: computing the probability of every one of the V^T possible sentences (for some length value T) which is obviously too slow.

The solution: beam search

beam search Instead, MT systems generally decode using **beam search**, a heuristic search method first proposed by [Lowerre \(1976\)](#). In beam search, instead of choosing the best token to generate at each timestep, we keep k possible tokens at each step. This fixed-size memory footprint k is called the **beam width**, on the metaphor of a flashlight beam that can be parameterized to be wider or narrower.

Thus at the first step of decoding, we compute a softmax over the entire vocabulary, assigning a probability to each word. We then select the k -best options from this softmax output. These initial k outputs are the search frontier and these k initial words are called **hypotheses**. A hypothesis is an output sequence, a translation-so-far, together with its probability.

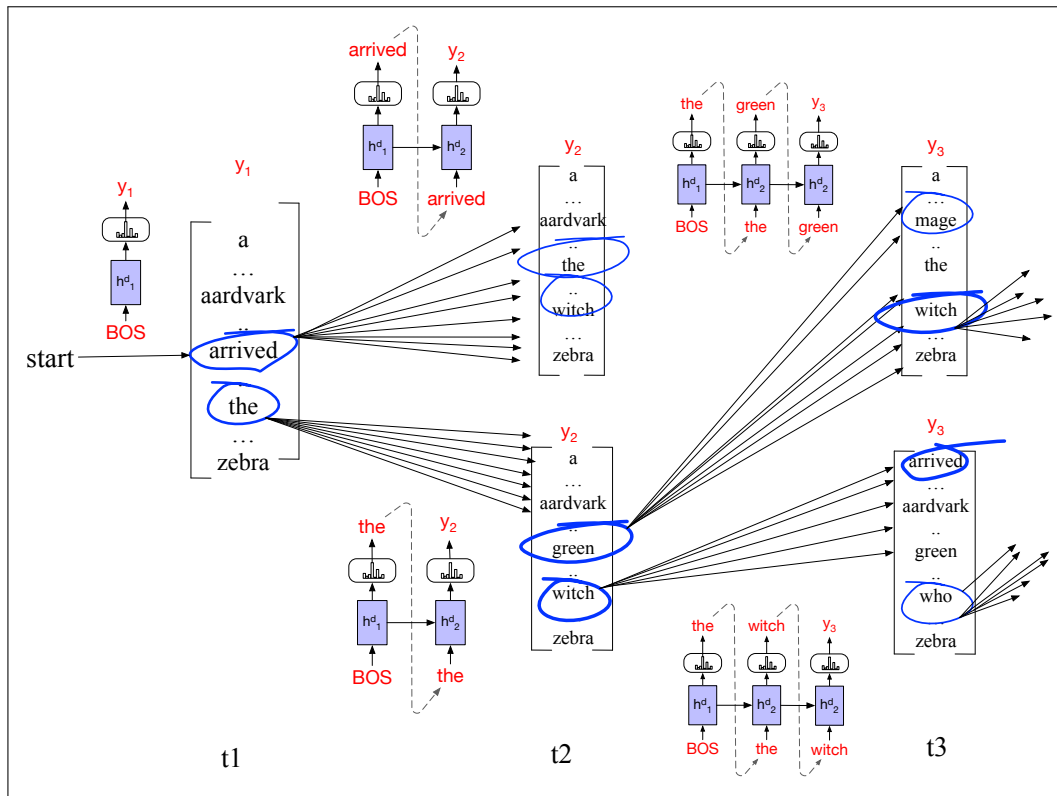


Figure 12.8 Beam search decoding with a beam width of $k = 2$. At each time step, we choose the k best hypotheses, form the V possible extensions of each, score those $k \times V$ hypotheses and choose the best $k = 2$ to continue. At time 1, the frontier has the best 2 options from the initial decoder state: *arrived* and *the*. We extend each, compute the probability of all the hypotheses so far (*arrived the*, *arrived aardvark*, *the green*, *the witch*) and again chose the best 2 (*the green* and *the witch*) to be the search frontier. The images on the arcs schematically represent the decoders that must be run at each step to score the next words (for simplicity not depicting cross-attention).

At subsequent steps, each of the k best hypotheses is extended incrementally

by being passed to distinct decoders, which each generate a softmax over the entire vocabulary to extend the hypothesis to every possible next token. Each of these $k \times V$ hypotheses is scored by $P(y_i|x, y_{<i})$: the product of the probability of the current word choice multiplied by the probability of the path that led to it. We then prune the $k \times V$ hypotheses down to the k best hypotheses, so there are never more than k hypotheses at the frontier of the search, and never more than k decoders. Fig. 12.8 illustrates this with a beam width of 2 for the beginning of *The green witch arrived*.

This process continues until an EOS is generated indicating that a complete candidate output has been found. At this point, the completed hypothesis is removed from the frontier and the size of the beam is reduced by one. The search continues until the beam has been reduced to 0. The result will be k hypotheses.

To score each node by its log probability, we use the chain rule of probability to break down $p(y|x)$ into the product of the probability of each word given its prior context, which we can turn into a sum of logs (for an output string of length t):

$$\begin{aligned}
 \text{score}(y) &= \log P(y|x) \\
 &= \log (P(y_1|x)P(y_2|y_1,x)P(y_3|y_1,y_2,x)\dots P(y_t|y_1,\dots,y_{t-1},x)) \\
 &= \sum_{i=1}^t \log P(y_i|y_1,\dots,y_{i-1},x)
 \end{aligned} \tag{12.15}$$

Thus at each step, to compute the probability of a partial sentence, we simply add the log probability of the prefix sentence so far to the log probability of generating the next token. Fig. 12.9 shows the scoring for the example sentence shown in Fig. 12.8, using some simple made-up probabilities. Log probabilities are negative or 0, and the max of two log probabilities is the one that is greater (closer to 0).

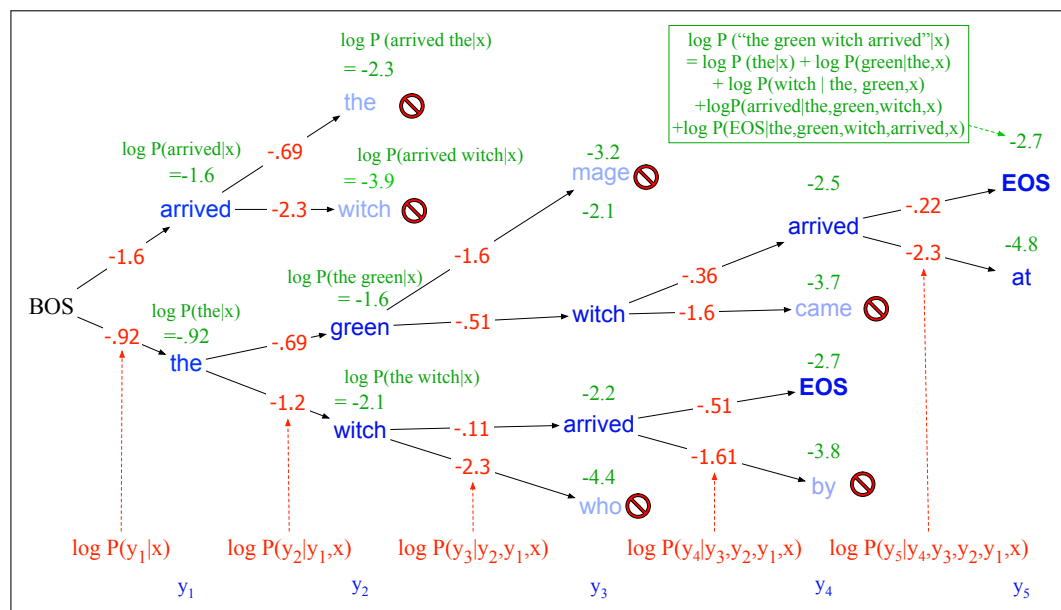


Figure 12.9 Scoring for beam search decoding with a beam width of $k = 2$. We maintain the log probability of each hypothesis in the beam by incrementally adding the logprob of generating each next token. Only the top k paths are extended to the next step.

Fig. 12.10 gives the algorithm. One problem with this version of the algorithm is that the completed hypotheses may have different lengths. Because language mod-

```

function BEAMDECODE(c, beam_width) returns best paths

  y0, h0 ← 0
  path ← ()
  complete_paths ← ()
  state ← (c, y0, h0, path) ;initial state
  frontier ← {state} ;initial frontier

  while frontier contains incomplete paths and beamwidth > 0
    extended_frontier ← {}
    for each state ∈ frontier do
      y ← DECODE(state)
      for each word i ∈ Vocabulary do
        successor ← NEWSTATE(state, i, yi)
        extended_frontier ← ADDTOBEAM(successor, extended_frontier,
                                      beam_width)

    for each state in extended_frontier do
      if state is complete do
        complete_paths ← APPEND(complete_paths, state)
        extended_frontier ← REMOVE(extended_frontier, state)
        beam_width ← beam_width - 1
    frontier ← extended_frontier

  return completed_paths

function NEWSTATE(state, word, word_prob) returns new state

function ADDTOBEAM(state, frontier, width) returns updated frontier

  if LENGTH(frontier) < width then
    frontier ← INSERT(state, frontier)
  else if SCORE(state) > SCORE(WORSTOF(frontier))
    frontier ← REMOVE(WORSTOF(frontier))
    frontier ← INSERT(state, frontier)
  return frontier

```

Figure 12.10 Beam search decoding.

els generally assign lower probabilities to longer strings, a naive algorithm would choose shorter strings for y . (This is not an issue during the earlier steps of decoding; since beam search is breadth-first, all the hypotheses being compared had the same length.) For this reason we often apply length normalization methods, like dividing the logprob by the number of words:

$$\text{score}(y) = \frac{1}{t} \log P(y|x) = \frac{1}{t} \sum_{i=1}^t \log P(y_i | y_1, \dots, y_{i-1}, x) \quad (12.16)$$

For MT we generally use beam widths k between 5 and 10, giving us k hypotheses at the end. We can pass all k to the downstream application with their respective scores, or if we just need a single translation we can pass the most probable hypothesis.

12.4.1 Minimum Bayes Risk Decoding

minimum
Bayes risk
MBR

Minimum Bayes risk or **MBR** decoding is an alternative decoding algorithm that

can work even better than beam search and also tends to be better than the other decoding algorithms like temperature sampling introduced in Section ??.

The intuition of minimum Bayes risk is that instead of trying to choose the translation which is most probable, we choose the one that is likely to have the least error. For example, we might want our decoding algorithm to find the translation which has the highest score on some evaluation metric. For example in Section 12.6 we will introduce metrics like chrF or BERTScore that measure the goodness-of-fit between a candidate translation and a set of reference human translations. A translation that maximizes this score, especially with a hypothetically huge set of perfect human translations is likely to be a good one (have minimum risk) even if it is not the most probable translation by our particular probability estimator.

In practice, we don't know the perfect set of translations for a given sentence. So the standard simplification used in MBR decoding algorithms is to instead choose the candidate translation which is most similar (by some measure of goodness-of-fit) with some set of candidate translations. We're essentially approximating the enormous space of all possible translations \mathcal{U} with a smaller set of possible candidate translations \mathcal{Y} .

Given this set of possible candidate translations \mathcal{Y} , and some similarity or alignment function util , we choose the best translation \hat{y} as the translation which is most similar to all the other candidate translations:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{c \in \mathcal{Y}} \text{util}(y, c) \quad (12.17)$$

Various util functions can be used, like chrF or BERTscore or BLEU. We can get the set of candidate translations by sampling using one of the basic sampling algorithms of Section ?? like temperature sampling; good results can be obtained with as few as 32 or 64 candidates.

Minimum Bayes risk decoding can also be used for other NLP tasks; indeed it was widely applied to speech recognition (Stolcke et al., 1997; Goel and Byrne, 2000) before being applied to machine translation (Kumar and Byrne, 2004), and has been shown to work well across many other generation tasks as well (e.g., summarization, dialogue, and image captioning (Suzgun et al., 2023)).

12.5 Translating in low-resource situations

For some languages, and especially for English, online resources are widely available. There are many large parallel corpora that contain translations between English and many languages. But the vast majority of the world's languages do not have large parallel training texts available. An important ongoing research question is how to get good translation with lesser resourced languages. The resource problem can even be true for high resource languages when we need to translate into low resource domains (for example in a particular genre that happens to have very little bitext).

Here we briefly introduce two commonly used approaches for dealing with this data sparsity: **backtranslation**, which is a special case of the general statistical technique called **data augmentation**, and **multilingual models**, and also discuss some socio-technical issues.

12.5.1 Data Augmentation

Data augmentation is a statistical technique for dealing with insufficient training data, by adding new synthetic data that is generated from the current natural data.

The most common data augmentation technique for machine translation is called **backtranslation**. Backtranslation relies on the intuition that while parallel corpora may be limited for particular languages or domains, we can often find a large (or at least larger) monolingual corpus, to add to the smaller parallel corpora that are available. The algorithm makes use of monolingual corpora in the **target** language by creating synthetic bitexts.

In backtranslation, our goal is to improve source-to-target MT, given a small parallel text (a bitext) in the source/target languages, and some monolingual data in the target language. We first use the bitext to train a MT system in the **reverse** direction: a target-to-source MT system. We then use it to translate the monolingual target data to the source language. Now we can add this synthetic bitext (natural target sentences, aligned with MT-produced source sentences) to our training data, and retrain our source-to-target MT model. For example suppose we want to translate from Navajo to English but only have a small Navajo-English bitext, although of course we can find lots of monolingual English data. We use the small bitext to build an MT engine going the other way (from English to Navajo). Once we translate the monolingual English text to Navajo, we can add this synthetic Navajo/English bitext to our training data.

Backtranslation has various parameters. One is how we generate the backtranslated data; we can run the decoder in greedy inference, or use beam search. Or we can do sampling, like the temperature sampling algorithm we saw in Chapter 8. Another parameter is the ratio of backtranslated data to natural bitext data; we can choose to upsample the bitext data (include multiple copies of each sentence). In general backtranslation works surprisingly well; one estimate suggests that a system trained on backtranslated text gets about 2/3 of the gain as would training on the same amount of natural bitext (Edunov et al., 2018).

12.5.2 Multilingual models

The models we’ve described so far are for bilingual translation: one source language, one target language. It’s also possible to build a **multilingual** translator.

In a multilingual translator, we train the system by giving it parallel sentences in many different pairs of languages. That means we need to tell the system which language to translate from and to! We tell the system which language is which by adding a special token l_s to the encoder specifying the source language we’re translating from, and a special token l_t to the decoder telling it the target language we’d like to translate into.

Thus we slightly update Eq. 12.9 above to add these tokens in Eq. 12.19:

$$\mathbf{h} = \text{encoder}(x, l_s) \quad (12.18)$$

$$y_{i+1} = \text{decoder}(\mathbf{h}, l_t, y_1, \dots, y_i) \quad \forall i \in [1, \dots, m] \quad (12.19)$$

One advantage of a multilingual model is that they can improve the translation of lower-resourced languages by drawing on information from a similar language in the training data that happens to have more resources. Perhaps we don’t know the meaning of a word in Galician, but the word appears in the similar and higher-resourced language Spanish.

12.5.3 Sociotechnical issues

Many issues in dealing with low-resource languages go beyond the purely technical. One problem is that for low-resource languages, especially from low-income countries, native speakers are often not involved as the curators for content selection, as the language technologists, or as the evaluators who measure performance (V et al., 2020). Indeed, one well-known study that manually audited a large set of parallel corpora and other major multilingual datasets found that for many of the corpora, less than 50% of the sentences were of acceptable quality, with a lot of data consisting of repeated sentences with web boilerplate or incorrect translations, suggesting that native speakers may not have been sufficiently involved in the data process (Kreutzer et al., 2022).

Other issues, like the tendency of many MT approaches to focus on the case where one of the languages is English (Anastasopoulos and Neubig, 2020), have to do with allocation of resources. Where most large multilingual systems were trained on bitexts in which English was one of the two languages, recent huge corporate systems like those of Fan et al. (2021) and Costa-jussà et al. (2022) and datasets like Schwenk et al. (2021) attempt to handle large numbers of languages (up to 200 languages) and create bitexts between many more pairs of languages and not just through English.

At the smaller end, V et al. (2020) propose a participatory design process to encourage content creators, curators, and language technologists who speak these low-resourced languages to participate in developing MT algorithms. They provide online groups, mentoring, and infrastructure, and report on a case study on developing MT algorithms for low-resource African languages. Among their conclusions was to perform MT evaluation by post-editing rather than direct evaluation, since having labelers edit an MT system and then measure the distance between the MT output and its post-edited version both was simpler to train evaluators and makes it easier to measure true errors in the MT output and not differences due to linguistic variation (Bentivogli et al., 2018).

12.6 MT Evaluation

Translations are evaluated along two dimensions:

- | | |
|-----------------|--|
| adequacy | 1. adequacy : how well the translation captures the exact meaning of the source sentence. Sometimes called faithfulness or fidelity . |
| fluency | 2. fluency : how fluent the translation is in the target language (is it grammatical, clear, readable, natural). |

Using humans to evaluate is most accurate, but automatic metrics are also used for convenience.

12.6.1 Using Human Raters to Evaluate MT

The most accurate evaluations use human raters, such as online crowdworkers, to evaluate each translation along the two dimensions. For example, along the dimension of **fluency**, we can ask how intelligible, how clear, how readable, or how natural the MT output (the target text) is. We can give the raters a scale, for example, from 1 (totally unintelligible) to 5 (totally intelligible), or 1 to 100, and ask them to rate each sentence or paragraph of the MT output.

We can do the same thing to judge the second dimension, **adequacy**, using raters to assign scores on a scale. If we have bilingual raters, we can give them the source sentence and a proposed target sentence, and rate, on a 5-point or 100-point scale, how much of the information in the source was preserved in the target. If we only have monolingual raters but we have a good human translation of the source text, we can give the monolingual raters the human reference translation and a target machine translation and again rate how much information is preserved. An alternative is to do **ranking**: give the raters a pair of candidate translations, and ask them which one they prefer.

Training of human raters (who are often online crowdworkers) is essential; raters without translation expertise find it difficult to separate fluency and adequacy, and so training includes examples carefully distinguishing these. Raters often disagree (source sentences may be ambiguous, raters will have different world knowledge, raters may apply scales differently). It is therefore common to remove outlier raters, and (if we use a fine-grained enough scale) normalizing raters by subtracting the mean from their scores and dividing by the variance.

As discussed above, an alternative way of using human raters is to have them **post-edit** translations, taking the MT output and changing it minimally until they feel it represents a correct translation. The difference between their post-edited translations and the original MT output can then be used as a measure of quality.

12.6.2 Automatic Evaluation

While humans produce the best evaluations of machine translation output, running a human evaluation can be time consuming and expensive. For this reason automatic metrics are often used as temporary proxies. Automatic metrics are less accurate than human evaluation, but can help test potential system improvements, and even be used as an automatic loss function for training. In this section we introduce two families of such metrics, those based on character- or word-overlap and those based on embedding similarity.

Automatic Evaluation by Character Overlap: chrF

The simplest and most robust metric for MT evaluation is called **chrF**, which stands for **character F-score** (Popović, 2015). chrF (along with many other earlier related metrics like BLEU, METEOR, TER, and others) is based on a simple intuition derived from the pioneering work of Miller and Beebe-Center (1956): a good machine translation will tend to contain characters and words that occur in a human translation of the same sentence. Consider a test set from a parallel corpus, in which each source sentence has both a gold human target translation and a candidate MT translation we'd like to evaluate. The chrF metric ranks each MT target sentence by a function of the number of character n-gram overlaps with the human translation.

Given the hypothesis and the reference, chrF is given a parameter k indicating the length of character n-grams to be considered, and computes the average of the k precisions (unigram precision, bigram, and so on) and the average of the k recalls (unigram recall, bigram recall, etc.):

chrP percentage of character 1-grams, 2-grams, ..., k-grams in the hypothesis that occur in the reference, averaged.

chrR percentage of character 1-grams, 2-grams, ..., k-grams in the reference that occur in the hypothesis, averaged.

The metric then computes an F-score by combining chrP and chrR using a weighting

parameter β . It is common to set $\beta = 2$, thus weighing recall twice as much as precision:

$$\text{chrF}\beta = (1 + \beta^2) \frac{\text{chrP} \cdot \text{chrR}}{\beta^2 \cdot \text{chrP} + \text{chrR}} \quad (12.20)$$

For $\beta = 2$, that would be:

$$\text{chrF2} = \frac{5 \cdot \text{chrP} \cdot \text{chrR}}{4 \cdot \text{chrP} + \text{chrR}}$$

For example, consider two hypotheses that we'd like to score against the reference translation *witness for the past*. Here are the hypotheses along with chrF values computed using parameters $k = \beta = 2$ (in real examples, k would be a higher number like 6):

REF: witness for the past,	
HYP1: witness of the past,	chrF2,2 = .86
HYP2: past witness	chrF2,2 = .62

Let's see how we computed that chrF value for HYP1 (we'll leave the computation of the chrF value for HYP2 as an exercise for the reader). First, chrF ignores spaces, so we'll remove them from both the reference and hypothesis:

REF: witnessforthepast, (18 unigrams, 17 bigrams)
 HYP1: witnessofthepast, (17 unigrams, 16 bigrams)

Next let's see how many unigrams and bigrams match between the reference and hypothesis:

unigrams that match: w i t n e s s f o t h e p a s t , (17 unigrams)
 bigrams that match: wi it tn ne es ss th he ep pa as st t, (13 bigrams)

We use that to compute the unigram and bigram precisions and recalls:

unigram P: $17/17 = 1$ unigram R: $17/18 = .944$
 bigram P: $13/16 = .813$ bigram R: $13/17 = .765$

Finally we average to get chrP and chrR, and compute the F-score:

$$\begin{aligned} \text{chrP} &= (17/17 + 13/16)/2 = .906 \\ \text{chrR} &= (17/18 + 13/17)/2 = .855 \\ \text{chrF2,2} &= 5 \frac{\text{chrP} * \text{chrR}}{4\text{chrP} + \text{chrR}} = .86 \end{aligned}$$

chrF is simple, robust, and correlates very well with human judgments in many languages (Kocmi et al., 2021).

Alternative overlap metric: BLEU

There are various alternative overlap metrics. For example, before the development of chrF, it was common to use a word-based overlap metric called **BLEU** (for BiLingual Evaluation Understudy), that is purely precision-based rather than combining precision and recall (Papineni et al., 2002). The BLEU score for a corpus of candidate translation sentences is a function of the **n-gram word precision** over all the sentences combined with a brevity penalty computed over the corpus as a whole.

What do we mean by n-gram precision? Consider a corpus composed of a single sentence. The unigram precision for this corpus is the percentage of unigram tokens

in the candidate translation that also occur in the reference translation, and ditto for bigrams and so on, up to 4-grams. BLEU extends this unigram metric to the whole corpus by computing the numerator as the sum over all sentences of the counts of all the unigram types that also occur in the reference translation, and the denominator is the total of the counts of all unigrams in all candidate sentences. We compute this n-gram precision for unigrams, bigrams, trigrams, and 4-grams and take the geometric mean. BLEU has many further complications, including a brevity penalty for penalizing candidate translations that are too short, and it also requires the n-gram counts be clipped in a particular way.

Because BLEU is a word-based metric, it is very sensitive to word tokenization, making it impossible to compare different systems if they rely on different tokenization standards, and doesn't work as well in languages with complex morphology. Nonetheless, you will sometimes still see systems evaluated by BLEU, particularly for translation into English. In such cases it's important to use packages that enforce standardization for tokenization like SACREBLEU (Post, 2018).

Statistical Significance Testing for MT evals

Character or word overlap-based metrics like chrF (or BLEU, or etc.) are mainly used to compare two systems, with the goal of answering questions like: did the new algorithm we just invented improve our MT system? To know if the difference between the chrF scores of two MT systems is a significant difference, we use the paired bootstrap test, or the similar randomization test.

To get a confidence interval on a single chrF score using the bootstrap test, recall from Section ?? that we take our test set (or devset) and create thousands of pseudo-testsets by repeatedly sampling with replacement from the original test set. We now compute the chrF score of each of the pseudo-testsets. If we drop the top 2.5% and bottom 2.5% of the scores, the remaining scores will give us the 95% confidence interval for the chrF score of our system.

To compare two MT systems A and B, we draw the same set of pseudo-testsets, and compute the chrF scores for each of them. We then compute the percentage of pseudo-test-sets in which A has a higher chrF score than B.

chrF: Limitations

While automatic character and word-overlap metrics like chrF or BLEU are useful, they have important limitations. chrF is very local: a large phrase that is moved around might barely change the chrF score at all, and chrF can't evaluate cross-sentence properties of a document like its discourse coherence (Chapter 24). chrF and similar automatic metrics also do poorly at comparing very different kinds of systems, such as comparing human-aided translation against machine translation, or different machine translation architectures against each other (Callison-Burch et al., 2006). Instead, automatic overlap metrics like chrF are most appropriate when evaluating changes to a single system.

12.6.3 Automatic Evaluation: Embedding-Based Methods

The chrF metric is based on measuring the exact character n-grams a human reference and candidate machine translation have in common. However, this criterion is overly strict, since a good translation may use alternate words or paraphrases. A solution first pioneered in early metrics like METEOR (Banerjee and Lavie, 2005) was to allow synonyms to match between the reference x and candidate \tilde{x} . More

recent metrics use BERT or other embeddings to implement this intuition.

For example, in some situations we might have datasets that have human assessments of translation quality. Such datasets consists of tuples (x, \tilde{x}, r) , where $x = (x_1, \dots, x_n)$ is a reference translation, $\tilde{x} = (\tilde{x}_1, \dots, \tilde{x}_m)$ is a candidate machine translation, and $r \in \mathbb{R}$ is a human rating that expresses the quality of \tilde{x} with respect to x . Given such data, algorithms like COMET (Rei et al., 2020) BLEURT (Sellam et al., 2020) train a predictor on the human-labeled datasets, for example by passing x and \tilde{x} through a version of BERT (trained with extra pretraining, and then finetuned on the human-labeled sentences), followed by a linear layer that is trained to predict r . The output of such models correlates highly with human labels.

In other cases, however, we don't have such human-labeled datasets. In that case we can measure the similarity of x and \tilde{x} by the similarity of their embeddings. The BERTSCORE algorithm (Zhang et al., 2020) shown in Fig. 12.11, for example, passes the reference x and the candidate \tilde{x} through BERT, computing a BERT embedding for each token x_i and \tilde{x}_j . Each pair of tokens (x_i, \tilde{x}_j) is scored by its cosine $\frac{x_i \cdot \tilde{x}_j}{|x_i| |\tilde{x}_j|}$. Each token in x is matched to a token in \tilde{x} to compute recall, and each token in \tilde{x} is matched to a token in x to compute precision (with each token greedily matched to the most similar token in the corresponding sentence). BERTSCORE provides precision and recall (and hence F_1):

$$R_{\text{BERT}} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\tilde{x}_j \in \tilde{x}} x_i \cdot \tilde{x}_j \quad P_{\text{BERT}} = \frac{1}{|\tilde{x}|} \sum_{\tilde{x}_j \in \tilde{x}} \max_{x_i \in x} x_i \cdot \tilde{x}_j \quad (12.21)$$

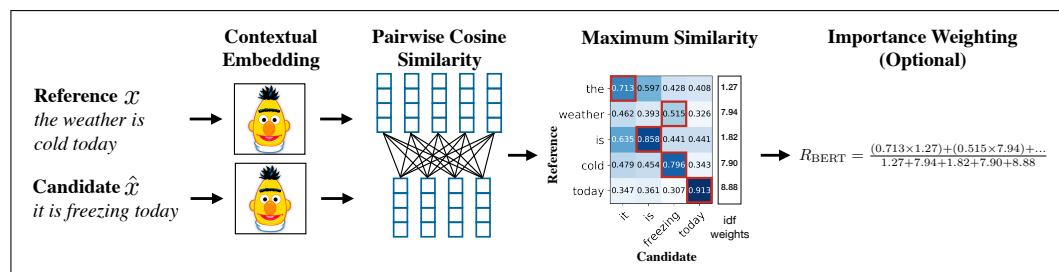


Figure 12.11 The computation of BERTSCORE recall from reference x and candidate \hat{x} , from Figure 1 in Zhang et al. (2020). This version shows an extended version of the metric in which tokens are also weighted by their idf values.

12.7 Bias and Ethical Issues

Machine translation raises many of the same ethical issues that we've discussed in earlier chapters. For example, consider MT systems translating from Hungarian (which has the gender neutral pronoun *ő*) or Spanish (which often drops pronouns) into English (in which pronouns are obligatory, and they have grammatical gender). When translating a reference to a person described without specified gender, MT systems often default to male gender (Schiebinger 2014, Prates et al. 2019). And MT systems often assign gender according to culture stereotypes of the sort we saw in Section ?? . Fig. 12.12 shows examples from Prates et al. (2019), in which Hungarian gender-neutral *ő is a nurse* is translated with *she*, but gender-neutral *ő is a CEO* is translated with *he*. Prates et al. (2019) find that these stereotypes can't completely be accounted for by gender bias in US labor statistics, because the biases are

amplified by MT systems, with pronouns being mapped to male or female gender with a probability higher than if the mapping was based on actual labor employment statistics.

Hungarian (gender neutral) source	English MT output
ő egy ápoló	she is a nurse
ő egy tudós	he is a scientist
ő egy mérnök	he is an engineer
ő egy pék	he is a baker
ő egy tanár	she is a teacher
ő egy esküvőszervező	she is a wedding organizer
ő egy vezérigazgató	he is a CEO

Figure 12.12 When translating from gender-neutral languages like Hungarian into English, current MT systems interpret people from traditionally male-dominated occupations as male, and traditionally female-dominated occupations as female (Prates et al., 2019).

Similarly, a recent challenge set, the WinoMT dataset (Stanovsky et al., 2019) shows that MT systems perform worse when they are asked to translate sentences that describe people with non-stereotypical gender roles, like “The doctor asked the nurse to help her in the operation”.

Many ethical questions in MT require further research. One open problem is developing metrics for knowing what our systems don’t know. This is because MT systems can be used in urgent situations where human translators may be unavailable or delayed: in medical domains, to help translate when patients and doctors don’t speak the same language, or in legal domains, to help judges or lawyers communicate with witnesses or defendants. In order to ‘do no harm’, systems need ways to assign **confidence** values to candidate translations, so they can abstain from giving incorrect translations that may cause harm.

confidence

12.8 Summary

Machine translation is one of the most widely used applications of NLP, and the encoder-decoder model, first developed for MT is a key tool that has applications throughout NLP.

- Languages have **divergences**, both structural and lexical, that make translation difficult.
- The linguistic field of **typology** investigates some of these differences; languages can be classified by their position along typological dimensions like whether verbs precede their objects.
- **Encoder-decoder** networks (for transformers just as we saw in Chapter 13 for RNNs) are composed of an **encoder** network that takes an input sequence and creates a contextualized representation of it, the **context**. This context representation is then passed to a **decoder** which generates a task-specific output sequence.
- **Cross-attention** allows the transformer decoder to view information from all the hidden states of the encoder.
- Machine translation models are trained on a **parallel corpus**, sometimes called a **bitext**, a text that appears in two (or more) languages.

- **Backtranslation** is a way of making use of monolingual corpora in the target language by running a pilot MT engine backwards to create synthetic bitexts.
- MT is evaluated by measuring a translation’s **adequacy** (how well it captures the meaning of the source sentence) and **fluency** (how fluent or natural it is in the target language). Human evaluation is the gold standard, but automatic evaluation metrics like **chrF**, which measure character n-gram overlap with human translations, or more recent metrics based on embedding similarity, are also commonly used.

Historical Notes

MT was proposed seriously by the late 1940s, soon after the birth of the computer (Weaver, 1949/1955). In 1954, the first public demonstration of an MT system prototype (Dostert, 1955) led to great excitement in the press (Hutchins, 1997). The next decade saw a great flowering of ideas, prefiguring most subsequent developments. But this work was ahead of its time—implementations were limited by, for example, the fact that pending the development of disks there was no good way to store dictionary information.

As high-quality MT proved elusive (Bar-Hillel, 1960), there grew a consensus on the need for better evaluation and more basic research in the new fields of formal and computational linguistics. This consensus culminated in the famously critical ALPAC (Automatic Language Processing Advisory Committee) report of 1966 (Pierce et al., 1966) that led in the mid 1960s to a dramatic cut in funding for MT in the US. As MT research lost academic respectability, the Association for Machine Translation and Computational Linguistics dropped MT from its name. Some MT developers, however, persevered, and there were early MT systems like *Météo*, which translated weather forecasts from English to French (Chandioux, 1976), and industrial systems like Systran.

In the early years, the space of MT architectures spanned three general models. In **direct translation**, the system proceeds word-by-word through the source-language text, translating each word incrementally. Direct translation uses a large bilingual dictionary, each of whose entries is a small program with the job of translating one word. In **transfer** approaches, we first parse the input text and then apply rules to transform the source-language parse into a target language parse. We then generate the target language sentence from the parse tree. In **interlingua** approaches, we analyze the source language text into some abstract meaning representation, called an **interlingua**. We then generate into the target language from this interlingual representation. A common way to visualize these three early approaches was the **Vauquois triangle** shown in Fig. 12.13. The triangle shows the increasing depth of analysis required (on both the analysis and generation end) as we move from the direct approach through transfer approaches to interlingual approaches. In addition, it shows the decreasing amount of transfer knowledge needed as we move up the triangle, from huge amounts of transfer at the direct level (almost all knowledge is transfer knowledge for each word) through transfer (transfer rules only for parse trees or thematic roles) through interlingua (no specific transfer knowledge). We can view the encoder-decoder network as an interlingual approach, with attention acting as an integration of direct and transfer, allowing words or their representations to be directly accessed by the decoder.

Vauquois
triangle

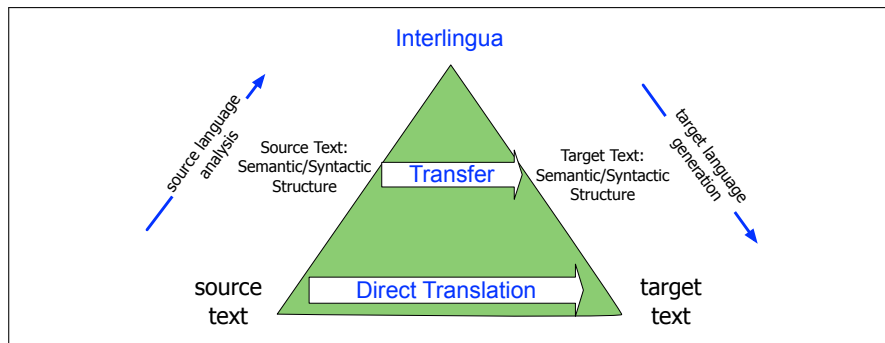


Figure 12.13 The Vauquois (1968) triangle.

Statistical methods began to be applied around 1990, enabled first by the development of large bilingual corpora like the **Hansard** corpus of the proceedings of the Canadian Parliament, which are kept in both French and English, and then by the growth of the web. Early on, a number of researchers showed that it was possible to extract pairs of aligned sentences from bilingual corpora, using words or simple cues like sentence length (Kay and Röscheisen 1988, Gale and Church 1991, Gale and Church 1993, Kay and Röscheisen 1993).

statistical MT
IBM Models
Candide

At the same time, the IBM group, drawing directly on the **noisy channel** model for speech recognition, proposed two related paradigms for **statistical MT**. These include the generative algorithms that became known as **IBM Models 1 through 5**, implemented in the **Candide** system. The algorithms (except for the decoder) were published in full detail—encouraged by the US government who had partially funded the work—which gave them a huge impact on the research community (Brown et al. 1990, Brown et al. 1993).

The group also developed a discriminative approach, called MaxEnt (for maximum entropy, an alternative formulation of logistic regression), which allowed many features to be combined discriminatively rather than generatively (Berger et al., 1996), which was further developed by Och and Ney (2002).

phrase-based
translation

By the turn of the century, most academic research on machine translation used statistical MT, either in the generative or discriminative mode. An extended version of the generative approach, called **phrase-based translation** was developed, based on inducing translations for phrase-pairs (Och 1998, Marcu and Wong 2002, Koehn et al. (2003), Och and Ney 2004, Deng and Byrne 2005, inter alia).

MERT

Once automatic metrics like BLEU were developed (Papineni et al., 2002), the discriminative log linear formulation (Och and Ney, 2004), drawing from the IBM MaxEnt work (Berger et al., 1996), was used to directly optimize evaluation metrics like BLEU in a method known as **Minimum Error Rate Training**, or **MERT** (Och, 2003), also drawing from speech recognition models (Chou et al., 1993). Toolkits like GIZA (Och and Ney, 2003) and **Moses** (Koehn et al. 2006, Zens and Ney 2007) were widely used.

Moses

transduction
grammars

There were also approaches around the turn of the century that were based on syntactic structure (Chapter 18). Models based on **transduction grammars** (also called **synchronous grammars**) assign a parallel syntactic tree structure to a pair of sentences in different languages, with the goal of translating the sentences by applying reordering operations on the trees. From a generative perspective, we can view a transduction grammar as generating pairs of aligned sentences in two languages. Some of the most widely used models included the **inversion transduction grammar** (Wu, 1996) and synchronous context-free grammars (Chiang, 2005),

inversion
transduction
grammar

Neural networks had been applied at various times to various aspects of machine translation; for example [Schwenk et al. \(2006\)](#) showed how to use neural language models to replace n-gram language models in a Spanish-English system based on IBM Model 4. The modern neural encoder-decoder approach was pioneered by [Kalchbrenner and Blunsom \(2013\)](#), who used a CNN encoder and an RNN decoder, and was first applied to MT by [Bahdanau et al. \(2015\)](#). The transformer encoder-decoder was proposed by [Vaswani et al. \(2017\)](#) (see the History section of Chapter 8).

Research on evaluation of machine translation began quite early. [Miller and Beebe-Center \(1956\)](#) proposed a number of methods drawing on work in psycholinguistics. These included the use of cloze and Shannon tasks to measure intelligibility as well as a metric of edit distance from a human translation, the intuition that underlies all modern overlap-based automatic evaluation metrics. The ALPAC report included an early evaluation study conducted by John Carroll that was extremely influential ([Pierce et al., 1966](#), Appendix 10). Carroll proposed distinct measures for fidelity and intelligibility, and had raters score them subjectively on 9-point scales. Much early evaluation work focuses on automatic word-overlap metrics like BLEU ([Papineni et al., 2002](#)), NIST ([Doddington, 2002](#)), **TER (Translation Error Rate)** ([Snover et al., 2006](#)), **Precision and Recall** ([Turian et al., 2003](#)), and **METEOR** ([Banerjee and Lavie, 2005](#)); character n-gram overlap methods like chrF ([Popović, 2015](#)) came later. More recent evaluation work, echoing the ALPAC report, has emphasized the importance of careful statistical methodology and the use of human evaluation ([Kocmi et al., 2021](#); [Marie et al., 2021](#)).

The early history of MT is surveyed in Hutchins [1986](#) and [1997](#); [Nirenburg et al. \(2002\)](#) collects early readings. See [Croft \(1990\)](#) or [Comrie \(1989\)](#) for introductions to linguistic typology.

Exercises

- 12.1** Compute by hand the chrF2,2 score for HYP2 on page [20](#) (the answer should round to .62).

- Anastasopoulos, A. and G. Neubig. 2020. [Should all cross-lingual embeddings speak English?](#) *ACL*.
- Artetxe, M. and H. Schwenk. 2019. [Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond](#). *TACL*, 7:597–610.
- Bahdanau, D., K. H. Cho, and Y. Bengio. 2015. Neural machine translation by jointly learning to align and translate. *ICLR 2015*.
- Banerjee, S. and A. Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). *Proceedings of ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization*.
- Bañón, M., P. Chen, B. Haddow, K. Heafield, H. Hoang, M. Esplà-Gomis, M. L. Forcada, A. Kamran, F. Kirefu, P. Koehn, S. Ortiz Rojas, L. Pla Sempere, G. Ramírez-Sánchez, E. Sarriás, M. Strelec, B. Thompson, W. Waites, D. Wiggins, and J. Zaragoza. 2020. [ParaCrawl: Web-scale acquisition of parallel corpora](#). *ACL*.
- Bar-Hillel, Y. 1960. The present status of automatic translation of languages. In F. Alt, ed., *Advances in Computers I*, 91–163. Academic Press.
- Bentivogli, L., M. Cettolo, M. Federico, and C. Federmann. 2018. [Machine translation human evaluation: an investigation of evaluation based on post-editing and its relation with direct assessment](#). *ICSLT*.
- Berger, A., S. A. Della Pietra, and V. J. Della Pietra. 1996. [A maximum entropy approach to natural language processing](#). *Computational Linguistics*, 22(1):39–71.
- Bickel, B. 2003. Referential density in discourse and syntactic typology. *Language*, 79(2):708–736.
- Bostrom, K. and G. Durrett. 2020. [Byte pair encoding is suboptimal for language model pretraining](#). *EMNLP*.
- Brown, P. F., J. Cocke, S. A. Della Pietra, V. J. Della Pietra, F. Jelinek, J. D. Lafferty, R. L. Mercer, and P. S. Roossin. 1990. [A statistical approach to machine translation](#). *Computational Linguistics*, 16(2):79–85.
- Brown, P. F., S. A. Della Pietra, V. J. Della Pietra, and R. L. Mercer. 1993. [The mathematics of statistical machine translation: Parameter estimation](#). *Computational Linguistics*, 19(2):263–311.
- Callison-Burch, C., M. Osborne, and P. Koehn. 2006. [Re-evaluating the role of BLEU in machine translation research](#). *EACL*.
- Chandioux, J. 1976. MÉTÉO: un système opérationnel pour la traduction automatique des bulletins météorologiques destinés au grand public. *Meta*, 21:127–133.
- Chiang, D. 2005. [A hierarchical phrase-based model for statistical machine translation](#). *ACL*.
- Chou, W., C.-H. Lee, and B. H. Juang. 1993. [Minimum error rate training based on \$n\$ -best string models](#). *ICASSP*.
- Comrie, B. 1989. *Language Universals and Linguistic Typology*, 2nd edition. Blackwell.
- Costa-jussà, M. R., J. Cross, O. Çelebi, M. Elbayad, K. Heafield, K. Heffernan, E. Kalbassi, J. Lam, D. Licht, J. Maillard, A. Sun, S. Wang, G. Wenzek, A. Youngblood, B. Akula, L. Barrault, G. M. Gonzalez, P. Hansanti, J. Hoffman, S. Jarrett, K. R. Sadagopan, D. Rowe, S. Spruit, C. Tran, P. Andrews, N. F. Ayan, S. Bhosale, S. Edunov, A. Fan, C. Gao, V. Goswami, F. Guzmán, P. Koehn, A. Mourachko, C. Ropers, S. Saleem, H. Schwenk, J. Wang, and NLLB Team. 2022. [No language left behind: Scaling human-centered machine translation](#). ArXiv.
- Croft, W. 1990. *Typology and Universals*. Cambridge University Press.
- Deng, Y. and W. Byrne. 2005. [HMM word and phrase alignment for statistical machine translation](#). *HLT-EMNLP*.
- Doddington, G. 2002. [Automatic evaluation of machine translation quality using \$n\$ -gram co-occurrence statistics](#). *HLT*.
- Dorr, B. 1994. [Machine translation divergences: A formal description and proposed solution](#). *Computational Linguistics*, 20(4):597–633.
- Dostert, L. 1955. The Georgetown-I.B.M. experiment. In *Machine Translation of Languages: Fourteen Essays*, 124–135. MIT Press.
- Dryer, M. S. and M. Haspelmath, eds. 2013. *The World Atlas of Language Structures Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig. Available online at <http://wals.info>.
- Edunov, S., M. Ott, M. Auli, and D. Grangier. 2018. [Understanding back-translation at scale](#). *EMNLP*.
- Fan, A., S. Bhosale, H. Schwenk, Z. Ma, A. El-Kishky, S. Goyal, M. Baines, O. Celebi, G. Wenzek, V. Chaudhary, N. Goyal, T. Birch, V. Liptchinsky, S. Edunov, M. Auli, and A. Joulin. 2021. [Beyond english-centric multilingual machine translation](#). *JMLR*, 22(107):1–48.
- ∀, W. Nekoto, V. Marivate, T. Matsila, T. Fasubaa, T. Kolawole, T. Fagbohunbe, S. O. Akinola, S. H. Muhammad, S. Kabongo, S. Osei, S. Freshia, R. A. Niyongabo, R. M. P. Ogayo, O. Ahia, M. Meressa, M. Adeyemi, M. Mokgesi-Seling, L. Okegbemi, L. J. Martinus, K. Tajudeen, K. Degila, K. Ogueji, K. Siminyu, J. Kreutzer, J. Webster, J. T. Ali, J. A. I. Orife, I. Ezeani, I. A. Dangana, H. Kamper, H. Elshahar, G. Duru, G. Kioko, E. Murhabazi, E. van Biljon, D. Whitenack, C. Onyefuluchi, C. Emezue, B. Dossou, B. Sibanda, B. I. Bassey, A. Olabiyi, A. Ramkilowan, A. Öktem, A. Akinfaderin, and A. Bashir. 2020. [Participatory research for low-resourced machine translation: A case study in African languages](#). *Findings of EMNLP*. The authors use the forall symbol to represent the whole Masakhane community.
- Gale, W. A. and K. W. Church. 1991. [A program for aligning sentences in bilingual corpora](#). *ACL*.
- Gale, W. A. and K. W. Church. 1993. [A program for aligning sentences in bilingual corpora](#). *Computational Linguistics*, 19:75–102.
- Goel, V. and W. Byrne. 2000. [Minimum bayes-risk automatic speech recognition](#). *Computer Speech & Language*, 14(2):115–135.
- Hutchins, W. J. 1986. *Machine Translation: Past, Present, Future*. Ellis Horwood, Chichester, England.

- Hutchins, W. J. 1997. [From first conception to first demonstration: The nascent years of machine translation, 1947–1954. A chronology.](#) *Machine Translation*, 12:192–252.
- Hutchins, W. J. and H. L. Somers. 1992. *An Introduction to Machine Translation*. Academic Press.
- Kalchbrenner, N. and P. Blunsom. 2013. [Recurrent continuous translation models.](#) *EMNLP*.
- Kay, M. and M. Röscheisen. 1988. Text-translation alignment. Technical Report P90-00143, Xerox Palo Alto Research Center, Palo Alto, CA.
- Kay, M. and M. Röscheisen. 1993. [Text-translation alignment.](#) *Computational Linguistics*, 19:121–142.
- Kocmi, T., C. Federmann, R. Grundkiewicz, M. Junczys-Dowmunt, H. Matsushita, and A. Menezes. 2021. [To ship or not to ship: An extensive evaluation of automatic metrics for machine translation.](#) ArXiv.
- Koehn, P. 2005. Europarl: A parallel corpus for statistical machine translation. *MT summit*, vol. 5.
- Koehn, P., H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. 2006. [Moses: Open source toolkit for statistical machine translation.](#) *ACL*.
- Koehn, P., F. J. Och, and D. Marcu. 2003. [Statistical phrase-based translation.](#) *HLT-NAACL*.
- Kreutzer, J., I. Caswell, L. Wang, A. Wahab, D. van Esch, N. Ulzii-Orshikh, A. Tapo, N. Subramani, A. Sokolov, C. Sikasote, M. Setyawan, S. Sarin, S. Samb, B. Sagot, C. Rivera, A. Rios, I. Papadimitriou, S. Osei, P. O. Suarez, I. Orife, K. Ogueji, A. N. Rubungo, T. Q. Nguyen, M. Müller, A. Müller, S. H. Muhammad, N. Muhammad, A. Mnyakeni, J. Mirzakhlov, T. Matangira, C. Leong, N. Lawson, S. Kudugunta, Y. Jernite, M. Jenny, O. Firat, B. F. P. Dossou, S. Dlamini, N. de Silva, S. Çabuk Ballı, S. Biderman, A. Battisti, A. Baruwā, A. Bapna, P. Baljekar, I. A. Azime, A. Awokoya, D. Ataman, O. Ahia, O. Ahia, S. Agrawal, and M. Adeyemi. 2022. [Quality at a glance: An audit of web-crawled multilingual datasets.](#) *TACL*, 10:50–72.
- Kudo, T. 2018. [Subword regularization: Improving neural network translation models with multiple subword candidates.](#) *ACL*.
- Kudo, T. and J. Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing.](#) *EMNLP*.
- Kumar, S. and W. Byrne. 2004. [Minimum Bayes-risk decoding for statistical machine translation.](#) *HLT-NAACL*.
- Lan, Z., M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations. *ICLR*.
- Lison, P. and J. Tiedemann. 2016. [Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles.](#) *LREC*.
- Lowerre, B. T. 1976. *The Harpy Speech Recognition System*. Ph.D. thesis, Carnegie Mellon University, Pittsburgh, PA.
- Marcu, D. and W. Wong. 2002. [A phrase-based, joint probability model for statistical machine translation.](#) *EMNLP*.
- Marie, B., A. Fujita, and R. Rubino. 2021. [Scientific credibility of machine translation research: A meta-evaluation of 769 papers.](#) *ACL*.
- McLuhan, M. 1964. *Understanding Media: The Extensions of Man*. New American Library.
- Miller, G. A. and J. G. Beebe-Center. 1956. [Some psychological methods for evaluating the quality of translations.](#) *Mechanical Translation*, 3:73–80.
- Nirenburg, S., H. L. Somers, and Y. Wilks, eds. 2002. *Readings in Machine Translation*. MIT Press.
- Och, F. J. 1998. *Ein beispiebsbasierter und statistischer Ansatz zum maschinellen Lernen von natürlichsprachlicher Übersetzung*. Ph.D. thesis, Universität Erlangen-Nürnberg, Germany. Diplomarbeit (diploma thesis).
- Och, F. J. 2003. [Minimum error rate training in statistical machine translation.](#) *ACL*.
- Och, F. J. and H. Ney. 2002. [Discriminative training and maximum entropy models for statistical machine translation.](#) *ACL*.
- Och, F. J. and H. Ney. 2003. [A systematic comparison of various statistical alignment models.](#) *Computational Linguistics*, 29(1):19–51.
- Och, F. J. and H. Ney. 2004. [The alignment template approach to statistical machine translation.](#) *Computational Linguistics*, 30(4):417–449.
- Papineni, K., S. Roukos, T. Ward, and W.-J. Zhu. 2002. [Bleu: A method for automatic evaluation of machine translation.](#) *ACL*.
- Pierce, J. R., J. B. Carroll, E. P. Hamp, D. G. Hays, C. F. Hockett, A. G. Oettinger, and A. J. Perlis. 1966. *Language and Machines: Computers in Translation and Linguistics*. ALPAC report. National Academy of Sciences, National Research Council, Washington, DC.
- Popović, M. 2015. [chrF: character n-gram F-score for automatic MT evaluation.](#) *Proceedings of the Tenth Workshop on Statistical Machine Translation*.
- Post, M. 2018. [A call for clarity in reporting BLEU scores.](#) *WMT 2018*.
- Prates, M. O. R., P. H. Avelar, and L. C. Lamb. 2019. Assessing gender bias in machine translation: a case study with Google Translate. *Neural Computing and Applications*, 32:6363–6381.
- Raffel, C., N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 21(140):1–67.
- Rei, R., C. Stewart, A. C. Farinha, and A. Lavie. 2020. [COMET: A neural framework for MT evaluation.](#) *EMNLP*.
- Schiebinger, L. 2014. Scientific research must take gender into account. *Nature*, 507(7490):9.
- Schwenk, H. 2018. [Filtering and mining parallel data in a joint multilingual space.](#) *ACL*.
- Schwenk, H., D. Dechelotte, and J.-L. Gauvain. 2006. [Continuous space language models for statistical machine translation.](#) *COLING/ACL*.
- Schwenk, H., G. Wenzek, S. Edunov, E. Grave, A. Joulin, and A. Fan. 2021. [CCMatrix: Mining billions of high-quality parallel sentences on the web.](#) *ACL*.
- Sellam, T., D. Das, and A. Parikh. 2020. [BLEURT: Learning robust metrics for text generation.](#) *ACL*.

- Slobin, D. I. 1996. Two ways to travel. In M. Shibatani and S. A. Thompson, eds, *Grammatical Constructions: Their Form and Meaning*, 195–220. Clarendon Press.
- Snover, M., B. Dorr, R. Schwartz, L. Micciulla, and J. Makhoul. 2006. A study of translation edit rate with targeted human annotation. *AMTA-2006*.
- Stanovsky, G., N. A. Smith, and L. Zettlemoyer. 2019. [Evaluating gender bias in machine translation](#). *ACL*.
- Stolcke, A., Y. Konig, and M. Weintraub. 1997. [Explicit word error minimization in N-best list rescoring](#). *EUROSPEECH*, volume 1.
- Suzgun, M., L. Melas-Kyriazi, and D. Jurafsky. 2023. [Follow the wisdom of the crowd: Effective text generation via minimum Bayes risk decoding](#). *Findings of ACL 2023*.
- Talmy, L. 1985. Lexicalization patterns: Semantic structure in lexical forms. In T. Shopen, ed., *Language Typology and Syntactic Description, Volume 3*. Cambridge University Press. Originally appeared as UC Berkeley Cognitive Science Program Report No. 30, 1980.
- Talmy, L. 1991. [Path to realization: A typology of event conflation](#). *BLS-91*.
- Thompson, B. and P. Koehn. 2019. [Vecalign: Improved sentence alignment in linear time and space](#). *EMNLP*.
- Turian, J. P., L. Shen, and I. D. Melamed. 2003. Evaluation of machine translation and its evaluation. *Proceedings of MT Summit IX*.
- Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. 2017. [Attention is all you need](#). *NeurIPS*.
- Vauquois, B. 1968. A survey of formal grammars and algorithms for recognition and transformation in machine translation. *IFIP Congress 1968*.
- Weaver, W. 1949/1955. Translation. In W. N. Locke and A. D. Boothe, eds, *Machine Translation of Languages*, 15–23. MIT Press. Reprinted from a memorandum written by Weaver in 1949.
- Wu, D. 1996. [A polynomial-time algorithm for statistical machine translation](#). *ACL*.
- Wu, Y., M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, Ł. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. S. Corrado, M. Hughes, and J. Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. ArXiv preprint arXiv:1609.08144.
- Zens, R. and H. Ney. 2007. [Efficient phrase-table representation for machine translation with applications to online MT and speech translation](#). *NAACL-HLT*.
- Zhang, T., V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi. 2020. [BERTscore: Evaluating text generation with BERT](#). *ICLR 2020*.
- Ziemski, M., M. Junczys-Dowmunt, and B. Pouliquen. 2016. [The United Nations parallel corpus v1.0](#). *LREC*.