

If you're still with me, thank you. I hope it was worth it. This is the end of our journey. We started with a goal of making a speech-to-text converter for our imaginary voice-controlled computer. In the process, we started from the most basic building blocks, counting and arithmetic, and reconstructed a transformer from scratch. My hope is that the next time you read an article about the latest natural language processing conquest, you'll be able to nod contentedly, having pretty darn good mental model of what's going on under the hood. Resources and credits ► The O.G. [paper](https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf), Attention is All You Need. \blacktriangleright A wildly helpful Python [implementation](https://nlp.seas.harvard.edu/2018/04/03/attention.html) of the transform

[Brandon](https://e2eml.school/) October 29, 2021

Tokenizing We made it all the way through the transformer! We covered it in enough detail that there should be no mysterious black boxes left. There are a few implementation details that we didn't dig into. You would need to know about them in order to build a working version for yourself. These last few tidbits aren't so much about how transformers work as they are about getting neural networks to behave well. [The Annotated Transformer](https://nlp.seas.harvard.edu/2018/04/03/attention.html) will help you fill in these gaps. We are not completely done yet though. There are still some important things to say about how we represent the data to start with. This is a topic that's close to my heart, but easy to neglect. It's not so much about the power of the algorithm as it is about thoughtfully interpreting the data and understanding what it means. We mentioned in passing that a vocabulary could be represented by a high dimensional one-hot vector, with one element associated with each word. In order to do this, we need to know exactly how many words we are going to be representing and what they are. A naïve approach is to make a list of all possible words, like we might find in Webster's Dictionary. For the English language this will give us several tens of thousands, the exact number depending on what we choose to include or exclude. But this is an oversimplification. Most words have several forms, including plurals, possessives, and compared can have alternative spellings. And unless your data has been very word to be carefully cleaned, it will contain typographical errors of all sorts. This doesn't even touch on the possibilities opened up by freeform text, neologisms, slang, jargon, and the vast universe of Unicode. An exhaustive list of all possible words would be infeasibly long. A reasonable fallback position would be to have individual characters serve as the building blocks, rather than words. An exhaustive list of characters is well within the capacity we have to compute. However there are a couple of problems with this. After we transform data into an embedding space, we assume the distance in that space has a sema interpretation, that is, we assume that points that fall close together have similar meanings, and points that are far away mean something very different. That allows us to implicitly extend what we learn about one word to its immediate neighbors, an assumption we rely on for computational efficiency and from which the transformer draws some ability to generalize. At the individual character level, there is very little semantic content. There are a few one character words in the English language for example, but not many. Emoji are the exception to this, but they are not the primary content of most of the data sets we are looking at. That leaves us in the unfortunate position of having an unhelpful embedding space. It might still be possible to work around this theoretically, if we could look at rich enough combinations of characters to build up semantically useful sequences like words, words stems, or word pairs. Unfortunately, the features that transformers create internally behave more like a collection of input pairs than an ordered set of inputs. That means that the representation of a word would be a collection of character pairs, without their order strongly represented. The transformer would be forced to continually work with anagrams, making its job much harder. And in fact experiments with character level representations have shown the transformers don't perform very well with them. Byte pair encoding Fortunately, there is an elegant solution to this. Called [byte pair encoding](https://en.m.wikipedia.org/wiki/Byte_pair_encoding). Starting with the character level representation, each character is assigned a code, its own unique byte. Then after scanning some representative data, the most common pair of bytes is grouped together and assigned a new byte, a new code. Ths new code is substituted back into the data, and the process is repeated. Codes representing pairs of characters can be combined with codes representing other characters or pairs of characters to get new codes representing longer sequences of characters. There's no limit to the length of character sequence a code can represent. They will grow as long as they need to in order to represent commonly repeated sequences. The cool part of byte pair encoding is that in infers which long sequences of characters to learn from the data, as opposed to dumbly representing all possible

sequences. it learns to represent long words like *transformer* with a single byte code, but would not waste a code on an arbitrary string of similar length, such as *ksowjmckder*. And because it retains all the byte codes for its single character building blocks, it can still represent weird misspellings, new words, and even foreign languages. When you use byte pair encoding, you get to assign it a vocabulary size, ad it will keep building new codes until reaches that size. The vocabulary size needs to be big enough, that the character strings get long enough to capture the semantic content of the the text. They have to mean something. Then they will be sufficiently rich to power transformers After a byte pair encoder is trained or borrowed, we can use it to pre-process out data before feeding it into the transformer. This breaks it the unbroken stream of text into a sequence of distinct chunks, (most of which are hopefully recognizebable words) and provides a concise code for each one. This is the process called tokenization. $\mathcal{A}^{\text{max}}_{\text{max}}$ Now recall that our original goal when we started this whole adveture was to translate from the audio signal or a spoken command to a text representation. So far all of our examples have been worked out with the assumption that we are working with characters and words of written language. We can extend this to audio too, but that will take an even bolder foray into signal preprocessing. The information in audio signals benefits from some heavy-duty preprocessing to pull out

the parts that our ears and brains use to understand speech. The method is called Melfrequecy cepstrum filtering, and it's every bit as baroque as the name suggests. Here's a well-illustrated [tutorial](http://www.speech.cs.cmu.edu/15-492/slides/03_mfcc.pdf) if you'd like to dig into the fascinating details. When the pre-processing is done, raw audio is turned into a a sequence of vectors, where each element represents the change of audio activity in a particular frequency range. It's dense (no elements are zero) and every element is real-valued. On the positive side, each vector makes a good "word" or token for the transformer because it means something. It can be directly translated into a set of sounds that is recognizeable as part of a word. On the other hand, treating each vector as a word is weird because each one is unique. It's extremely unlikely that the same set of vector values will ever occur twice, because there are so many subtly different combination of sounds. Our previous strategies of onehot representation and byte pair encoding are of no help. The trick here is to notice that dense real-valued vectors like this is what we end up with *after* embedding words. Transformers love this format. To make use of it, we can use the results of the ceptrum pre-processing as we would the embedded words from a text example. This saves us the steps of tokenization an embedding. It's worth noting that we can do this with any other type of data we want too. Lots of recorded data comes in the form of a sequence of dense vectors. We can plug them right in to a transformer's encoder as if they were embedded words. Wrap up

► Jay Alammar's insightful transformer [walkthough.](https://jalammar.github.io/illustrated-transformer/) ► Lukasz Kaiser's (one of the authors) [talk](https://www.youtube.com/watch?v=rBCqOTEfxvg) explaining how transformers work. ► The [illustrations](https://docs.google.com/presentation/d/1Po-GY7X-mXmPKHr8Vh29S4tFPv23TjeY-jq-yShlivM/edit?usp=sharing) in Google Slides.

[All text and images CC0 except where noted.](https://creativecommons.org/publicdomain/zero/1.0/) The opinions here are wholly my own.