



DEEP LEARNING

Master Positional Encoding • Part I

We present a "derivation" of the fixed positional encoding used in Transformers, helping you get a full intuitive understanding.

Jonathan Kernes

Feb 15, 2021 21 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

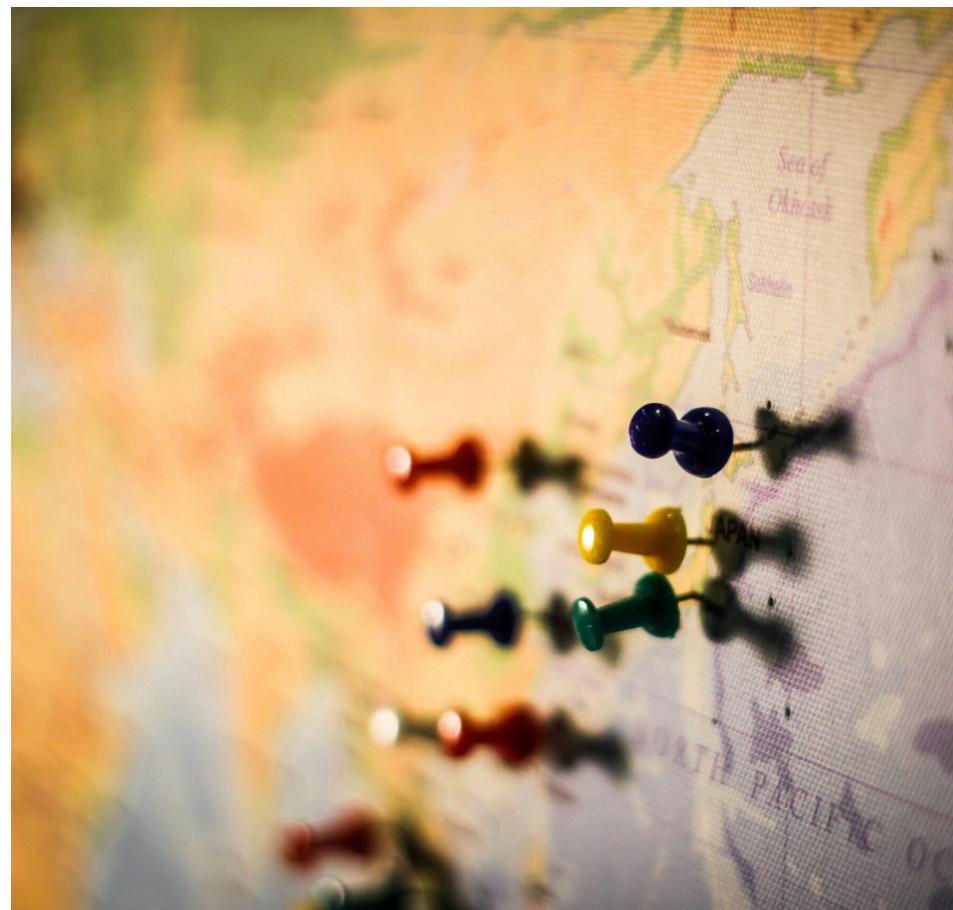


Photo by [T.H. Chia](#) on [Unsplash](#)

This is Part I of two posts on positional encoding (Part II is now available [here!](#)):

- **Part I:** the intuition and "derivation" of the fixed sinusoidal positional encoding.
- **Part II:** how do we, and how *should* we actually inject positional information into an attention model (or any other model that may need a positional embedding). This focuses on relative encodings.

As far as I know, there are some original insights in the field of positional encoding viewed as a manifold. But, the field is young, so I doubt I'm the first to notice this; if anyone knows more, please let me know!

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Part I: the fixed sinusoidal position encoding

Introduction

Imagine you've finally managed to work your way through understanding [Transformers](#). You've drawn those weight matrices mapping keys to queries, and convinced yourself of the meaning of the logits tensor. You've even gone a step further and figured out how multi-head attention works, admiring the efficiency that a little tensor reshaping can provide. You learned how to mask your input sequences by directly manipulating logits and how to gracefully transfer sequences from encoder to decoder. You're so satisfied with yourself you place yourself on the couch, throw your feet up and treat yourself to a nice cookie... No. Two chocolate chip cookies. You're

Contented with your newfound knowledge, you fire up the embedding layer. "It's easy," you muse to yourself, turning words into vectors since the backstreet boy in you turns your eyes towards the bottom of the model diagram, something unfamiliar catches your gaze. The **position encoding** layer. You immediately stop eating your

falling to the floor, as you suddenly remember. *Positional information doesn't show up anywhere.*

Not yet panicking, you read ahead hoping you'll find see an easy solution. Instead all you get are a few sentences telling you about some fixed sinusoid function and the directions that when referring to the usual embedding layer "...the two can be summed." *Furious*, turning the page back and forth, you realize there [LATEST](#) Having taken a physics class once before, you know what happened. You've been "it's trivial" ed.

Positional encoding is basically the "it's trivial" one of many dreaded 2–3 word phrases like "it's not the point", where the intent is most definitely the opposite of what you're saying. It's *not* fine, it *does* matter, and it's *not* trivial.

Now, that we all acknowledge the situation, let's figure out what positional encoding is. We are going to "deriving" it by taking a physicist's approach. Come up with a few theories/guesses for how to accomplish our task, test them upon them until we get an encoding that mostly does what we want. We will put off where and how to pop in the encoding layer to Part II. For now, we are just trying to figure out how to sequence some positional meaning.

Discrete positional encoding

The most difficult part of this might just be trying to figure out exactly that we're trying to do. **If this subsection makes no sense, just directly skip to the next one.** / The following goal:

[EDITOR'S PICKS](#)
[DEEP DIVES](#)

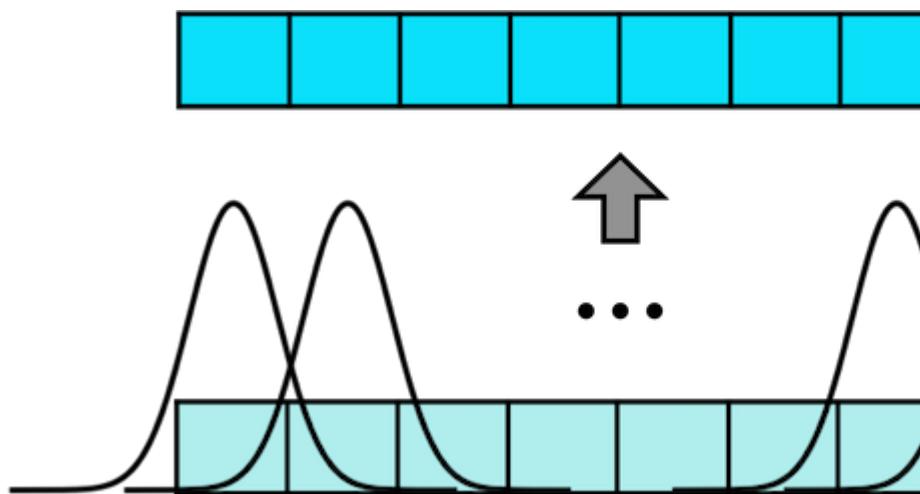
[CONTRIBUTE](#)
[NEWSLETTER](#)
[Sign in](#)

[Contributor Portal](#)

A positional encoding is a finite dimensional representation of the location or "position" of items in a sequence. Given some sequence $A = [a_0, \dots, a_{n-1}]$, the positional encoding must be some type of tensor that we can feed to a model to tell it where some value a_i is in the sequence A.

Admittedly, this is vague; this is because the "feed to a model" line is also vague. To fix this, let's assume we want to try the following task:

Given a sequence A of length T, output another length T where the output for each index is the it's neighbors.



Source: the current author. The bottom layer is our input encoding vector. The over we feed this into, producing an output vector of the same size. Each dark blue is average of its lighter counterpart.

This model is meaningless. I repeat. This model is just to give us an idea of an example of one way you can add positional encoding to a neural network.

We are looking for a positional encoding *tensor* not an encoding *network*. This means that that our position tensor should be a representative of the position.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

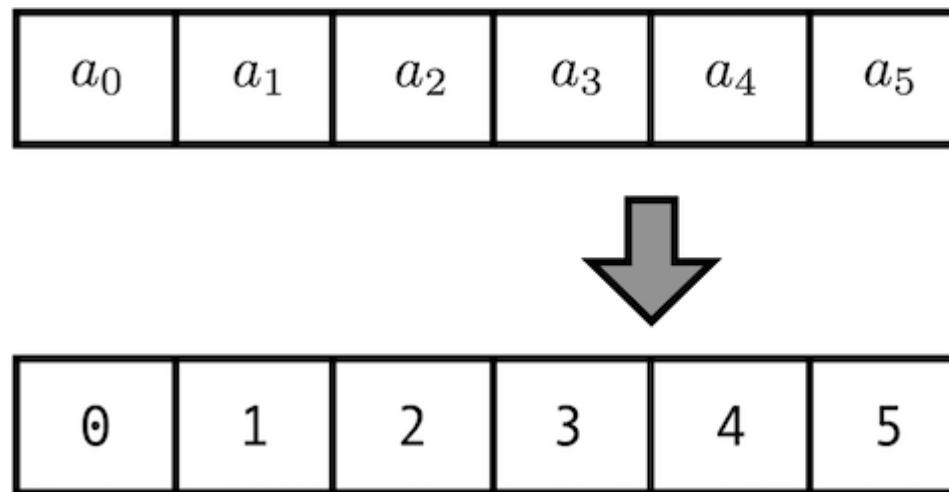
our network will take care of whatever operations need to be done to that.

This encoding is going to have the same dimension as the sequence length, which could be very long. Also, the sequence length is not even a fixed quantity. Feeding a variable length vector into a model is a big problem. We assume that there is some maximum value called the **max_position** for which our model just can't handle. For all of our sequences. For all of our solutions we fix our encoding and truncate down as needed for inputs.

Hopefully everything will be clear with an example.

Guess #1: just count.

Let's start with the simplest positional encoding we can create. We just count! We map every element in the sequence to the following encoding:



Source: current author. Sequence_length = 8. We simply create a positional encoding for every number.

Pretty basic, we just created a new vector where each number. This is an **absolute positional encoding**. the sequence is labeled by its position relative to the coordinates, assumed to be the start of the current

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

continue to work with absolute encodings, as they can be used to derive **relative encodings**.

Now let's be critical. Here's what's wrong with our initial guess: the scale of these numbers is huge. If we have a sequence of 500 tokens, we'll end up with a 500 in our vector. In general, neural nets like their weights to hover around zero, and usually be equally balanced positive and negative. If not, you open yourself up to all sorts of problems, like exploding gradients and un

LATEST
EDITOR'S PICKS

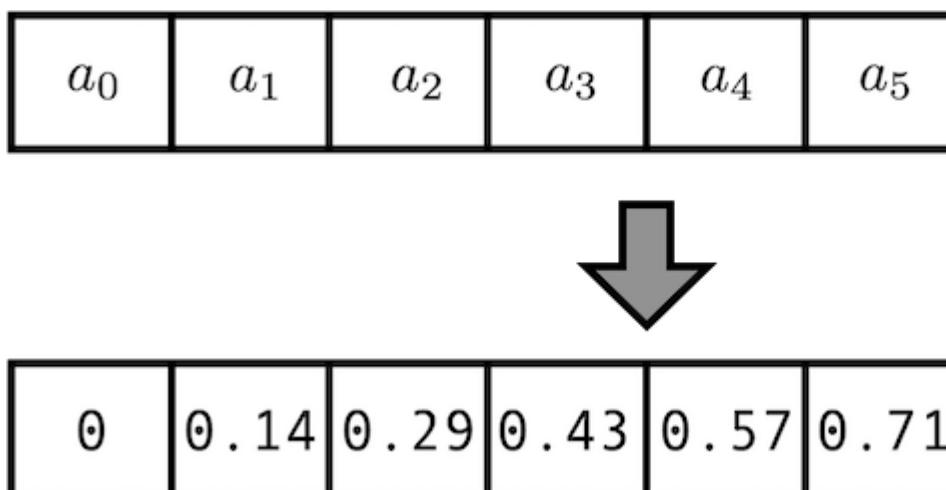
DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



Source: current author. We normalize each entry by dividing by the sequen

But oh no! we've introduced another problem. You deal with arbitrary sequence length. That's because each entry is divided by the sequence length. A position of say 0.8 means a totally different thing to a sequence of length 5 than it does for one of length 20. (For length 5, 0.8 would be the 4th element. For sequence length 20, 0.8 represents the 16th element!). **Naively normalizing work with variable sequence lengths.**

What we need to be able to do, is find a way to count without ever using numbers greater than 1. That's a big clue, and you may be able to guess the answer.

Use binary numbers.

Guess #3: just count! but using binary instead of decimal

Instead of writing say 35 for the 35th element, we represent it via its binary form 100011. Cool, everyone. But, in our excitement, we forgot that the number of elements in our vector doesn't matter whether we use binary or decimal; it's always just gain anything... Instead, we actually need to do two things: 1) convert our integer to binary, and 2) convert our scalar into a vector.

Whoa, scalar to a vector? What does that mean? I positional encoding *vector* now becomes a position. Each number gets its own **binary vector** as opposed to a single scalar value. This is a pretty big conceptual shift, but by increasing dimensionality, we are able to both keep arbitrary numbers while also restricting numbers to the range [0,1]! Our vector looks like this:

LATEST

EDITOR'S PICKS

DEEP DIVES

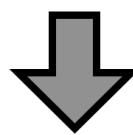
CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

a_0	a_1	a_2	a_3	a_4	a_5	a_6	a_7
-------	-------	-------	-------	-------	-------	-------	-------



value: 0 1 2 3 4 5 6 7

d_{model}	0	0	0	0	1	LATEST
	0	0	1	1	0	EDITOR' S PICKS
	0	1	0	1	0	DEEP DIVES

sequence length

Source: current author. The matrix version of positional encoding. Each binary vector is a binary number. e.g. index 5 holds 100 in binary, which is the value 4.

LATEST
EDITOR' S PICKS
DEEP DIVES
CONTRIBUTE
NEWSLETTER
[Sign in](#)
[Contributor Portal](#)

We have some freedom. We can choose the dimension of the embedding space that holds the binary vector to be any power of 2. The only thing this dimension does is tell us how many bits we need. For simplicity, let's call it d_{model} , the embedding dimension that our network uses. To make sure we understand the vector representation, let's say we want to represent the word "apple" as a binary vector. We would be represented as the binary vector:

35 \longleftrightarrow [0, ..., 0, 1, 0, 0, 0, 1, 1]

Analysis time. What are some problems with our current approach?

1. We still haven't fully normalized. Remember that the values in the vector should be positive and negative roughly equally. This means we need to rescale [0,1] \rightarrow [-1,1] via $f(x) = 2x - 1$.

2. Our binary vectors come from a discrete function, and not a discretization of a continuous function.

The latter point is subtle, and it has to do with using a *vector* to represent a *scalar* position. Let's compare two encodings, that try to measure a distance from 0 → 3. The first encoding will be really simple: only use the x-axis dimension and put the position there represented by a variable x in [0,3]

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

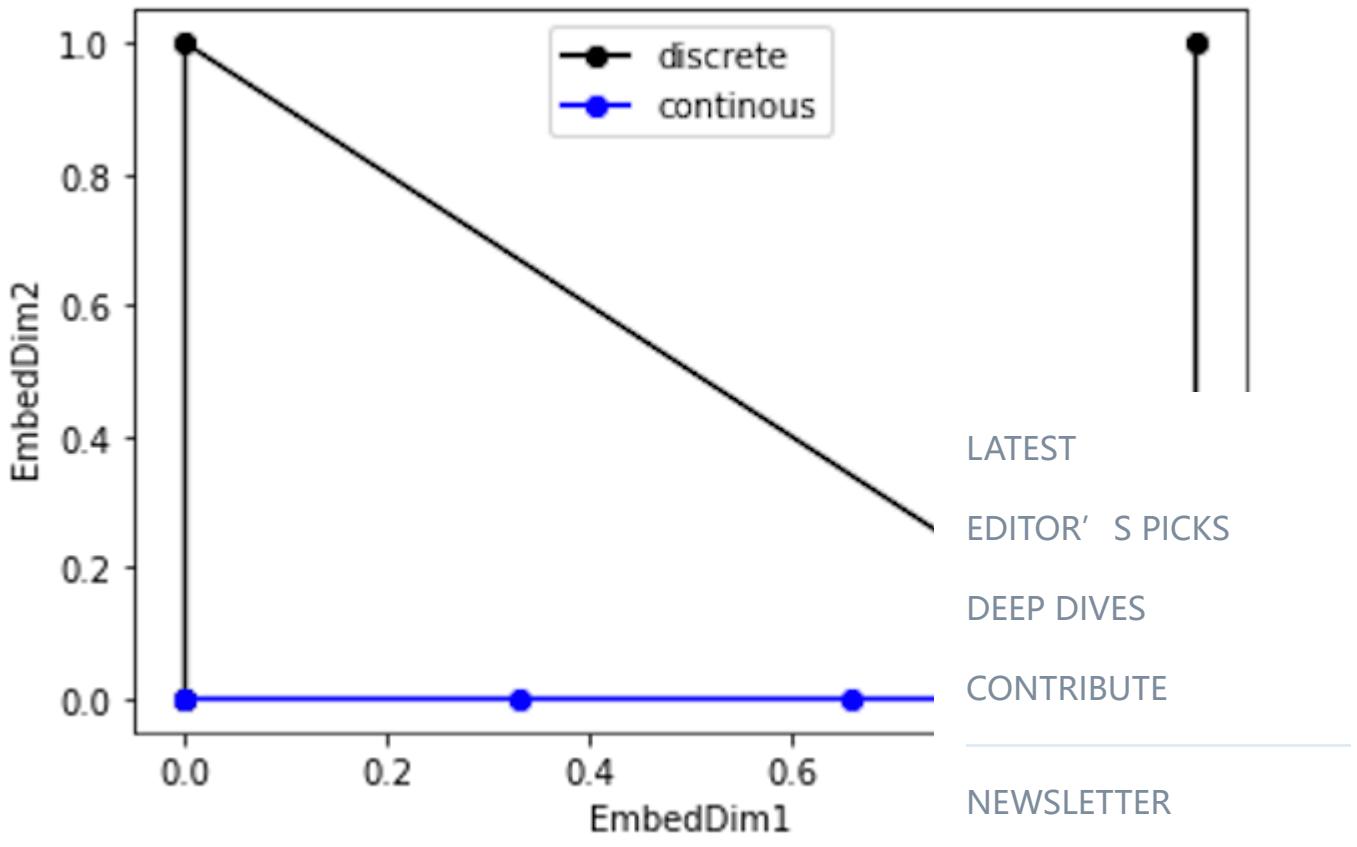
Sign in

Contributor Portal

Now let's compare it to our binary vector discretization. This is discrete, we can only represent 4 positions: give 4 points:

Binary vector encoding: [0,0], [0,1], [1,0], [1,1]

Any time we discretize something, we have inclined to think that something can be **interpolated**, meaning that we can draw a continuous function by connecting the dots. Take a look at the interpolations for our two encodings:



Source: current author. The space curve showing flow of position. The discrete curve is a step function from $(0,1) \rightarrow (1,1)$, and is much more jagged than the smooth continuous curve.

You can see that the continuous encoding curve is smooth and easy to estimate intermediate positions. For the binary vector, things are a mess. This will lead us to our next hint: **1. the binary vector a discretization of something**

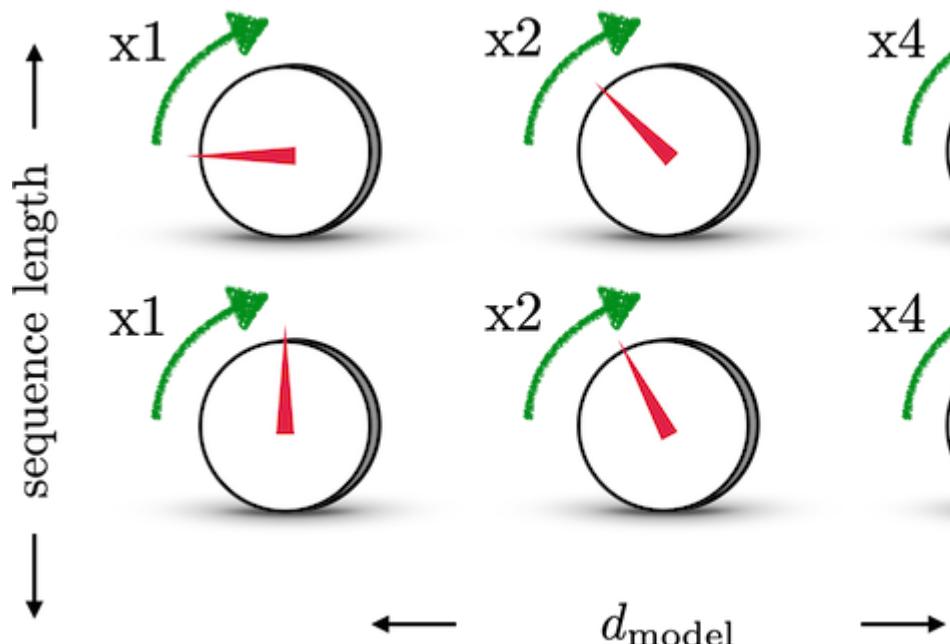
So with pictures, we want a function that connects two points in a smooth way that looks natural. For anyone who has done calculus, you know what we are really doing is finding an embedding function.

Yes, if you use positional encoding, you can sound like you're telling people that you're using manifolds. What this means, in terms, is that we are looking for a curve in the $\text{dim}_1 \times \text{dim}_2$ space, that as you walk along it slowly increases your EmbedDim2 value in a smooth, continuous way. Doing this properly brings us to our next hint:

Guess #4: use a *continuous* binary vector

To make our binary vector continuous, we need a function that interpolates a back and forth cycle $0 \rightarrow 1 \rightarrow 0 \rightarrow 1 \rightarrow 0 \dots$ etc. Cycles you say? Time to bring in trig functions. Let's use a sine function. Also, a sine function exists on $[-1,1]$ so it's also properly normalized, which is an added bonus!

Here's the image/analogy you should now have. Every matrix entry in the positional encoding matrix represents a "dial". Typically your dial might do something like adjust volume (volume is analogous to position). Like normal dials we use in the encoding matrix can continuously go on (1). To get a precise volume (position), we use a vector (vector) of dials with different sensitivity. The first will adjust volume ever so slightly, perhaps by 1 unit, a barely noticeable change. The second will be more powerful, adjusting volume by 2 units. The third will adjust by 4 units, the fourth by 8 units, etc.



Source: current author. Each position in the sequence (column) is represented by a vector of dials. This visualization shows how the sensitivity of each dial changes over the sequence length. As the sequence length increases, the sensitivity of the dials increases exponentially, allowing for more precise control over the final output.

In this way, you can get a precise volume over various dimensions. The strength increases exponentially. Instead of building a sequence that contains say 512 sound levels, we can build 8

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

as much precision ($2^8=256$ if you're wondering why 8 dials). Super efficient huh!

NOTATION WARNING: we transposed the positional encoding matrix relative to the previous section. We now have dim M = (sequence_length, d_model), which is the usual way this is written.

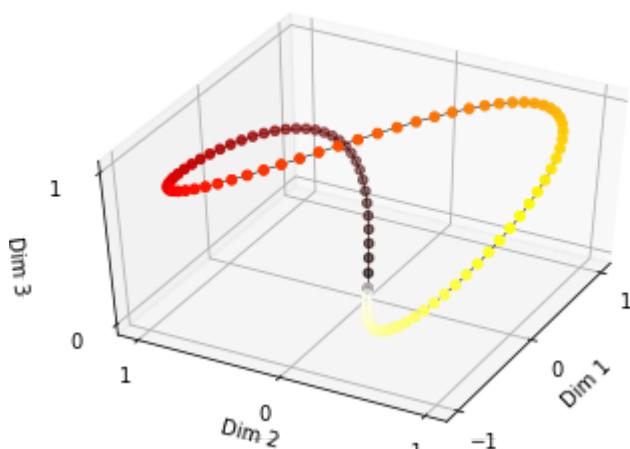
We now want our sine functions to undergo one full cycle over the right moments. For the first entry, it should start at zero every time we move one step in the sequence. The frequency of $\pi/2$ for the first dial, $\pi/4$ for the second, $\pi/8$ for the third, and on and on...

We now have our first prototype positional encoding of a **matrix denoted by M**, where the y-axis of the matrix corresponds to the discrete position $_xi$ in the sequence (0, 1, ..., n-1), and the x-axis corresponds to the vector embedding dimension. Elements of the matrix are given by the formula:

$$M_{ij} = \sin(2\pi i / 2^j) = \sin(x_i \omega_j)$$

Here, ω is the dial frequency, which is monotonically decreasing with respect to the dimension index j , giving the sequence position.

Let's now compare our continuous-discretized encoding with a plain old discretized encoding.



Source: current author. The colors indicate the absolute position that that coordinate has in each dimension. For example, $\omega=pi/2, pi/4, pi/8$ for dimensions 1, 2 and 3 respectively.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Look, it's all continuous-y and manifold-y.

Analysis time. We are very close to the final answer. There are two things we still want to fix.

First, there is a discontinuity problem. The above plot is a smooth curve like we wanted. However, it is also a closed curve. Consider moving from position n , to $n+1$. This would be the same as moving from position 0 to position d_{model} . The curve transitions yellow \rightarrow brown. By construction, this is equivalent to position 1. This was necessary to make the curve continuous, but it's also not what we want! $n+1$ is far apart, not close together. For anyone who has studied topology or transforms, this problem should be familiar. Only frequencies are unique.

To fix this, we observe that there's no reason that our embedding vector line up with a binary vector of frequencies ($\omega_0, \omega_1, \dots, \omega_{d_{\text{model}}}$). As long as they are *monotonically increasing*, which they will increase, and only increase. We can lower all of them this way keep ourselves far away from the boundaries.

$$M_{ij} = \sin(x_i \omega_0^{j/d_{\text{model}}})$$

The slowed down dial version. ω_0 is the min frequency. Check the edge case $j=d_{\text{model}}$ gives the smallest.

The original authors choose a minimum frequency for this purpose. they then increase frequencies exponentially with $\omega_{\min}^{\{range(d)/d\}}$ to achieve the monotonicity. Why choose 10,000? I truly have no idea, I'm guessing it made this work the best.

The second problem isn't a problem insomuch as it's a limitation. There is no simple operation that allows us to add

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

positional units from a position vector. Such an operation, known as translation, would allow us to write the positional encoding vector of index i as a function of the vector for index j .

Why is this a very desirable property to have? Imagine we have a network that is trying to translate the sentence "I am going to eat." The combination "is/am/are" + "going" + "to" + "verb" is a very common grammatical structure, with a fixed positional structure. It ends up at index 1, "to" at index 2, etc. In this case, if we want the network to learn to pay attention to the word "I" that occurs before "am going to". "I" is located 4 units to the left of "verb". Since our attention layer uses linear transformations to form the keys, queries, and values, it would be nice to have the positional encodings such that the linear transformation for the position vector located 4 units to the left, so that it matches the position vector of "verb". The query and key vectors would then align perfectly.

In summary, we want to see if we can modify our positional encoding so that it can be translated via linear transformation.

Final guess, #5: Use both sine and cosine to make the matrix periodic

The math is going to seem a little heavy, but it's not too bad. We just need to keep track of notation. First, let's suppose we have a sequence of length seq_len . We can represent the positional encoding matrix as:

$$PE = \begin{pmatrix} \mathbf{v}^{(0)} \\ \vdots \\ \mathbf{v}^{(seq_len-1)} \end{pmatrix}$$

The row-vector \mathbf{v} is a vector of sines evaluated at different frequencies, but with varying frequencies

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

$$\mathbf{v}^{(i)} = [\sin(\omega_0 x_i), \dots, \sin(\omega_{n-1} x_i)]$$

Each row-vector represents the positional encoding vector of a single, discrete position. The positional encoding matrix **PE**, is a vector of these vectors, and so it is a (seq_len, n) dimensional matrix (seq_len is the sequence length). We now want to find a linear transformation **T(dx)** such that the following equation holds:

$$\text{PE}(x + \Delta x) = \text{PE}(x) \cdot \mathbf{T}(\Delta x)$$

How will we do this? First, notice that since all of the values are sines, the positions x are actually *angles*. From this, we know that any operation \mathbf{T} that shifts the argument must be some kind of rotation. Rotations can be represented by applying a linear transformation to a (cosine, sine) pair. Using standard results for rotation matrices, we can make the following identity:

$$\begin{pmatrix} \cos(\theta + \phi) \\ \sin(\theta + \phi) \end{pmatrix} = \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix} \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix}$$

to build a rotation.

Note, this formula has the linear operator acting on the right as we have so far shown. It's conceptually easier to think of the left operation so we will write that answer here. Keep in mind that you will need to take the transpose of the current **PE** matrix.

To build an encoding that supports translations via linear transformations, we do the following. Create a duplicate encoding instead of sine. Now, build a new positional encoding matrix as follows: alternating between the cosine and sine rows. Insert the first column and append it to the final **PE** matrix.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

replacing every $\sin(\dots)$ with a $[\cos(\dots) \sin(\dots)]$ pair. In our new **PE** matrix, the row-vectors look like the following:

$$\mathbf{v}^{(i)} = [\cos(\omega_0 x_i), \sin(\omega_0 x_i), \dots, \cos(\omega_{n-1} x_i), \sin(\omega_{n-1} x_i)]$$

We now build the full linear transformation by using a bunch of block-diagonal linear transformations. Each block will have a different matrix, since the frequencies that the block acts on are di
in order to translate the _k_th dial with frequency
we would need a total angle shift of delta=omega
can now be written as:

$$\mathbf{T}(\Delta x) = \begin{pmatrix} \begin{bmatrix} \cos(\omega_0 \Delta x) & -\sin(\omega_0 \Delta x) \\ \sin(\omega_0 \Delta x) & \cos(\omega_0 \Delta x) \end{bmatrix} & \dots \\ \vdots & \ddots \\ 0 & \dots \end{pmatrix}$$

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

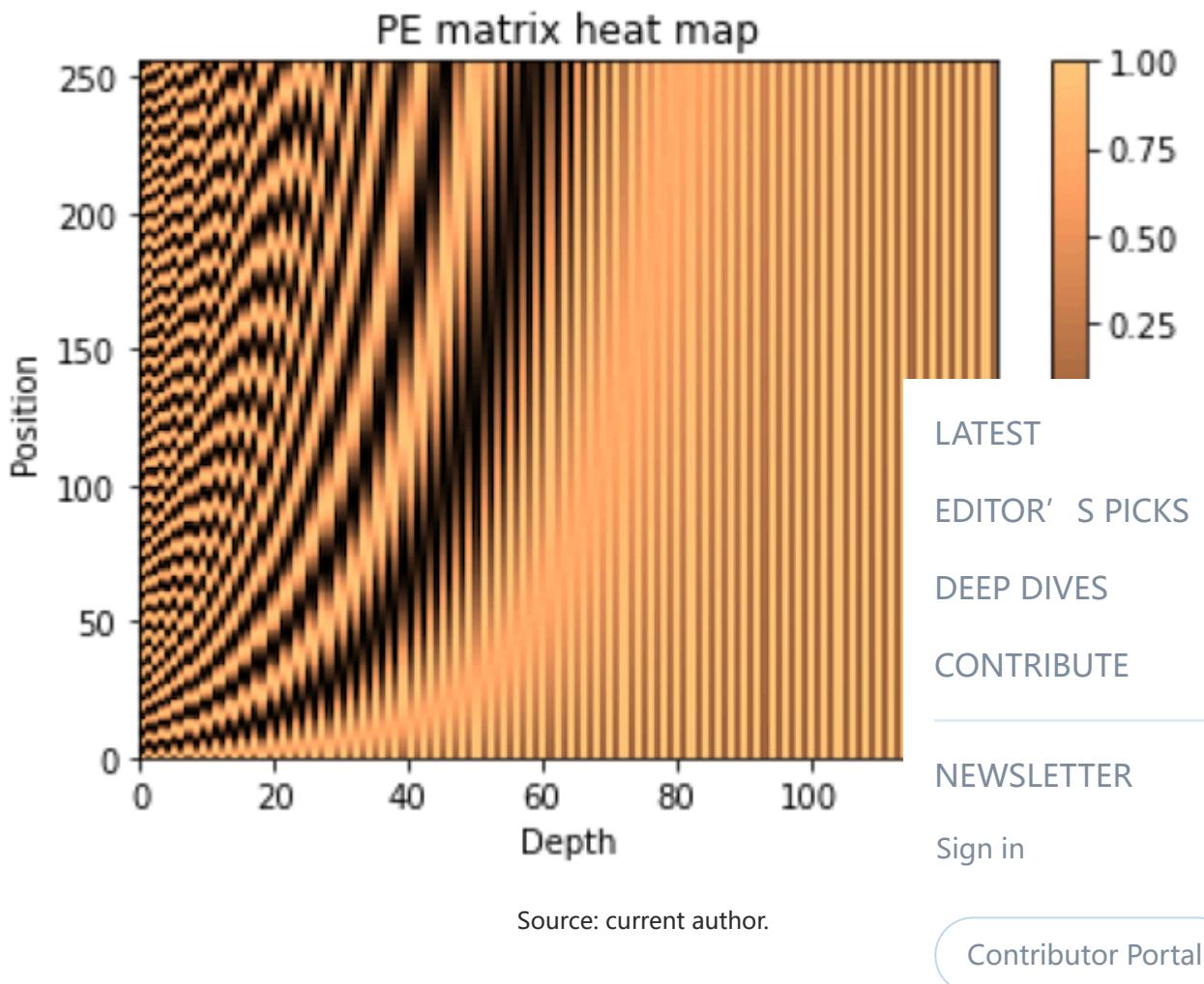
NEWSLETTER

Sign in

If you take the transpose of this, you can directly i
previous $\text{PE}(x+dx)=\text{PE}(x)*T$ equation, thereby prov
the existence of a translation matrix!

Contributor Portal

That wraps up mostly everything. Let's now take
PE matrix actually looks like.



Here's how to interpret the above PE matrix heatmap:

- 1. Most of the PE matrix is not needed for encoding.** The plot that are darker indicate dials that have been activated, meaning that the sines and cosines in that region are moving away from their initial values of 0 and 1 respectively. This shows how much of the embedding space is being used to store positional information. As you can see by following the curve-ish thing, activating a dial one step deeper becomes exponentially more difficult (the black part of e^x).
- 2. Vertical lines at greater depth vary less than at shallower depth.** Pick a fixed depth, then move upward. You'll see the color shifts light → dark → light... The frequency of variation *decreases* with depth, showing our intuition that positions at greater depth are more sensitive.

Final answer. No more guesses

At this point we've done everything we need to do, so let's summarize:

1. Positional encoding is represented by a matrix. For sequence length T and model depth D, this would be a (T, D) tensor.
2. The "positions" are just the indices in the sequence.
3. Each row of the PE matrix is a vector that represents the interpolated position of the discrete value assigned to the index.
4. The row-vector is an alternating series of sine and cosine frequencies that decrease according to a geometric progression.
5. There exists a matrix multiplication that can sample any row-vector we want.

Here are two ways to create a positional encoding: one using numpy and one using only TensorFlow operators. The numpy version is cleaner and easier to understand, but the TensorFlow version is more efficient. The TensorFlow version uses broadcasting to create the PE matrix angles via array slices of the angle matrix, using subscripts, which makes it more difficult to understand.

A couple of notes to understand the code. First we use broadcasting to create the PE matrix angles via array slices of the angle matrix, using subscripts, which makes it more difficult to understand.

Bonus property

There is an additional property that we would like the positional encoding to satisfy, and it is related to how positions are used in an actual attention model. Attention is performed such that each embedding associated with each position is

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

has some representative key or query. These keys and queries allow us to quickly determine if the current position should "attend" on another position, by performing the dot product $\text{query}_i @ \text{key}_j$. If the dot product is large, then it implies the key at j matches the query at i . In order for this to work with positional encoding vectors, we would like that positions close to each other return large key-query dot products, while those far away return small ones.

figure below:

LATEST

EDITOR'S PICKS

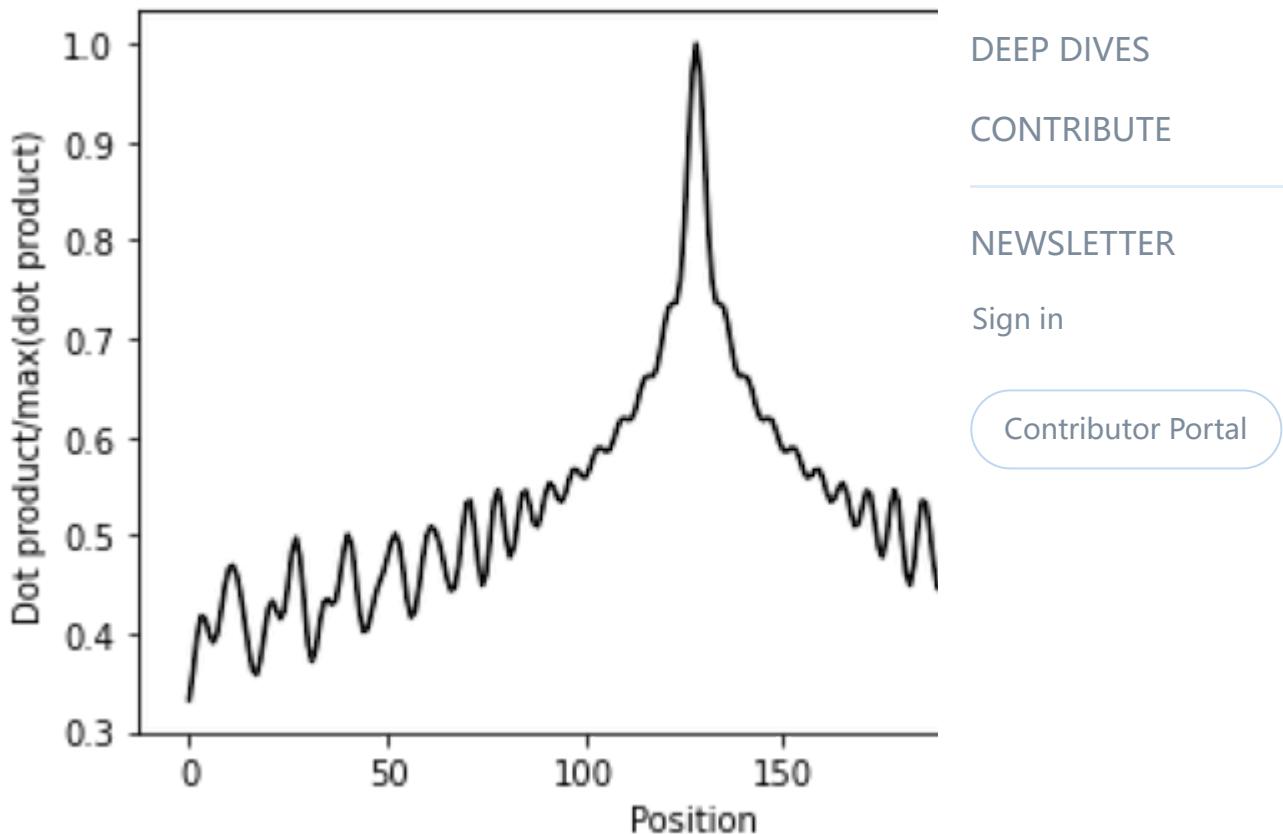
DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



Source: current author. We show the dot product of vector $v^{(128)}$ with all other vectors in the sequence with parameters $d_{\text{embed}}=128$, $\text{max_position}=256$. Dot products are nor

This figure shows the dot product between a part of a positional encoding vector representing the 128th position, and all other positional encoding vectors. Notice, that without any scaling, the dot product is very large when comparing $v^{(128)}$ to itself, and steadily decreases as the position becomes further away. A nice bonus feature!

The reasons for this can be traced back to the orthogonality of sines and cosines. As we get further away, the frequencies between positional encoding vectors are no longer correlated, and the summation of random sines and cosines is known to tend to zero.

Conclusion

Thanks for sticking with it until the end. Like all things that are explained the least tend to be the but oh well.

We learned that positional encoding is a means of location of objects in a sequence into information network (or other model) can understand and use approach where we just give the model integers of position, we learned that this awful for a number numbers can be too large, there's not an obvious information, it doesn't extend well to sequences

We next made a series of guesses, acting like phys produce a working theory. By continually adjusting guesses to incorporate more desired characteristics landed on sinusoidal positional encoding matrix. vectors to represent sequence positions. It can be allows for translations, and deal with varying seq

The details and intuition behind implementation explored in a future article. For example, in that article questions like "why do we just add positional encoding to sequence embeddings?" Here, it seems like a decent answer is that the PE matrix is sparse. Only a small so it seems like the model may treat the addition

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

concatenation, reserving some part of the embedding for position, and the rest for tokens.

References

1. Vaswani, Ashish, et al. "Attention is all you need." *arXiv preprint arXiv:1706.03762* (2017).

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

WRITTEN BY

Jonathan Kernes

See all from Jonathan Kernes

Contributor Portal

Topics:

Deep Learning

Editors Pick

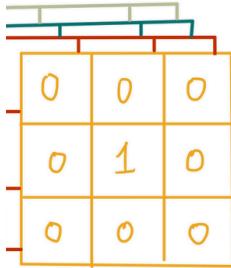
Hands On Tutorials

Transformers

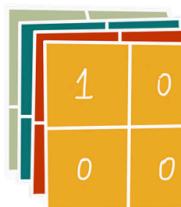
Share this article:



Related Articles



max
pooling



ARTIFICIAL INTELLIGENCE

Implementing Convolutional Neural Networks in TensorFlow

Step-by-step code guide to building a Convolutional Neural Network

Shreya Rao

August 20, 2024 6 min read



DEEP LEARNING

Deep Dive into xLSTMs by TensorFlow

Explore the world of xLSTMs - a probabilistic present-day LSTM

Srijanie Dey, PhD

July 9, 2024 13 min read

LATEST
EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

[Sign in](#)

[Contributor Portal](#)



NLP

Check Your Biases

Symbolic Engines and Unexpected Results - A Personal Coding Experience

Veronica Villa

March 17, 2022 9 min read



CINEMA

Evaluating Cinema: Which syntactical features are most important?

This article explores the relationship between a movie's genre, leveraging

Christabelle Pabal

January 20, 2024

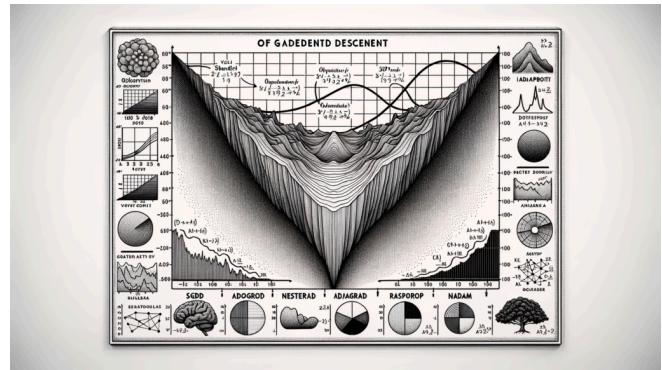
DEEP LEARNING

Speeding Up the Vision Transformer with BatchNorm

How integrating Batch Normalization in an encoder-only Transformer architecture can lead to reduced training time...

Anindya Dey, PhD

August 6, 2024 28 min read



DATA SCIENCE

The Math Behind Optimizers: Understanding SGD, Adam, and RMSprop

This is a bit different from what you might say.

Peng Qian

August 17, 2024

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



DATA SCIENCE

Latest picks: Time Series Forecasting with Deep Learning and Attention Mechanism

Your daily dose of data science

TDS Editors

November 4, 2020 1 min read



DATA SCIENCE

Stacked Ensemble Advanced Predictions With H2O.ai

And how I place the largest machine learning models in them!

Sheila Teo

December 18, 2022



MACHINE LEARNING

Large Language Models, MirrorBERT - Transforming Models into Universal Lexical and Sentence...

Discover how mirror augmentation generates data and aces the BERT performance on semantic similarity tasks

Vyacheslav Efimov

December 12, 2023 7 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



Your home for data science and AI. The world's leading publication for data analytics, data engineering, machine learning, and artificial intelligence professionals.

© Insight Media Group, LLC 2025

[ABOUT](#) • [PRIVACY POLICY](#) • [TERMS OF USE](#)

Towards Data Science is now independent!

[COOKIES SETTINGS](#)

[Sign up to our newsletter](#)

Email address*

First name*

Last name*

Job title*

Job level*

Please Select



Company name*

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Subscribe Now

- I consent to receive newsletters and other communications from Science publications.*