

# The Double Helix inside the NLP Transformer

Anonymous ACL submission

## Abstract

This study introduces a novel framework for exploring the information processing within NLP Transformers. We categorize information into four distinct layers: positional, syntactic, semantic, and contextual. Challenging the conventional integration of positional data into semantic embeddings, we propose a more effective “Linear-and-Add” method. Our analysis uncovers an intrinsic separation of positional elements in deeper layers, revealing that these components form a helix-like pattern in both encoder and decoder stages. Notably, our approach enables the identification of Part-of-Speech (PoS) clusters within conceptual dimensions. These insights offer a new perspective on information processing in the complex architecture of NLP Transformers, potentially guiding future developments in the field.

## 1 Introduction

Large Language Models (LLMs), such as ChatGPT, have become a focal point of recent research, primarily due to the Transformer architecture, which is central to all modern LLMs (Radford et al., 2018; Radford et al., 2019; Radford et al., 2019; Brown et al., 2020; Ouyang et al., 2022; OpenAI, 2023). Introduced by Vaswani et al. (2017), the Transformer employs an attention mechanism to efficiently process different parts of input data simultaneously. This has led to superior performance over older models across several NLP tasks, including machine translation and question answering. While the Transformer has been instrumental to advancements in AI, the intricacies of its function remain complex. Our work aims to demystify the Transformer's operations and provide a clear framework for its analysis.

## 2 Methodology

### 2.1 From Words to Concepts

Communication aims to transmit concepts, which are nuanced by context, a term like “server” can signify different entities in technology and food catering. Recognizing this, Transformers shift focus from words to concepts, which encapsulate multiple layers of information: positional (word location), syntactic (grammatical role), semantic (inherent meaning), and contextual (relation to surrounding words). Rather than interacting with words, Transformers process “tokens”—units that can represent word parts, punctuation, or specific syntax—allowing for versatility transcending linguistic variability.

### 2.2 The Meaning of “Meaning”

Dictionary definitions often contain circular reasoning, explaining words using other undefined words. This reveals the operational nature of “meaning” - a word's meaning is simply its relationships to other related words. However, words can be ambiguous, having multiple interconnected meanings. More universal “concepts” better capture distinct meanings. While a word like “server” has overloaded meanings, concepts require clearer explanations using multiple words and contextual understanding. As another example, “project” as a noun or verb carries different connotations.

We can visualize concepts as “mini-galaxies” of meaningful words orbiting a central point. Tracking this conceptual center is more practical than tracing all orbital word associations. Using concepts enables language translation by mapping between conceptual vector representations across languages. In embedding spaces like Word2Vec (Mikolov et al., 2013), concepts manifest as weighted superpositions of word vectors. Words

80 are likewise messy superpositions of underlying 126 embedding space dimension, typically equal to the  
81 concepts. We can freely encode words into 127 word embedding dimension. The number 10,000,  
82 concepts and decode concepts back into words as 128 chosen as a scalar, represents a typical book chapter  
83 needed. Just as neural networks have deep layers to 129 length in words.  
84 construct representations, language also converts 130 This grouping of sine and cosine values differs  
85 primitive concepts into richer, syntactic concepts. 131 from the original even/odd convention in the  
86 Subunits of the neural network can process 132 Transformer paper (Vaswani et al., 2017). Each  
87 dedicated conceptual subspaces, merging and 133 sine and cosine pair represents the phase of a  
88 projecting the results to manage dimensionality. 134 “hand” on a clock, with periods ranging from 6.28  
89 This is the foundation for Transformer attention - 135 for lower dimensions to  $6.28 \times 10^4$  for larger digits  
90 understanding words in conceptual, syntactic, and 136 ( $i \sim d_{\text{model}}/2$ ). Two key aspects are noteworthy. First,  
91 contextual terms. 137 a sophisticated neural network should capture the  
138 intricacies of the clock-like embedding system, and  
139 not solely depend on Euclidean distance to measure

### 92 2.3 Positional Encoding

93 The Transformer algorithm marks a significant 140 the proximity of positions. Secondly, it's  
94 advancement in how it handles the position of a 141 advantageous to consider an “effective position”  
95 token in a sentence. Unlike Recurrent Neural 142 reflecting a range of nearby PE(pos,  $i$ ) values,  
96 Networks (RNNs), including Long Short-Term 143 similar to the fuzzy representation of concepts by a  
97 Memory machines (LSTMs), which treat positions 144 collection of tokens. This position encoding  
98 as indices, the Transformer considers them as 145 scheme, integrated with neural networks and its  
99 additional information. This approach enables 146 inherent multi-slice fuzziness, is expected to be  
100 parallel processing of all tokens simultaneously. 147 robust and stable, due to its built-in redundancy in  
101 A real number can be represented as either the 148 digits.

102 amplitude or the phase of a complex number. Given 149  
103 that the semantic embedding space resembles a 150  
104 hypersphere and aligns more with phase encoding, 151  
105 applying phase encoding to positions is promising. 152  
106 Analog clocks, using a hand's angle on a circle to 153  
107 denote time, serve as a simple analogy for phase 154  
108 encoding. Although a single hand could suffice, 155  
109 multiple hands are often used in analog clocks to 156  
110 improve resolution. In angular measurements like 157  
111 those in clocks, we often employ modular 158  
112 arithmetic. For instance, we measure seconds 159  
113 within the range of 0 to 60 seconds. However, such 160  
114 discontinuity poses challenges in neural networks, 161  
115 which typically do not perform modular arithmetic 162  
116 inherently. A more effective method involves using 163  
117 the sine and cosine of an angle, providing a 164  
118 continuous representation using two real numbers 165  
119 rather than a single angular value. To encode a 166  
120 word's position, a clock system can be utilized. 167  
121 Consider the following empirical positional 168  
122 encoding for  $0 \leq i < d_{\text{model}}/2$ :

$$123 \begin{cases} \text{PE}(\text{pos}, i) = \sin\left(\frac{\text{pos}}{10,000 \frac{d_{\text{model}}}{2i}}\right), \\ \text{PE}\left(\text{pos}, \frac{d_{\text{model}}}{2} + i\right) = \cos\left(\frac{\text{pos}}{10,000 \frac{d_{\text{model}}}{2i}}\right), \end{cases} \quad (1)$$

124 Here, “pos” represents the word's position in the  
125 text, “ $T$ ” is the dimension index, and “ $d_{\text{model}}$ ” is the

126 embedding space dimension, typically equal to the  
127 word embedding dimension. The number 10,000,  
128 chosen as a scalar, represents a typical book chapter  
129 length in words.

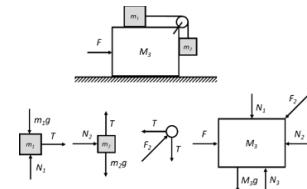
130 This grouping of sine and cosine values differs  
131 from the original even/odd convention in the  
132 Transformer paper (Vaswani et al., 2017). Each  
133 sine and cosine pair represents the phase of a  
134 “hand” on a clock, with periods ranging from 6.28  
135 for lower dimensions to  $6.28 \times 10^4$  for larger digits  
136 ( $i \sim d_{\text{model}}/2$ ). Two key aspects are noteworthy. First,  
137 a sophisticated neural network should capture the  
138 intricacies of the clock-like embedding system, and  
139 not solely depend on Euclidean distance to measure  
140 the proximity of positions. Secondly, it's  
141 advantageous to consider an “effective position”  
142 reflecting a range of nearby PE(pos,  $i$ ) values,  
143 similar to the fuzzy representation of concepts by a  
144 collection of tokens. This position encoding  
145 scheme, integrated with neural networks and its  
146 inherent multi-slice fuzziness, is expected to be  
147 robust and stable, due to its built-in redundancy in  
148 digits.

### 149 2.4 Purpose of Using Positional Encoding

150 Transformers represent a significant innovation by  
151 treating tokens quasi-independently, which enables  
152 parallel processing. Each token independently  
153 passes through the same Transformer neural  
154 network. This process is similar to using “free body  
155 diagrams” in physics, where each object in a  
156 complex system is analyzed separately with  
157 external influences considered as forces (see  
158 Fig.1). In Transformers, tokens navigate through  
159 the neural network, influenced by other tokens only  
160 via attention and feedforward layers, paralleling  
161 the application of Newton's second law

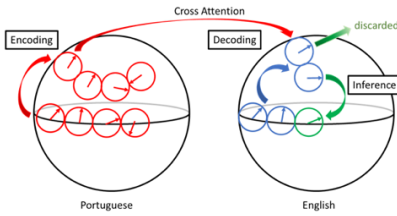
$$162 F_i = \sum_{j=1}^n F_{ij} = m_i a_i \quad (2)$$

163 To facilitate parallel processing, positional  
164 information is embedded in vectors rather than in  
165 the indices of the tokens. For clarity, consider only  
166 the highest two digits of positional encoding, which  
167 form a circular path. Fig. 2 illustrates the  
168 Transformer algorithm operations, pointing out the  
169 relative dominance of positional over semantic  
170 encoding.



172

Fig. 1. Free body diagrams.



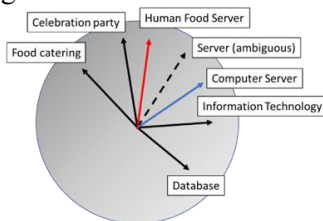
173

174 Fig. 2. Operations diagram of the Transformer  
175 algorithm.

176 The “free body diagram” analogy highlights that all  
177 operations – encoding, decoding, and inference –  
178 occur for each individual token. This design allows  
179 parallel processing of tokens and makes the neural  
180 network flexible to sentence length. The  
181 algorithm’s structure remains unaffected by the  
182 number of tokens, akin to the irrelevance of the  
183 number of bodies in mechanical laws of physics.  
184 During training, Transformers optimize for  
185 efficiency by predicting the next token for each  
186 output token in a training batch simultaneously.  
187 However, in inference mode, the prediction is  
188 instead sequential: one token at a time. Once a  
189 token is predicted, it need not be predicted again.  
190 This disparity between training and inference also  
191 occurs in the behavior of dropout layers (Srivastava  
192 et al., 2014).

193 **2.5 Attention Mechanism as “Two-Body  
194 Forces”**

195 The key innovation of Transformers lies in treating  
196 tokens quasi-independently, allowing for parallel  
197 processing. Analogous to “free body diagrams” in  
198 physics, where each object is analyzed  
199 independently, with external influences from other  
200 objects entering only as forces (as shown in Fig. 3),  
201 Transformers analyze each token in a similar  
202 fashion. This approach facilitates a context-  
203 sensitive representation, enabling a word like  
204 “server” to be interpreted differently based on its  
205 surrounding text.



206

207 Fig. 3. Semantic proximity and context of word.

208 The formal formula for the naïve self-attention is:

209 
$$\text{attention}(Q, K, V) = \text{softmax}\left(\frac{QK^t}{\sqrt{d_k}}\right)V \quad (3)$$

210 Here,  $Q, K, V$  are the “query vector”, “key vector  
211 collection” and “value vector collection”.

212  $Q = w_i$  query word’s embedding (horizontal vector)

213  $K = [w_{ij}]$  matrix from collection of key vectors

214  $V = [w_{ij}]$  same as  $K$  in the case of simple attention

215  $d_k$ : subspace dimension (for normalization)

216 In simple attention,  $V$  equals  $K$ . The output is a  
217 “shift” vector,  $w$  influence of all words ( $K$ ) in the  
218 paragraph, akin to the gravitational force drawing  
219 objects toward a common center. Real-world  
220 Transformers use a more complex form of self-  
221 attention to capture syntactic nuances. Syntax plays  
222 a vital role in understanding words, as  
223 demonstrated by different interpretations of the  
224 verb “left” as an adjective or verb. To incorporate  
225 syntactic information, we can project query and  
226 key vectors onto syntactic subspaces.

227 Extending the physical analogy, attention can be  
228 seen as a “force” shifting embedding vectors.  
229 While naive self-attention is akin to gravitational or  
230 electrostatic forces, more intricate forces are  
231 needed to capture the diverse properties of tokens.  
232 In physics, different forces act on particles based  
233 on their properties, such as mass or charge.  
234 Similarly, each attention mechanism can be viewed  
235 as a distinct force coming from a different type of  
236 subspace charge.

237 The formula for the query-key-value approach  
238 for generalized attention is:

239 
$$\text{attention}(QW_h^Q, KW_h^K, VW_h^V) =$$
  
240 
$$\text{softmax}\left(\frac{QW_h^Q W_h^K K^t}{\sqrt{d_k}}\right) VW_h^V \quad (4)$$

241 The projection-rotation matrices  $W_h^Q, W_h^K, W_h^V$ ,  
242 along with the dimension  $d_k$  of the subspace where  
243 we perform the dot product, are used to calculate,  
244 and determine the attention strengths. Using  
245 multiple “heads” in multi-headed attention allows  
246 for specific syntactic roles to be identified,  
247 applying appropriate shifts to each query word.

248 In self-attention,  $Q, K$ , and  $V$  belong to the same  
249 language, whereas in cross-attention, they belong  
250 to different languages. After processing through  
251 multiple attention heads and a feedforward layer,  
252 the Transformer effectively translates a sequence of  
253 words into concepts. This process is comparable to  
254 the translation of DNA into proteins, where  
255 information is distributed across multiple slices,  
256 not isolated to individual tokens.

257 Overall, Transformers use a sophisticated  
258 mechanism to convert word sequences into

conceptual sequences, adapting to different grammatical structures of languages, much like the intricate process of translating genetic information into functional proteins.

## 2.6 Shift Invariance and Generalized Inner Product

Shift invariance is a crucial concept when considering the relationships between words in a text. The relationship between a word  $A$  and a word  $B$  should remain consistent, even if additional text is prepended to the original sentence where they appear. This principle, known as shift invariance, is essential for nearby words unless strong contextual dependencies are present. Without shift invariance, the Transformer algorithm would need different implementations of attention heads for different positions.

In the standard Transformer model, the positional encoding vector is added to the semantic embedding vector, which seems to compromise translational invariance. The embedding of shifted words yields altered vectors, and the generalized inner product of these vectors does not maintain translational invariance:

$$\begin{aligned} \langle Q, K \rangle &\neq \langle Q + P_{Q1}, K + P_{K1} \rangle \\ &\neq \langle Q + P_{Q2}, K + P_{K2} \rangle \end{aligned} \quad (5)$$

Where  $P_{Q1}$  and  $P_{K1}$  are the positional encoding vectors of the words  $Q$  and  $K$ , and  $P_{Q2}$  and  $P_{K2}$  are their shifted versions. Despite its apparent lack of explicit shift invariance, the Transformer accomplishes invariance through two mechanisms. Firstly, a submanifold exists in the  $Q, K$  space, preserving the value of  $\langle Q, K \rangle$  for certain alignments of positional vectors. Secondly, the projection matrices  $W_i^Q, W_i^K, W_i^V$  can project onto subspaces effectively orthogonal to the positional encoding vectors. Additionally, the input word embedding is trained alongside other neural weights, aiding the Transformer in separating positional and semantic dimensions. If the “conceptual dimensions” represented by vectors  $Q$  and  $K$  are orthogonal to the positional vectors the positional vectors  $P_{Q1}, P_{K1}, P_{Q2}$  and  $P_{K2}$ , we have:

$$\begin{aligned} \langle Q + P_{Q1}, K + P_{K1} \rangle &= \langle Q, K \rangle + \langle P_{Q1}, P_{K1} \rangle \\ \langle Q + P_{Q2}, K + P_{K2} \rangle &= \langle Q, K \rangle + \langle P_{Q2}, P_{K2} \rangle \end{aligned} \quad (6)$$

Given shift invariance in positional encoding (i.e.,  $\langle P_{Q1}, P_{K1} \rangle = \langle P_{Q2}, P_{K2} \rangle$ ), the overall Transformer algorithm can achieve shift invariance. Positional

encoding schemes like a multi-handed clock on a hypertorus satisfy this requirement, as do helix-shaped encodings. The most general topology for positional encoding is a hypertorus, with the helix and straight line as limiting cases, making a multi-handed clock with sinusoidal functions a natural choice.

The Transformer model's capacity to implement shift invariance is remarkable, given the general non-invariance of inner products. This seems to be largely due to the training of the semantic embedding layer in conjunction with other layers, allowing the model to separate semantic and positional dimensions. Conceptual dimensions eventually become orthogonal to positional ones. During encoding, the Transformer allocates certain dimensions for positional encoding while reserving others for semantic or conceptual encoding. This automatic separation of dimensions is one of the Transformer model's most notable features.

## 3 Data and Codes

To build the model, we used Google's Transformer codes ([tensorflow](#)) as a base. We ran the model in a Docker container on an NVIDIA-SMI 470.103.01 powered DGX server, equipped with Tesla V100-SXM2 GPUs and CUDA Version 11.4.

We employed the Portuguese-English translation dataset from TensorFlow Datasets, consisting of approximately 52,000 training, 1,200 validation, and 1,800 test instances. Each instance contains a Portuguese-English sentence pair. For efficiency, we tokenized the dataset and arranged it into ragged batches. Our model processes tokenized Portuguese and English sequence pairs (pt, en) as inputs. The target labels are the corresponding English sequences, offset by one token, ensuring that the target label for each position in the English sequence is the next token. We used `tf.keras.layers.Embedding` layer ([Gal et al., 2016](#)) to convert input Portuguese and target English tokens into vectors.

Unlike static embeddings in Word2Vec and GloVe ([Mikolov et al., 2013](#); [Gal et al., 2016](#); [Pennington et al., 2014](#)), our Transformer model generates dynamic word embeddings. It starts with randomly initialized weights and refines them during training. The dimensionality in word embeddings, referring to the number of features, was set to 128 for this project. The encoder and decoder architectures are similar, each with  $N = 4$  attention layers. Each layer comprises an 8-head

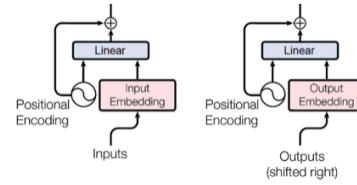
358 multi-head self-attention mechanism and a fully  
 359 connected feed-forward network, enabling concept  
 360 synthesis. The decoder includes an additional  
 361 multi-head cross-attention sub-layer to focus on  
 362 various parts of the encoder's output, utilizing the  
 363 conceptual information for enhanced performance  
 364 in the translation task.

## 365 4 Results

### 366 4.1 Questioning the Adding of Positional 367 Encoding Vector

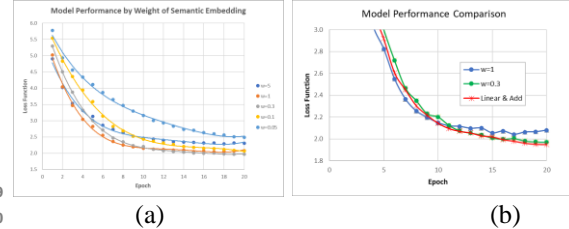
368 In the conventional Transformer algorithm, the  
 369 positional encoding vector is directly added with  
 370 the semantic embedding vector, a process that lacks  
 371 justified rationale. This direct addition is  
 372 problematic, as these vectors originate from  
 373 distinct domains - akin to adding apples with  
 374 oranges. To elucidate this issue, we adapted  
 375 Google's Transformer code and replaced the direct  
 376 sum of the semantic embedding  $S$  vector and  
 377 positional encoding vector  $P$ ,  $S + P$ , with  $wS + P$ ,  
 378 where  $w$  is a weight factor. We evaluated the loss  
 379 function across varying values of  $w$  with results  
 380 depicted in Figure 5(a). Our findings revealed an  
 381 optimal weight of  $w = 0.3$ , challenging the  
 382 convention of a direct sum. This indicates that the  
 383 positional information significantly outweighs the  
 384 semantic information, portraying a picture of small  
 385 "semantic hyperspheres" situated along a broad  
 386 positional encoding trajectory.

387 An enhanced method involves utilizing a full  
 388 linear neural network instead of a simple weighted  
 389 sum, as demonstrated in Figure 4. By integrating a  
 390 dense linear layer (sans activation function) and a  
 391 dropout layer, a more sophisticated combination is  
 392 achieved. Adding back the positional encoding  
 393 vector post-linear layer further enhances model  
 394 stability and mitigates the risk of local minima  
 395 during training. We have termed this methodology  
 396 "Linear & Add," reflecting the sequence of linear  
 397 layer application followed by addition. Figure 5  
 398 contrasts the model performance using three  
 399 distinct methods of combining positional encoding  
 400 with semantic embedding: (1) direct addition, (2)  
 401 weighted sum with a reduced semantic embedding  
 402 significance ( $w=0.3$ ), and (3) the application of a  
 403 dense layer followed by summation. Consistent  
 404 with our hypothesis, the "Linear & Add" method  
 405 outperforms the others.



406

407 Figure 4 "Linear & Add" method to concatenate the  
 408 positional encoding and the semantic embedding vectors



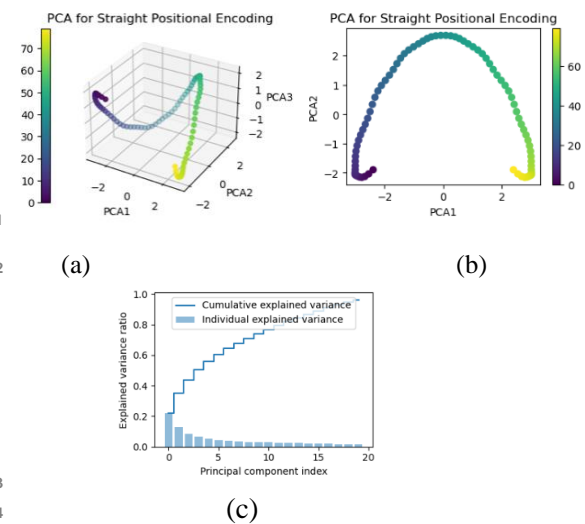
409

410 (a) (b)  
 411 Fig. 5. (a) Model performance by weight of semantic  
 412 embedding. (b) Model performance comparison of three  
 413 different approaches for combining positional encoding  
 414 with semantic embedding.

### 415 4.2 The First Helix: Positional Information 416 After the Encoding Stage of a 417 Transformer

418 To understand the positional encoding in the  
 419 Transformer's input layer, we conducted Principal  
 420 Component Analysis (PCA). Applying PCA to the  
 421 linear-and-add positional encoding vectors (up to  
 422 the 80<sup>th</sup> position) yielded a 2D plot of the first two  
 423 components (Fig. 6b), a 3D plot of the first three  
 424 components (Fig. 6a), and a bar chart showing the  
 425 variance explained (Fig. 6c). The analysis revealed  
 426 that the positional encoding follows a 2D path  
 427 resembling the "Arch of St. Louis" (the Gateway  
 428 Arch). Notice the variance decays slowly across the  
 429 principal components.

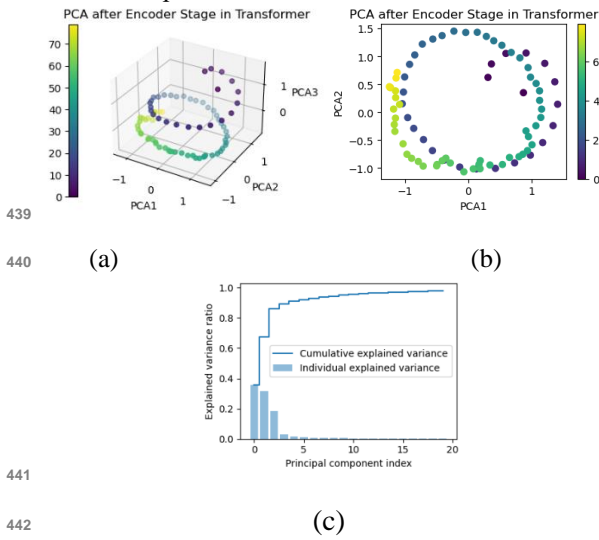
430



431

432 (a) (b)  
 433 (c)  
 434 Fig. 6. Principal Component Analysis of positional  
 435 encoding. (a) 3D plot of the first three PCA components.  
 436

437 (b) 2D plot of the first two PCA components. (c) 438 Variance explained.



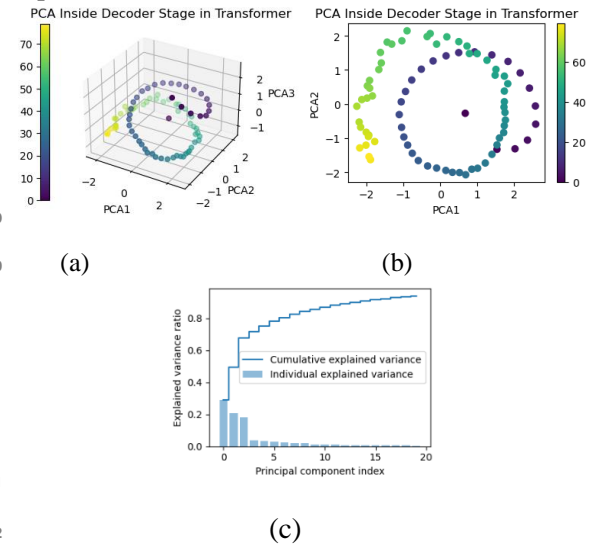
443 Fig. 7. Principal Component Analysis after the encoding 444 stage. (a) 3D plot of the first three PCA components. (b) 445 2D plot of the first two PCA components. (c) Variance 446 explained.

447 We repeated this analysis post-encoding stage of 448 the Transformer. By averaging the embedding 449 vectors of 1,000 randomly selected sentences at 450 identical positions, we isolated positional 451 information, effectively removing semantic, 452 contextual, and syntactic information. PCA on the 453 resulting “average sentence” showed that the 454 residual positional encoding forms a helical shape 455 (Fig. 7a and 7b). In contrast to the input positional 456 encoding that spans all 128 dimensions, the post- 457 encoding positional vectors predominantly use 458 three dimensions (Fig. 7c). This dimensional 459 reduction facilitates shift invariance, allowing the 460 Transformer to capture semantic, syntactic, and 461 contextual information in the remaining 125 462 dimensions.

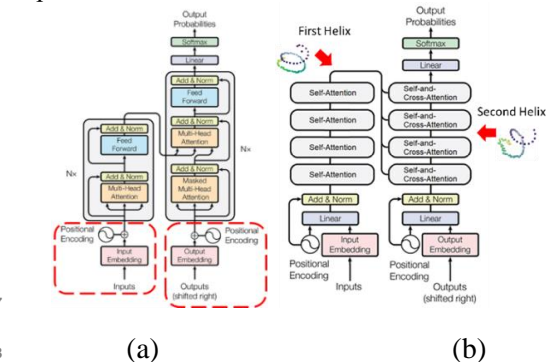
463 One might question the wisdom of combining 464 positional and conceptual information from the 465 outset. While separating them in the input stage is 466 valid, merging them there offers flexibility in 467 handling varying sentence lengths. This approach 468 eliminates the need for preset dimensions for 469 positional information, allowing the algorithm to 470 dynamically allocate dimensions for positional and 471 conceptual elements. This adaptability is key for 472 handling diverse text lengths and complexities, 473 effectively allowing the algorithm to “set its own 474 clock”. Nevertheless, exploring alternative 475 methods for merging information remains a vital 476 area of research (Wang et al., 2020).

### 477 4.3 The Second Helix: Positional 478 Information Deep Inside the Decoder 479 Stage of a Transformer

480 In exploring the decoding stage of the transformer 481 algorithm, we aimed to identify a helical pattern 482 akin to what was observed in the encoding stage. 483 Surprisingly, we found such a helix, but not where 484 we initially expected. In the decoding stage, as 485 tokens are generated to form English sentences, the 486 positional information resembles the original 487 positional encoding, tracing a path similar to the 488 “Arch of St. Louis.” A more detailed examination 489 of the four layers of decoder-attention in Google's 490 transformer algorithm revealed a distinct 3D helix 491 in the second layer, as depicted in Figures 8a and 492 8b. Principal Component Analysis confirmed the 493 three-dimensional nature of this helix. This 494 discovery suggests that the helical pattern is a core 495 characteristic of the transformer architecture, 496 raising interesting questions about its role in the 497 algorithm's language processing and generation 498 capabilities.



503 Fig. 8. Principal Component Analysis in the decoder 504 stage. (a) 2D plot of the first two PCA components. (b) 505 3D plot of the first three PCA components. (c) Variance 506 explained.



509 Fig. 9. (a) Google Transformer model architecture. (b)  
 510 Identification of two helix patterns in the transformer  
 511 model architecture.

#### 512 4.4 Accidental Mingling

513 In the transformer algorithm, positional and  
 514 semantic embeddings are combined, which could  
 515 potentially result in indistinguishable embeddings  
 516 for different words. To illustrate, if token A and  
 517 token B have positional embeddings of 3 and 4, and  
 518 semantic embeddings of 6 and 5, respectively, their  
 519 combined embeddings would both sum to 9. This  
 520 overlap might obscure distinct meanings.  
 521 However, the transformer model employs  
 522 strategies to prevent such overlaps.  
 523 Renormalization in deep learning (De Mello Koch  
 524 et al., 2020) emphasizes slight differences between  
 525 embeddings. The high dimensionality of the  
 526 embedding space and the improbability of overlaps  
 527 in linguistic translation reduce the likelihood of  
 528 collisions. Crucially, the transformer's semantic  
 529 embeddings are trained to be orthogonal to  
 530 positional embeddings, which minimizes the risk  
 531 of mixing up different words, ensuring accurate  
 532 language processing.

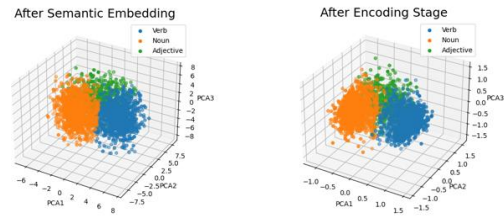
#### 533 4.5 Running Text Density of Words

534 Different languages often need varying word  
 535 counts to express identical ideas. For instance,  
 536 Spanish usually uses more words than English for  
 537 the same concept. This discrepancy poses a  
 538 question about transformers handling the “text  
 539 density” differences. The adaptation likely takes  
 540 place at the transformer's cross-attention layer,  
 541 which discerns text relationships and modulates the  
 542 output's text density to match the target language's  
 543 norms.

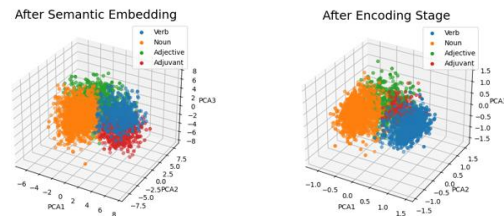
#### 544 4.6 Part of Speech

545 Through analyzing a large sample of sentences and  
 546 calculating average embedding vector values, we  
 547 successfully identified positional information  
 548 vectors for each transformer stage. Isolating the  
 549 conceptual dimensions - syntactic, semantic, and  
 550 contextual information – can be achieved by  
 551 deducting the positional vector from the  
 552 embedding vector. Dimension reduction can then  
 553 be performed using PCA, reducing the dimensions  
 554 from 128 to 5. This allowed us to delve deeper into  
 555 the language structure, uncovering meaningful  
 556 relationships within the language model.

557 Our study focused on the transformer's encoder  
 558 side (Portuguese), examining tokens of four or  
 559 more characters, excluding suffixes, punctuation,  
 560 and sentence markers. For clustering, we employed  
 561 K-Means++ (sklearn.cluster.KMeans) with the  
 562 elbow method to determine the optimal number of  
 563 clusters. This has revealed distinct clusters  
 564 corresponding to parts of speech (PoS) attributes  
 565 like nouns, adjectives, and verbs, as highlighted in  
 566 Fig. 10.



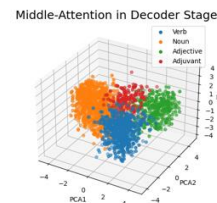
567  
 568 Fig. 10. Tokens formed distinct clusters by their part of  
 569 speech (PoS) attributes, such as nouns, adjectives, and  
 570 verbs. (Left) Pre-attention stage after semantic  
 571 embedding. (Right) Post-attention stage after encoding.



572  
 573 Fig. 11. A second level of clustering split Verboid cluster  
 574 into Verb and Adjuvant clusters. (Left) Pre-attention  
 575 stage after semantic embedding. (Right) Post-attention  
 576 stage after encoding.



577  
 578 Fig. 12. Word Clouds After Encoding Stage. (Top row)  
 579 3 clusters clustering. (Bottom row) Verbs and Adjuvants  
 580 clusters from sub-clustering of “Verboids”



581  
 582 Fig. 13: In mid-attention encoder stage English tokens  
 583 formed distinct clusters by their PoS attributes, verbs,  
 584 nouns, adjectives, and adjuvants.

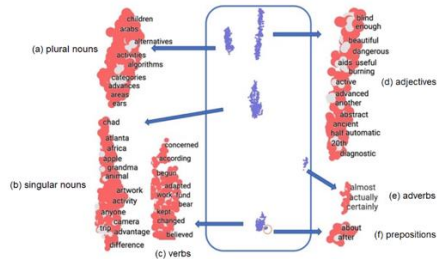


Fig.14. PoS clusters visualized with t-SNE

Interestingly, the verb cluster included numerous “functional words” such as pronouns, prepositions, conjunctions, determiners, and some adverbs. We termed these “adjuvants” - a term derived from Latin, meaning “to help” - due to their supportive role in sentence construction. Applying a second layer of K-Means++ clustering to this “verboid” cluster distinguished the adjuvants as a separate group, as shown in Fig. 11. Fig. 12 displays word clouds post-encoding, showing meaningful grouping. Our findings also indicated that the self-attention layers enhance cluster delineation. This underscores the self-attention mechanism's effectiveness in capturing token relationships and their PoS functions. We observed that points fluctuated between (-6.0, 6.0) in the original semantic layer, narrowing to (-1.0, 1.0) post-encoding, reflecting a “renormalization” process akin to the “renormalization group” in physics (De Mello Koch et al., 2020). Analysis of the decoder side (English) revealed challenges in forming well-defined clusters due to the decoder's dual role in concept formation/translation as well as decoding to tokens. Optimal PoS clustering was achieved in the mid-attention stage (Fig. 13) using PCA and two-step K-means++ clustering. Additionally, we explored t-distributed stochastic neighbor embedding (t-SNE) (Maaten et al., 2008) for visualizing PoS clusters without PCA reduction. Applied directly to 128-dimensional embeddings, distinct PoS clusters emerged after over 10,000 iterations, as illustrated in Fig. 14. Lastly, for post-attention stage analysis, K-Means++ alone was insufficient for clear PoS clustering. However, applying the technique to di-grams of current and next tokens yielded discernible PoS clusters.

#### 4.7 Additional Syntactical, Semantic, and Contextual Information

We observed that removing positional information from vectors reveals “conceptual” information, with single-delta vectors indicating some PoS

cluster structure. For tokens appearing multiple times, we calculated their “semantic vector” by averaging their single-delta vectors. Subtracting the “semantic vector” from the single-delta embedding yields a residual “double-delta vector”, likely representing the syntactical and contextual dimensions of each token. However, comprehending the full extent of syntactical information requires examining token sequences, not just individual tokens. It's possible that single-delta or double-delta vectors contain additional syntactical information, but confirming this necessitates analysis of multiple consecutive tokens. We leave this detailed examination for future research.

## 5 Conclusions

The Transformer architecture has significantly advanced natural language processing with its innovative design and exceptional performance. Central to its function is the encoding of input words into “concepts”, facilitating efficient sequence-to-sequence translation. A key feature of the Transformer is positional encoding, which not only allows for processing of variable-length sequences but also supports parallel processing on modern hardware, akin to the use of free-body diagrams in physics. Remarkably, the Transformer's principles can be applied to diverse data types, such as images (Dosovitskiy et al., 2020) and protein folding (Jumper et al., 2021), demonstrating its versatility and adaptability in various deep learning applications.

This paper has provided an in-depth analysis of the Transformer, examining its handling of positional, syntactical, semantic, and contextual information. We highlighted the unique positional information mapping in the encoding and decoding stages, drawing parallels to the double-helix structure of DNA. The significance of the semantic embedding layer in token clustering by parts of speech was also emphasized. These insights demystify the transformer algorithm, presenting it as a comprehensible tool for AI practitioners. Our goal is to equip professionals with a thorough understanding of the transformer's mechanisms, enabling them to leverage its potential and spur innovation in natural language processing.



## 678 References

- 679 Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD,  
680 Dhariwal P, et al. 2020. *Language models are few-*  
681 *shot learners*. Advances in neural information  
682 processing systems. 33:1877-901.  
683 arXiv:2005.14165
- 684 Ellen De Mello Koch, Robert De Mello Koch and Ling  
685 Cheng. 2020. *Is Deep Learning a Renormalization*  
686 *Group Flow?*. IEEE Access, vol. 8, pp. 106487-  
687 106505, doi: 10.1109/ACCESS.2020.3000901..
- 688 A. Dosovitskiy, L. Beyer, A. Kolesnikov, D.  
689 Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,  
690 M. Minderer, G. Heigold, S. Gelly, et al. 2020. *An*  
691 *image is worth 16x16 words: Transformers for*  
692 *image recognition at scale*. arXiv:2010.11929.
- 693 Yarin Gal and Zoubin Ghahramani. 2016. *A*  
694 *Theoretically Grounded Application of Dropout in*  
695 *Recurrent Neural Networks*. Advances in Neural  
696 Information Processing Systems. arXiv:1512.05287
- 697 Jumper, J., Evans, R., Pritzel, A. et al. 2021. *Highly*  
698 *accurate protein structure prediction with*  
699 *AlphaFold*. Nature 596, 583–589.
- 700 Laurens van der Maaten, L.J.P.; Hinton, G.E. 2008.  
701 *Visualizing High-Dimensional Data Using t-SNE*.  
702 Journal of Machine Learning Research 9:2579-  
703 2605.
- 704 Mikolov, T., Chen, K., Corrado, G., and Dean, J. 2013.  
705 *Efficient estimation of word representations in*  
706 *vector space*. arXiv:1301.3781.
- 707 OpenAI. 2023. *GPT-4 Technical Report*.  
708 arXiv:2303.08774
- 709 Ouyang L, Wu J, Jiang X, Almeida D, Wainwright CL,  
710 Mishkin P, et al. 2022. *Training language models to*  
711 *follow instructions with human feedback*.  
712 arXiv:220302155.
- 713 Pennington, J., Socher, R., & Manning, C. D. 2014.  
714 *GloVe: Global Vectors for Word Representation*. In  
715 Proceedings of the 2014 Conference on Empirical  
716 Methods in Natural Language Processing (EMNLP)  
717 pp. 1532-1543
- 718 Radford A, Narasimhan K, Salimans T, Sutskever I, et  
719 al. 2018. *Improving language understanding by*  
720 *generative pre-training*, OpenAI.
- 721 Radford A, Wu J, Child R, Luan D, Amodei D,  
722 Sutskever I, et al. 2019. *Language models are*  
723 *unsupervised multitask learners*. OpenAI blog.  
724 1(8):9.
- 725 Radford A, Wu J, Amodei D, Amodei D, Clark J,  
726 Brundage M, et al. 2019. *Better language models*  
727 *and their implications*. OpenAI Blog [https://openai](https://openai.com/blog/better-language-models.1(2))  
728 [com/blog/better-language-models.1\(2\)](https://openai.com/blog/better-language-models.1(2)).
- 729 Srivastava N, Hinton G, Krizhevsky A, Sutskever I and  
730 Salakhutdinov R. 2014. *Dropout: a simple way to*  
731 *prevent neural networks from overfitting*. The  
732 Journal of Machine Learning Research.  
733 1;15(1):1929-58.
- 734 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob  
735 Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz  
736 Kaiser and Illia Polosukhin. 2017. *Attention Is All*  
737 *You Need*. Advances in Neural Information  
738 Processing Systems.
- 739 Wang, Yu-An, and Yun-Nung Chen. 2020. *What do*  
740 *position embeddings learn? an empirical study of*  
741 *pre-trained language model positional encoding*.  
742 Proceedings of the 2020 Conference on Empirical  
743 Methods in Natural Language Processing (EMNLP)  
744 <https://www.tensorflow.org/text/tutorials/transformer>  
745  
746