

Attention Is All You Need

The abstract introduces the Transformer, a novel network architecture proposed as an alternative to the dominant sequence transduction models that rely on complex recurrent or convolutional neural networks (RNNs or CNNs) with encoder-decoder structures. While the best-performing contemporary models already incorporated attention mechanisms to connect encoders and decoders, the Transformer distinguishes itself by being based solely on attention mechanisms, thereby completely eliminating recurrence and convolutions.

Experiments conducted on two machine translation tasks demonstrate the Transformer's superiority in quality, its enhanced parallelizability, and a significant reduction in training time compared to existing models. Specifically:

- On the WMT 2014 English-to-German translation task, the Transformer achieved a BLEU score of 28.4. This marked an improvement of over 2 BLEU points compared to the previous best results, including those achieved by ensemble models.
- On the WMT 2014 English-to-French translation task, the model established a new state-of-the-art BLEU score of 41.8 for a single model. This was achieved after only 3.5 days of training on eight GPUs, representing a small fraction of the training costs associated with the leading models reported in the literature at the time.

Furthermore, the abstract highlights the Transformer's capability to generalize well to other tasks beyond machine translation. It was successfully applied to English constituency parsing, performing effectively with both large and limited training datasets.

Recurrent neural networks (RNNs), including long short-term memory (LSTM) and gated recurrent units (GRUs), have been foundational in sequence modeling and transduction tasks like language modeling and machine translation. However, these models inherently process sequences step-by-step. Specifically, they generate a sequence of hidden states h_t where each state is a function of the previous state h_{t-1} and the input at the current position t . This sequential nature is a primary limitation because it prevents parallelization within training examples. This lack of parallelization becomes a significant bottleneck for longer sequences, as memory constraints also limit the ability to batch examples effectively. While techniques like factorization and conditional computation have improved computational efficiency and model performance, the fundamental constraint of sequential computation persists in these recurrent

architectures.

Attention mechanisms have emerged as a powerful tool in sequence modeling, enabling models to capture dependencies between elements in input or output sequences regardless of their distance. However, in most existing applications, these attention mechanisms are used in conjunction with a recurrent network, thus not fully alleviating the sequential processing constraint.

To address these limitations, the paper introduces the Transformer, a novel model architecture. The Transformer's core innovation is its complete departure from recurrence. Instead, it relies entirely on an attention mechanism to draw global dependencies between input and output. This design choice allows for significantly more parallelization during training. The paper states that the Transformer can achieve state-of-the-art results in translation quality, even with relatively short training times (e.g., twelve hours on eight P100 GPUs).

Background

The primary goal motivating alternatives to sequential computation, such as Recurrent Neural Networks (RNNs), is the reduction of the computational steps required. Models like the Extended Neural GPU [16], ByteNet [18], and ConvS2S [9] address this by employing convolutional neural networks (CNNs) as their fundamental components. These architectures compute hidden representations for all input and output positions in parallel. However, in these CNN-based models, the number of operations needed to relate signals from two arbitrary positions within the input or output sequences scales with the distance between these positions. Specifically, this operational dependency grows linearly with distance for ConvS2S and logarithmically for ByteNet. This characteristic makes it more challenging for these models to learn dependencies between distant positions [12].

Self-attention, also known as intra-attention, is an attention mechanism that operates by relating different positions within a single sequence to derive a comprehensive representation of that sequence. Prior to the Transformer, self-attention had been successfully applied to various tasks, including reading comprehension, abstractive summarization, textual entailment, and the learning of task-independent sentence representations [4, 27, 28, 22]. Additionally, end-to-end memory networks utilized recurrent attention mechanisms, as opposed to sequence-aligned recurrence, demonstrating good performance on tasks like simple-language question answering and language modeling [34].

The Transformer model diverges significantly from and builds upon these prior concepts. Unlike the aforementioned CNN-based models (Extended Neural GPU, ByteNet, ConvS2S) where the operations to relate positions scale with distance, the Transformer reduces this to a constant number of operations. This approach facilitates learning long-range dependencies. However, this constant-time operation comes at the cost of reduced effective resolution, which arises from averaging attention-weighted positions. The Transformer mitigates this effect through the use of Multi-Head Attention, a mechanism detailed in section 3.2 of the paper.

Crucially, the Transformer is presented as the first transduction model to rely entirely on self-attention for computing representations of its input and output. This is a key distinction from previous models that might have used attention mechanisms in conjunction with sequence-aligned RNNs or convolutions. The Transformer, by contrast, eschews both sequence-aligned RNNs and convolutional layers, basing its architecture solely on self-attention. The paper further aims to motivate the use of self-attention and discuss its advantages over models such as [17, 18] and [9].

The Transformer model architecture, as depicted in Figure 1, relies on stacked self-attention mechanisms and point-wise, fully connected layers for both its encoder and decoder components.

Encoder and Decoder Stacks

The core of the Transformer consists of an encoder and a decoder, each being a stack of $N=6$ identical layers.

Encoder: Each of the $N=6$ layers in the encoder is composed of two sub-layers:

1. A multi-head self-attention mechanism.
2. A simple, position-wise fully connected feed-forward network.

A residual connection is employed around each of these two sub-layers, followed by layer normalization. The output of each sub-layer is thus $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{sublayer}(x)$ is the function implemented by the sub-layer itself. This residual connection can be seen as a way to allow gradients to propagate more easily through deep networks, analogous to identity mappings in transformations which preserve information from the previous step. All sub-layers in the model, including the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$. This consistent dimensionality throughout the model facilitates these residual connections, ensuring operations occur within a vector space of fixed dimension.

Decoder: The decoder also consists of a stack of $N=6$ identical layers. Each decoder layer, in addition to the two sub-layers found in an encoder layer (multi-head self-attention and position-wise feed-forward network), incorporates a third sub-layer. This third sub-layer performs multi-head attention over the output of the encoder stack. Similar to the encoder:

- Residual connections are used around each sub-layer, followed by layer normalization.
- All sub-layers and embedding layers produce outputs of dimension $d_{model} = 512$.

A key modification in the decoder's self-attention sub-layer is the prevention of positions from attending to subsequent positions. This masking, combined with the fact that output embeddings are offset by one position relative to the input, ensures the auto-regressive property: predictions for a position i can only depend on known outputs at positions less than i . This is mathematically akin to defining terms in a sequence recursively, where each term y_i is a function of $y_1, \dots, y_{\{i-1\}}$.

Attention

An attention function maps a query and a set of key-value pairs to an output. The query, keys, values, and output are all vectors. The output is computed as a weighted sum of the values. The weight assigned to each value is determined by a compatibility function of the query with the corresponding key. This process is analogous to calculating a weighted average, where the weights (derived from the compatibility function) signify the importance of each value vector with respect to the query vector.

Scaled Dot-Product Attention

The specific attention mechanism used is "Scaled Dot-Product Attention" (Figure 2, left).

- Inputs: Queries and keys of dimension d_k , and values of dimension d_v .
- Process:
 1. Dot products of the query with all keys are computed. This dot product acts as the raw compatibility score, measuring the similarity between query and key vectors.
 2. Each dot product is then scaled by dividing by $\sqrt{d_k}$. This scaling factor is crucial because for large values of d_k , the dot products can grow large in magnitude, pushing the softmax function into regions with extremely small gradients, hindering learning. The $\sqrt{d_k}$ scaling, derived from the assumption that components of query and key vectors are independent random variables

with mean 0 and variance 1 (making their dot product have variance d_k), helps normalize the variance of these dot products.

3. A softmax function is applied to these scaled dot products to obtain non-negative weights that sum to one, representing the attention distribution over the values.
4. The output is the weighted sum of the values using these weights.

When processing multiple queries simultaneously (packed into a matrix Q), keys (matrix K), and values (matrix V), the output is computed as: $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$. This formulation relies entirely on matrix multiplication operations, which are highly optimized and efficient.

Compared to additive attention (which uses a feed-forward network for compatibility), dot-product attention is faster and more space-efficient. While additive attention can outperform unscaled dot-product attention for larger d_k , the scaling factor in Scaled Dot-Product Attention addresses this issue.

Multi-Head Attention

Instead of a single attention function using d_{model} -dimensional inputs, Multi-Head Attention (Figure 2, right) employs multiple attention "heads" in parallel.

- Process:
 1. The queries, keys, and values are linearly projected h times using different learned linear transformations (parameter matrices W_i^Q, W_i^K, W_i^V) into lower dimensions: d_k for queries and keys, and d_v for values. In this model, $d_k = d_v = d_{\text{model}} / h$. These projections can be viewed as mapping the original vectors into h different representation subspaces.
 2. Scaled Dot-Product Attention is applied in parallel to each of these h projected sets of queries, keys, and values, yielding h output vectors (head_i). Each head can thus learn to focus on different aspects or relationships in the data.
 3. The h output vectors (each d_v -dimensional) are concatenated.
 4. The concatenated vector is then projected again using another learned linear transformation (parameter matrix W^O) to produce the final output, which is d_{model} -dimensional.

The formula is: $\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$ where $\text{head}_i = \text{Attention}(Q W_i^Q, K W_i^K, V W_i^V)$ and $W_i^Q \in \mathbb{R}^{(d_{\text{model}} \times d_k)}$, $W_i^K \in \mathbb{R}^{(d_{\text{model}} \times d_k)}$, $W_i^V \in \mathbb{R}^{(d_{\text{model}} \times d_v)}$, and $W^O \in \mathbb{R}^{(h \cdot d_v \times d_{\text{model}})}$.

In this work, $h = 8$ parallel attention layers (heads) are used. For each head, $d_k = d_v = d_{\text{model}} / h = 512 / 8 = 64$. This multi-head approach allows the model to jointly attend to information from different representation subspaces at different positions, which might be inhibited by the averaging effect in single-head attention. The total computational cost remains similar to a single-head attention with full d_{model} dimensionality because of the reduced dimension used within each head.

Applications of Attention in our Model

Multi-head attention is utilized in three distinct ways within the Transformer:

1. Encoder-Decoder Attention:

- Queries originate from the previous decoder layer.
- Memory keys and values are sourced from the output of the encoder stack.
- This configuration enables every position in the decoder to attend to all positions in the input sequence. This mechanism is fundamental for sequence-to-sequence tasks, allowing the model to, for instance, align parts of a source sentence with parts of its translation.

2. Encoder Self-Attention:

- Keys, values, and queries all originate from the output of the previous layer within the encoder itself.
- This allows each position in the encoder to attend to all other positions in the previous layer of the encoder. This helps the encoder build contextual representations of each input token by considering its relationship with all other tokens in the input sequence.

3. Decoder Self-Attention (Masked):

- Keys, values, and queries all originate from the output of the previous layer within the decoder itself.
- This allows each position in the decoder to attend to all positions in the decoder up to and including its own position.
- Crucially, a mask is applied within the scaled dot-product attention to prevent "leftward" information flow (i.e., attending to subsequent positions). This is achieved by setting the softmax input values corresponding to illegal (future) connections to $-\infty$. This masking preserves the auto-regressive property, ensuring that generating the token at position i only depends on previously generated tokens. This is mathematically similar to operations involving lower or upper triangular matrices, where dependencies are strictly ordered.

Position-wise Feed-Forward Networks

In addition to attention sub-layers, each layer in both the encoder and decoder stacks contains a fully connected feed-forward network (FFN).

- This FFN is applied to each position separately and identically. This means the same set of transformations (defined by the network's weights and biases) is applied to the representation of each token, independent of other tokens' representations at this stage.
- The FFN consists of two linear transformations with a Rectified Linear Unit (ReLU) activation function in between: $\text{FFN}(x) = \max(0, x \cdot W_1 + b_1) \cdot W_2 + b_2$
- While the linear transformations are identical across positions within a layer, different layers use different parameters (W_1 , b_1 , W_2 , b_2).
- This can also be described as two convolutions with a kernel size of 1.
- The input and output dimensionality of the FFN is $d_{\text{model}} = 512$. The inner-layer has a dimensionality of $d_{\text{ff}} = 2048$. This expansion and subsequent contraction of dimensionality allows the network to learn more complex transformations for each position.

Embeddings and Softmax

- Embeddings: Learned embeddings are used to convert input tokens (e.g., words) and output tokens into vectors of dimension $d_{\text{model}} = 512$. This mapping from a discrete vocabulary to a continuous vector space allows the model to capture semantic relationships between tokens. In the embedding layers, these learned embedding weights are multiplied by $\sqrt{d_{\text{model}}}$.
- Softmax: To convert the decoder's final output (a vector of dimension d_{model} for each position) into predicted next-token probabilities, a standard learned linear transformation is applied, followed by a softmax function. The linear transformation projects the d_{model} -dimensional vector to a vector whose dimension is the size of the vocabulary, and the softmax function then normalizes these scores into a probability distribution across all possible tokens.
- Weight Sharing: The same weight matrix is shared between the two embedding layers (input token embeddings and output token embeddings) and the pre-softmax linear transformation. This reduces the total number of parameters in the model and can improve generalization.

Positional Encoding

Since the Transformer model architecture does not include recurrence (like RNNs) or convolution (like CNNs), it has no inherent mechanism to understand the order or

position of tokens in a sequence. To address this, "positional encodings" are introduced.

- **Function:** These encodings inject information about the relative or absolute position of tokens into the model. They are added to the input embeddings at the bottom of both the encoder and decoder stacks. The positional encodings have the same dimension d_{model} as the embeddings, allowing them to be summed directly. This is analogous to explicitly providing an index for each element in a mathematical sequence.
- **Method:** Sinusoidal functions of different frequencies are used: $PE(\text{pos}, 2i) = \sin(\text{pos} / 10000^{(2i/d_{\text{model}})})$ $PE(\text{pos}, 2i+1) = \cos(\text{pos} / 10000^{(2i/d_{\text{model}})})$. Here, pos is the token's position in the sequence (e.g., 0, 1, 2, ...), and i is the dimension index within the encoding vector (from 0 to $d_{\text{model}}/2 - 1$). Each dimension of the positional encoding thus corresponds to a sinusoid. The wavelengths of these sinusoids form a geometric progression, ranging from 2π to $10000 * 2\pi$.
- **Significance:** This choice of sinusoidal functions was motivated by the hypothesis that it would allow the model to easily learn to attend by relative positions. Specifically, for any fixed offset k , $PE_{\{\text{pos}+k\}}$ can be represented as a linear function of $PE_{\{\text{pos}\}}$. This is a direct consequence of trigonometric addition formulas (e.g., $\sin(A+B) = \sin A \cos B + \cos A \sin B$), meaning that the positional encoding of $\text{pos}+k$ can be obtained from $PE_{\{\text{pos}\}}$ via a linear transformation (a rotation in the 2D plane formed by the \sin and \cos pair for each frequency). This property is mathematically elegant and useful for modeling relative positions.
- While experiments with learned positional embeddings yielded nearly identical results, the sinusoidal version was chosen because it might enable the model to generalize better to sequence lengths not encountered during training. The use of sinusoids of varying frequencies is somewhat analogous to Fourier analysis, where signals are decomposed into constituent frequencies.

This section compares self-attention layers with recurrent and convolutional layers, which are commonly used for transforming a variable-length sequence of symbol representations (x_1, \dots, x_n) to another sequence of equal length (z_1, \dots, z_n) , where $x_i, z_i \in \mathbb{R}^d$. The preference for self-attention is motivated by three key aspects: computational complexity, parallelization capability, and the efficiency in learning long-range dependencies.

These aspects are summarized in Table 1:

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types.

- n is the sequence length.

- d is the representation dimension.
- k is the kernel size of convolutions.
- r is the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

The desiderata for choosing self-attention are:

1. Total Computational Complexity per Layer:

- Self-Attention: $O(n^2 \cdot d)$.
- Recurrent: $O(n \cdot d^2)$.
- Convolutional: $O(k \cdot n \cdot d^2)$ for standard convolutions. Separable convolutions can reduce this to $O(k \cdot n \cdot d + n \cdot d^2)$.
- Self-attention layers are faster than recurrent layers when the sequence length n is smaller than the representation dimensionality d . This is often the case for sentence representations used in machine translation (e.g., word-piece or byte-pair encodings). For instance, if n represents the number of terms in a sequence and d the richness of each term's description, an $O(n^2 \cdot d)$ complexity signifies a quadratic cost in sequence length, while $O(n \cdot d^2)$ is quadratic in descriptive richness. If $n < d$, the former becomes more efficient. This is analogous to choosing an algorithm whose complexity depends favorably on the smaller dimension of the input data.
- Even with $k = n$, the complexity of a separable convolution is comparable to the combination of a self-attention layer and a point-wise feed-forward layer, an approach adopted in the Transformer model.
- For very long sequences, restricted self-attention, which only considers a neighborhood of size r , can improve computational performance to $O(r \cdot n \cdot d)$.

2. Amount of Computation That Can Be Parallelized: This is measured by the

minimum number of sequential operations required.

- Self-Attention: $O(1)$. This means all positions can be processed in parallel, as if computing elements of a vector where each element's computation is independent of others.
- Recurrent: $O(n)$. These layers require sequential operations, as the computation for a position i depends on the computation for position $i-1$. This is akin to calculating terms in a recurrence relation $a_i = f(a_{i-1})$, where each step must await the completion of the previous.
- Convolutional: $O(1)$ (assuming the kernel operations themselves can be parallelized across positions).
- The $O(1)$ sequential operations for self-attention allow for significantly more parallelization compared to recurrent layers. This is like having an algorithm where all fundamental computational units can be executed simultaneously, given sufficient processing resources, as opposed to an algorithm that has an inherent sequential bottleneck.

3. Path Length Between Long-Range Dependencies: Learning long-range dependencies is crucial in sequence transduction. Shorter paths for forward and backward signals between any two positions in the input and output sequences facilitate this learning.

- Self-Attention: $O(1)$. It connects all input positions to all output positions with a constant path length. This is analogous to a complete graph where any two nodes are directly connected, allowing information to travel in a single step.
- Recurrent: $O(n)$. The path length between distant positions scales linearly with the distance, meaning information has to traverse many intermediate steps, potentially degrading or being lost. This resembles information flow in a linear chain or a path graph.
- Convolutional: A single convolutional layer with kernel width $k < n$ does not connect all pairs of input/output positions. To connect all positions, a stack of $O(n/k)$ convolutional layers is needed for contiguous kernels, or $O(\log_k(n))$ layers for dilated convolutions. This creates hierarchical paths; for dilated convolutions, the path length grows logarithmically, which is more efficient than linear but less than constant. Think of this as information propagating through layers of local neighborhoods, gradually expanding its reach.
- Restricted Self-Attention: Increases the maximum path length to $O(n/r)$. This is a trade-off for improved computational efficiency with very long sequences, where direct global connections might be too costly.

As a side benefit, self-attention mechanisms can also yield more interpretable models, as the attention distributions can reveal how the model relates different parts of the input and output sequences. Individual attention heads often learn to perform distinct tasks related to syntactic and semantic structures within sentences.

Training

This section outlines the training methodology for the Transformer models.

Training Data and Batching

The models were trained using two primary datasets:

1. WMT 2014 English-German Dataset:

- Comprised approximately 4.5 million sentence pairs.
- Sentences were encoded using byte-pair encoding (BPE). This technique can be conceptualized as a data-driven method to find efficient sub-word units, akin to identifying "prime components" of words, resulting in a shared source-target vocabulary of about 37,000 tokens.

2. WMT 2014 English-French Dataset:

- A significantly larger dataset containing 36 million sentences.
- Tokens were split into a 32,000 word-piece vocabulary.

For training, sentence pairs were batched together based on approximate sequence length. Each training batch was constructed to contain approximately 25,000 source tokens and 25,000 target tokens. This batching strategy is an efficiency measure, analogous to grouping elements of similar scale in numerical computations.

Hardware and Schedule

The training was conducted on a single machine equipped with 8 NVIDIA P100 GPUs, leveraging their parallel processing capabilities.

- Base Models:

Utilized hyperparameters described throughout the paper.

Each training step took approximately 0.4 seconds.

These models were trained for a total of 100,000 steps, which corresponded to 12 hours of training.

- Big Models: (Details for this configuration are specified as being from "the bottom line of table 3")

Step time was 1.0 second.

These larger models were trained for 300,000 steps, requiring 3.5 days.

The training process, involving a set number of steps, can be likened to an iterative numerical method converging towards a solution.

Optimizer

The Adam optimizer was employed for training, a sophisticated algorithm for gradient-based optimization well-suited for high-dimensional parameter spaces.

- Adam Parameters: $\beta_1 = 0.9$, $\beta_2 = 0.98$, and $\epsilon = 10^{-9}$.
- Learning Rate (**lr**): The learning rate was varied dynamically over the course of training according to the following formula: $lr = \frac{d_{model}}{\min(step_num, warmup_steps) \cdot \sqrt{step_num}}$. Here, d_{model} represents the dimensionality of the model, and $step_num$ is the current training step number.
- Learning Rate Schedule:
 1. The learning rate increases linearly for the first $warmup_steps$. This initial phase, with $warmup_steps$ set to 4000, can be compared to cautiously initiating a search in a complex landscape with small steps.
 2. After the warm-up phase, the learning rate decreases proportionally to the inverse square root of $step_num$. This strategy of gradually reducing the step size is common in optimization algorithms as they approach an optimum, helping to prevent overshooting and refine the solution.

Regularization

Three types of regularization techniques were utilized during training to prevent overfitting and improve generalization, conceptually similar to adding penalty terms or constraints in mathematical optimization problems.

1. Residual Dropout:

- Dropout was applied to the output of each sub-layer before this output was added to the sub-layer's input (the residual connection) and subsequently normalized.
- Additionally, dropout was applied to the sums of the embeddings and the positional encodings in both the encoder and decoder stacks.

- For the base model, a dropout rate of $P_{\text{drop}} = 0.1$ was used. Dropout acts by randomly "ignoring" a fraction of neuron outputs during training, which can be thought of as forcing the network to learn more robust features, not overly reliant on any single pathway, much like testing a system's stability by randomly removing components.

2. Label Smoothing:

- A label smoothing value of $\lambda_s = 0.1$ was employed.
- This technique modifies the target labels from being hard (e.g., one-hot encoded) to being "softer". Instead of aiming for 100% confidence on the correct label, it encourages the model to assign a small portion of the probability mass to other labels.
- The authors note that while label smoothing hurts perplexity (meaning the model becomes less confident in its specific predictions), it leads to improvements in accuracy and BLEU score (a metric for translation quality). This can be seen as preventing the model from becoming over-confident, leading to a smoother decision boundary and better generalization, akin to considering a small neighborhood around a point solution in an optimization problem.

The experimental results demonstrate the Transformer architecture's efficacy and efficiency, establishing new benchmarks in natural language processing tasks and providing insights crucial for subsequent Large Language Model (LLM) development.

6.1 Machine Translation

The Transformer model's performance was rigorously evaluated on WMT 2014 English-to-German (EN-DE) and English-to-French (EN-FR) translation tasks. The results, detailed in Table 2, showcased state-of-the-art performance.

On the English-to-German task, the "Transformer (big)" model achieved a BLEU score of 28.4. This score surpassed all previously reported models, including ensembles, by more than 2.0 BLEU points, establishing a new state-of-the-art. This achievement is significant, comparable to a breakthrough in solving a long-standing mathematical conjecture. The training for this model took 3.5 days on 8 P100 GPUs. Even the "Transformer (base model)" outperformed all previous models and ensembles, yet at a substantially lower training cost.

On the English-to-French task, the "Transformer (big)" model achieved a BLEU score of 41.8 (as per Table 2, while the text mentions 41.0 when comparing against single

models). This score also represented a new state-of-the-art, outperforming previous top single models at less than a quarter of their training cost. For this task, the Transformer (big) model used a dropout rate $p_{\text{drop}} = 0.1$.

The following table, extracted from Table 2, highlights the Transformer's performance against previous state-of-the-art models and its training efficiency:

Model	EN-DE BLEU	EN-FR BLEU	Training Cost (FLOPs) EN-DE	Training Cost (FLOPs) EN-FR
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	-
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	-

Note: Training costs for EN-FR for Transformer models are not explicitly detailed in the FLOPs comparison part of Table 2.

The substantial reduction in training cost (measured in FLOPs) for achieving superior or competitive results can be analogized to discovering a more efficient algorithm in mathematics, reducing computational complexity (e.g., from exponential to polynomial time).

Training utilized techniques like residual dropout ($p_{\text{drop}} = 0.1$ for the base model) applied to sub-layer outputs and summed embeddings/positional encodings, and label smoothing ($\epsilon_{\text{ls}} = 0.1$), which, while slightly increasing perplexity, improved accuracy and BLEU scores. For inference, beam search (beam size 4, length penalty $= 0.6$) was employed. Base models averaged the last 5 checkpoints, while big models averaged the last 20.

6.2 Model Variations

To understand the contribution of different architectural components, variations of the base Transformer model were tested on the English-to-German development set (newstest2013). These experiments, summarized in Table 3, provide insights analogous to exploring a parameter space to find an optimal configuration or understanding how different terms in a complex equation contribute to the overall solution.

Key findings include:

- (A) Number of Attention Heads (h) and Key/Value Dimensions (d_k, d_v):

Varying the number of attention heads (h) while keeping the computational load per layer roughly constant showed that a single attention head ($h=1, d_k=d_v=512$) performed 0.9 BLEU worse (24.9 BLEU) than the baseline 8 heads (25.8 BLEU). This is like choosing an inappropriate dimensionality for a projection; if too low, information is lost.

Performance also degraded with an excessive number of heads ($h=32, d_k=d_v=16$, resulting in 25.4 BLEU). An optimal balance exists. The base model used $h=8$ with $d_k=d_v=64$.

- (B) Attention Key Size (d_k):

Reducing the key dimension d_k (e.g., to 16 or 32, from the base of 64) while keeping d_v appropriately scaled, negatively impacted model quality (e.g., $d_k=16$ gave 25.1 BLEU).

This suggests that the mechanism for determining compatibility (dot product) benefits from richer key representations and that a more sophisticated compatibility function might be advantageous. Reducing d_k too much is like trying to define a precise relationship using vectors that have lost too much of their distinguishing information.

- (C) Model Size ($d_{\text{model}}, d_{\text{ff}}$):

Generally, larger models (increased d_{model} or feed-forward inner-layer dimensionality d_{ff}) performed better, achieving higher BLEU scores, though with more parameters. For instance, increasing d_{model} from 256 (24.5 BLEU) to 1024 (26.0 BLEU) or d_{ff} from 1024 (25.4 BLEU for $d_{\text{model}}=256$) to 4096 (26.2 BLEU for $d_{\text{model}}=512$) improved results. This is expected, much like how higher-capacity functions can model more complex data, but also increases the risk of overfitting if not regularized.

- (D) Dropout and Label Smoothing ($P_{\text{drop}}, \epsilon_{\text{ls}}$):

Dropout was confirmed to be very helpful in preventing overfitting. For example, removing residual dropout ($P_{\text{drop}} = 0.0$) decreased BLEU to 24.6 from the baseline 25.8. This acts as a regularization technique, akin to adding a penalty term in mathematical optimization to prevent solutions from becoming too complex.

The effect of label smoothing (ϵ_{ls}) was also explored; $\epsilon_{\text{ls}} = 0.0$ gave 25.3 BLEU, slightly lower than $\epsilon_{\text{ls}} = 0.1$ (25.8 BLEU) for the base settings.

- (E) Positional Encoding:

Replacing the fixed sinusoidal positional encodings with learned positional embeddings resulted in nearly identical performance (25.7 BLEU vs. 25.8 BLEU for the base model). This suggests that while positional information is crucial, the specific method of sinusoidal encoding is a robust choice, or the model can readily learn an equally effective representation. This is like finding that a problem can be solved effectively using either a predefined basis (like Fourier series) or an adaptively learned one.

The "Transformer (big)" model ($N=6$, $d_{\text{model}}=1024$, $d_{\text{ff}}=4096$, $h=16$, $P_{\text{drop}}=0.3$ - though EN-FR big model used 0.1) achieved 26.4 BLEU on this development set with 300K steps.

6.3 English Constituency Parsing

To assess the Transformer's generalizability beyond translation, it was applied to English constituency parsing. This task is challenging due to strong structural constraints on the output and its significantly greater length compared to the input. Traditionally, RNN-based sequence-to-sequence models had not achieved state-of-the-art results in small-data regimes for this task.

A 4-layer Transformer ($d_{\text{model}}=1024$) was trained on:

1. The Wall Street Journal (WSJ) section of the Penn Treebank (~ 40K sentences) with a 16K token vocabulary.
2. A semi-supervised setting using the WSJ data plus ~ 17M sentences from larger corpora, with a 32K token vocabulary.

Minimal hyperparameter tuning was performed. The results on WSJ Section 23 (F1 score) are shown below, demonstrating the Transformer's strong generalization capabilities:

Parser	Training	WSJ 23 F1
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Vinyals & Kaiser et al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7

Dyer et al. (2016) [8]	generative	93.3
------------------------	------------	------

Despite the lack of task-specific tuning, the Transformer performed surprisingly well:

- With WSJ-only training, it achieved an F1 score of 91.3, outperforming the BerkeleyParser [29] (90.4 F1).
- In the semi-supervised setting, it achieved an F1 score of 92.7. These results were better than all previously reported models except for the Recurrent Neural Network Grammar [8].

The Transformer's success in this distinct task, without extensive adaptation, highlighted the fundamental power of its attention-based architecture. This is analogous to a powerful mathematical tool developed for one domain (e.g., number theory) finding unexpected and effective applications in another (e.g., graph theory or physics).

Significance for LLM Development: These collective results were pivotal for the trajectory of LLM development.

1. New Performance Standards: Establishing new SOTA in machine translation, particularly with enhanced efficiency (lower FLOPs), demonstrated the viability of attention as a primary mechanism, encouraging the field to move beyond recurrent architectures for many sequence modeling tasks. This is akin to a paradigm shift in mathematical approaches to a class of problems.
2. Architectural Insights: The model variation experiments provided crucial data on how components like attention heads, model dimensions, and regularization (dropout) affect performance. This empirical evidence guided the design of more complex and deeper Transformer models. It's like meticulously analyzing the components of a mathematical structure to understand their individual and collective roles.
3. Generalization Capability: Success in English constituency parsing underscored that the Transformer architecture was not merely a specialized translation model but a more general-purpose sequence processing engine. This versatility was a key factor in its adoption and adaptation for a wide array of NLP tasks, forming the bedrock of modern LLMs like BERT and GPT. This is similar to the discovery of a fundamental mathematical concept, like vector spaces or group theory, that proves applicable and unifying across diverse fields.

The combination of high performance, computational efficiency relative to previous SOTA, and strong generalization capabilities firmly positioned the Transformer as a foundational architecture for the subsequent rapid advancements in Large Language Models.

The paper concludes by presenting its central contribution: the Transformer, the pioneering sequence transduction model predicated entirely on attention mechanisms. This model revolutionizes traditional encoder-decoder architectures by replacing recurrent layers with multi-headed self-attention. This architectural shift can be likened to moving from a sequential processing model, where each step depends strictly on the previous (analogous to an iterative function $f_n = g(f_{n-1}, x_n)$), to a paradigm allowing for more parallel computation of representations by considering all elements of a sequence simultaneously. The "multi-headed" aspect of self-attention itself can be understood as employing multiple, distinct "projection subspaces" or analytical perspectives in parallel to capture diverse types of relationships within the data.

The Transformer demonstrated significant practical advantages. For translation tasks, it trains considerably faster than architectures reliant on recurrent or convolutional layers. It also set new state-of-the-art performance levels on both the WMT 2014 English-to-German and WMT 2014 English-to-French translation benchmarks. Notably, on the English-to-German task, the authors' best Transformer model surpassed the performance of all previously reported ensemble models.

The authors express strong optimism regarding the future of attention-based models and outline several directions for future work:

- **Broader Applications:** Extending the application of attention-based models to tasks beyond machine translation.
- **Multi-Modal Extension:** Adapting the Transformer to handle problems involving input and output modalities other than text, such as images, audio, and video. This involves generalizing the sequence-to-sequence framework to richer, potentially multi-dimensional data structures.
- **Efficient Attention:** Investigating local, restricted attention mechanisms. The goal is to efficiently process very large inputs and outputs by focusing computational resources, which is conceptually similar to identifying and leveraging sparsity in large mathematical systems or concentrating on the most significant terms in a series expansion.
- **Less Sequential Generation:** Pursuing methods to make the output generation process less sequential, thereby enabling more parallel computation during inference.

To foster further research and development, the code used for training and evaluating their models has been made publicly available at <https://github.com/tensorflow/tensor2tensor>.

Attention visualizations offer insights into the model's internal mechanisms, demonstrating how it processes language. These visualizations are optional for understanding the core model but provide deeper insights into its learned behaviors.

Figure 3: Long-Distance Dependencies This figure illustrates the attention mechanism's capability to handle long-distance dependencies within the encoder self-attention in layer 5 of 6.

- **Observation:** When focusing on the verb "making," many of the attention heads (represented by different colors) connect it to words that are distant in the sequence.
- **Specific Insight:** The visualization shows that the model successfully completes the phrase "making...more difficult" by attending to "more difficult" even though several words separate it from "making." This highlights the model's ability to capture relationships between words that are not immediately adjacent, a crucial aspect for understanding complex sentences. This could be likened to how, in a mathematical sequence, a term's properties might be determined by terms far preceding it, not just its immediate neighbors.

Figure 4: Anaphora Resolution This figure showcases how two specific attention heads, also in layer 5 of 6, appear to be involved in anaphora resolution. Anaphora resolution is the task of identifying what a pronoun or other referring expression refers to.

- **Context:** The example sentence is "The Law will never be perfect, but its application should be just..."
- **Visualization:**

The top panel displays the full attention patterns for one of these heads (head 5).

The bottom panel isolates the attention from the word "its" for two heads (heads 5 and 6).

- **Key Insight:** The attentions for the word "its" are notably "very sharp." This indicates that the model has learned to strongly associate the pronoun "its" with its antecedent, "Law." This precision is vital for correct interpretation. Mathematically, this sharp focus is analogous to a function that maps an input (the pronoun) to a very specific output (its antecedent) with high confidence or a strong, unambiguous signal, rather than a diffuse mapping.

Figure 5: Task Specialization of Attention Heads This figure demonstrates that different attention heads within the encoder self-attention (at layer 5 of 6) learn to focus on different types of structural information within a sentence.

- **Observation:** The visualization provides two examples from two different attention

heads, showing distinct patterns of attention across the sentence.

- Key Insight: The fact that these heads exhibit different behaviors suggests they have "clearly learned to perform different tasks." This implies a form of specialization, where each head might become sensitive to particular syntactic or semantic relationships. This is akin to how a set of basis vectors in a vector space can each capture a different dimension or aspect of the data, or how different mathematical operators are designed for distinct functions. One head might focus on verb-object relationships, while another might track noun modifiers, collectively contributing to a comprehensive understanding of the sentence structure.

The Transformer model, introduced in the paper "Attention Is All You Need," marked a significant departure from previous approaches to tasks involving sequences of data, such as translating languages.

Core Concepts: Moving Beyond Sequential Processing with "Attention"

Prior to the Transformer, dominant AI models for sequence tasks, like machine translation, often relied on recurrence. This means they processed information—such as the words in a sentence—one by one, in a strict order. Imagine reading a very long mathematical proof by only looking at one line at a time, with your understanding of the current line heavily dependent on your memory of the immediately preceding one. This sequential nature has two main drawbacks:

1. It's inherently slow because processing the next item often requires finishing the current one. This limits how much computation can be done in parallel.
2. For long sequences (like lengthy sentences or entire documents), information from the beginning can get diluted or "forgotten" by the time the model reaches the end, making it hard to capture long-range dependencies.

The Transformer's revolutionary idea was to dispense almost entirely with this recurrence. Instead, its architecture is built upon a mechanism called "attention."

The attention mechanism allows the model to process all elements in a sequence (e.g., all words in a source sentence) simultaneously. When producing an output (e.g., the next word in a translated sentence), it can "look back" at any or all parts of the input sequence and assign different "weights" or "importance scores" to them. This means it can dynamically decide which input words are most relevant for generating the current output word, regardless of their position in the sentence.

- Mathematical Analogy for Attention: Consider a mathematician trying to solve a complex problem or derive a new theorem. They might have numerous axioms, definitions, and previously proven lemmas at their disposal. At any step in their

reasoning, they don't just consider the last statement they wrote; they can draw upon any of these foundational pieces of knowledge, intuitively weighing which ones are most pertinent to the current part of the problem. The attention mechanism in a Transformer formalizes a similar idea for language: it learns to identify and prioritize the most relevant pieces of information from the entire input to inform each part of the output.

- **Number Theory Analogy:** When analyzing the properties of a large integer, one might look at its prime factorization (e.g., $N = p_1^{a_1} * p_2^{a_2} * \dots * p_k^{a_k}$). The properties of N (like its number of divisors or its value modulo another number) depend on all these prime factors collectively. The attention mechanism allows the model to consider all "elemental components" (words or tokens) of an input sequence simultaneously, much like considering all prime factors at once to understand the whole number, rather than processing information sequentially.

The Breakthrough: Why "Attention Is All You Need" Was Significant

The shift from recurrence to an attention-centric model was a breakthrough for several key reasons:

1. **Massive Parallelization and Speed:** By not being constrained to process words one after another, the Transformer can perform computations for all words in a sequence largely in parallel. This makes it exceptionally well-suited for modern hardware like GPUs, which excel at parallel processing.
 - **Analogy:** Imagine a team of collaborators working on different parts of a large mathematical problem simultaneously, rather than one person working through it step-by-step. This greatly accelerates the overall progress. Training Transformer models became significantly faster than training recurrent models of comparable capability.
2. **Superior Handling of Long-Range Dependencies:** Because the attention mechanism can directly connect any two words in a sequence, regardless of how far apart they are, Transformers are much better at capturing relationships between distant parts of the text. In recurrent models, the "signal" from an early word had to pass through many intermediate steps to influence a later word, often diminishing in strength.
 - **Analogy:** If a crucial lemma for a proof is stated on page 5 of a 100-page mathematical treatise, and it's needed for a step on page 95, an attention mechanism can directly "refer" back to page 5. A purely sequential processor might have "forgotten" the details of page 5 by the time it reaches page 95.

Impact of the Transformer Architecture

The introduction of the Transformer had a profound and wide-ranging impact:

- Revolutionized Machine Translation: The paper demonstrated that Transformer models could achieve state-of-the-art results in machine translation, often surpassing previous benchmarks in quality while requiring significantly less time to train.
- Foundation for Modern Large Language Models (LLMs): The Transformer architecture became the backbone for most modern LLMs. These models are capable of understanding and generating human-like text for a vast array of tasks, including question answering, text summarization, code generation, and creative writing.
- Transformed Human-AI Interaction: The capabilities unlocked by Transformers have led to more sophisticated and natural interactions with AI systems. Chatbots became more conversational, search engines could understand more nuanced queries, and tools for content creation and information retrieval became significantly more powerful.

In summary, the Transformer's core innovation was to show that by relying primarily on attention mechanisms and moving away from sequential recurrence, AI models could process language more efficiently (through parallelization) and more effectively (by better capturing long-range context). This paradigm shift has been instrumental in the rapid advancements seen in natural language processing and AI in recent years.