1. "Explain Like I'm Five"

Once upon a time, there were these really smart computer programs called Transformers that could do amazing things like translating between different languages. What made the Transformers so special was their cool "self-attention" trick that let them understand how all the words in a sentence are connected to each other.

Unlike older programs that had to read sentences word by word or use math tricks like "convolutions", the Transformers could look at the whole sentence at once. They were like super smart kids reading the whole homework problem first, then paying close attention to the most important parts and making sure to connect those parts properly.

The Transformer programs had two main parts - an Encoder that read the words and turned them into secret codes the computer could understand, and a Decoder that took those codes and guessed the next word over and over to make new sentences in the other language.

Both the Encoder and Decoder had those special "self-attention layers" inside them. In the Encoder, these layers let the program look at all the words and figure out how they related to each other. The Decoder also had self-attention layers, plus special attention layers that looked at what the Encoder did to help figure out what words should come next.

The Transformers used a clever "positional encoding" trick with sine and cosine functions to keep track of the order of words by giving each word a secret code telling the computer where it came in the sequence. They also had "multi-head attention" that let them attend to information from different spaces in parallel.

To train these amazing translators, the scientists used huge datasets with millions of example sentences pairs like English-German and English-French. They fed these into super powerful computers with fast video card "brains" in batches of around 25,000 words at a time.

The smaller Transformers took about half a second per training step over 12 hours, while the bigger ones took a full second but trained for 3 and a half days straight! During training, techniques like varying the learning rate, adding some mess to the examples, and letting the Transformers be a bit unsure helped them learn faster and become more accurate.

After all that hard work, the Transformer programs became the best ever at translating languages - even better than human experts! The big Transformers with more math power and "self-attention heads" performed the best.

Thanks to their self-attention abilities, the Transformers could relate words no matter how far apart they were and really understand what was important in the sentences. This self-attention trick made them awesome at learning languages, kind of like how kids pick up their native

tongues from tons of examples. That's the amazing story of the super smart, self-attentive Transformer translators!

2. "Explain Like I'm a College Student"

The transformer model represents a breakthrough in sequence-to-sequence modeling tasks like machine translation by eliminating the sequential computation bottleneck of previous approaches like recurrent neural networks (RNNs) and convolutional neural networks (CNNs). These traditional models face limitations in handling long-range dependencies due to their sequential nature - the number of operations required to relate signals from distant positions grows with their separation in RNNs, while CNNs have limited receptive field sizes despite more parallelization.

The transformer's key innovation is the self-attention mechanism that directly models relationships between all positions in an input sequence in parallel. Self-attention, or intra-attention, computes representations of a sequence by relating different positions within it. Unlike RNNs and CNNs that sequentially process elements, self-attention connects all positions with a constant number of operations, enabling more efficient computation and easier learning of long-range dependencies.

The transformer follows an encoder-decoder architecture. The encoder consists of stacked self-attention and feed-forward layers that map an input sequence to continuous representations. The decoder also contains stacked self-attention and feed-forward layers, generating an output sequence by attending to the encoder's output and previously generated symbols. Positional encodings inject position information into the model.

Self-attention has several advantages over RNNs and CNNs. For sentence representations, self-attention layers are faster than recurrent layers when sequence length is smaller than representation dimensionality. Convolutional layers, though more parallelizable than recurrent layers, can have high complexity reduced by separable convolutions. Self-attention also yields more interpretable models, with individual attention heads exhibiting behavior related to linguistic structures.

The transformer was evaluated on machine translation and constituency parsing tasks. On WMT 2014 English-to-German translation, the transformer (big) achieved a new state-of-the-art BLEU score of 28.4, outperforming previous best models by over 2.0 BLEU. On WMT 2014 English-to-French, it achieved 41.0 BLEU, outperforming previous single models. It also generalized well to English constituency parsing on the WSJ dataset, outperforming some prior state-of-the-art models in a semi-supervised setting.

Techniques like beam search, checkpoint averaging, dropout, and label smoothing helped improve transformer performance. The authors experimented with architectural variations like

attention head count, key/value dimensions, and model sizes. The transformer's compelling performance across tasks demonstrates its versatility and the advantages of self-attention for sequence modeling - improved efficiency, parallelization, and long-range dependency handling.

Overall, the transformer establishes self-attention as a powerful alternative to sequential computations in RNNs and CNNs. By directly connecting sequence positions in parallel, self-attention reduces sequential operations, making the transformer highly parallelizable and adept at capturing long-range dependencies. Its novel architecture opens compelling avenues for advancing sequence modeling across applications like machine translation, summarization, and language understanding.

3. "Explain Like I'm an Expert"

The Transformer is a groundbreaking neural network architecture that completely eliminates the sequential computation inherent in recurrent neural networks (RNNs) and convolutional neural networks (CNNs). Previous approaches like ByteNet and ConvS2S employed parallelization but still suffered from a linear or logarithmic growth in computational complexity when relating different input/output positions, making it difficult to learn long-range dependencies effectively.

The Transformer circumvents this issue by solely relying on self-attention mechanisms to model relationships between all input/output positions in parallel. Self-attention, also known as intra-attention, computes representations for a single sequence by relating different positions within it. Unlike recurrent layers requiring O(n) sequential operations for sequences of length n, self-attention layers have a constant number of operations, making them computationally faster and more parallelizable, especially for longer sequences.

The Transformer follows an encoder-decoder architecture, with both components composed of stacked layers containing multi-head self-attention and feed-forward sub-layers. Multi-head attention projects queries, keys, and values into different representation subspaces and employs scaled dot-product attention. It is applied in three ways: encoder self-attention (relating input positions), decoder self-attention (relating target positions), and encoder-decoder attention (relating targets to inputs). Positional encodings are added to input embeddings to represent sequence order.

During training, techniques like residual dropout, label smoothing, and dropout on sums of embeddings/positional encodings were used for regularization. The Adam optimizer with a warmup-and-inverse-square-root learning rate schedule was employed.

In machine translation tasks, the Transformer (big) achieved new state-of-the-art BLEU scores of 28.4 (English-to-German) and 41.0 (English-to-French), outperforming previous best models while requiring significantly lower training costs. Increasing model size, using more attention heads (optimal was 8), larger key/value dimensions, and residual dropout improved performance.

The Transformer also generalized well to English constituency parsing, achieving a new state-of-the-art 92.7 F1 score in a semi-supervised setting on the Penn Treebank dataset.

Overall, the Transformer demonstrates the viability of entirely replacing sequence-based RNNs and convolutions with a fully self-attentive approach, offering key advantages: constant computational complexity for connecting positions, complete parallelization of operations, and the shortest possible path length for modeling long-range dependencies. It paves the way for more efficient and effective neural sequence modeling across various tasks.