# 1 Attention

**Motivation**: One of the core mechanisms inside of current LLMs.

**Warning**: Convolution takes in a single node, while attention takes in all nodes.

## 1.1 Transformer Architecture

**Notes**:
- **Transformer Layer:**
    - Attention mechanism (multi-headed)
    - Positional encodings
- With massive unsupervised datasets:
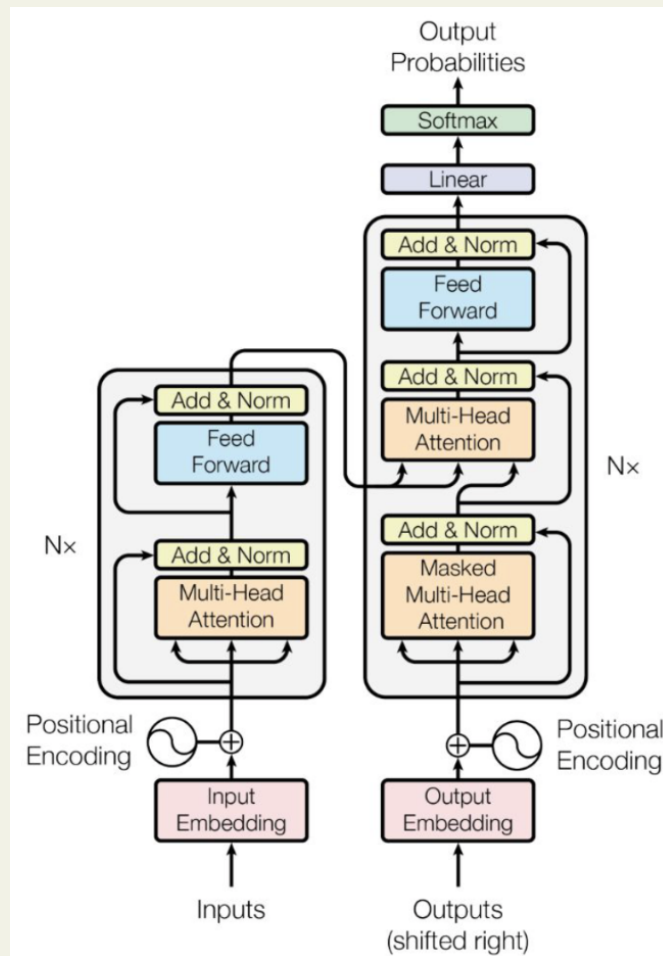    - Masked self-supervised training
    - Contrastive training



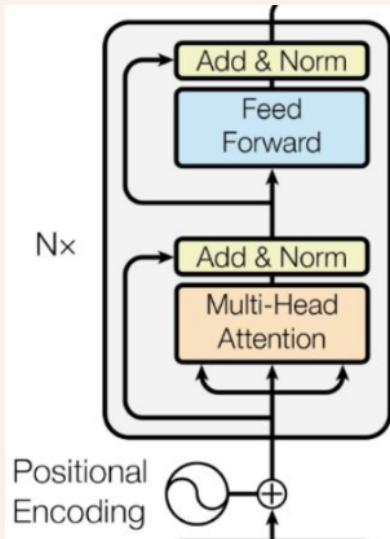Figure 1: LS: Encoder, RS: Decoder

### 1.1.1   Transformer Layer

**Summary**:



Figure 2

| Component | Description |
|---|---|
| Positional encoding | Learn to map integer positions into a vectorized representation. $$\mathrm{PE}_{(\mathrm{pos},i)} = \begin{cases} \sin\left(\dfrac{\mathrm{pos}}{10000^{\frac{i}{d_{model}}}}\right) & \text{if } i \text{ is even} \\ \cos\left(\dfrac{\mathrm{pos}}{10000^{\frac{i-1}{d_{model}}}}\right) & \text{if } i \text{ is odd} \end{cases}$$ |
| Multi-Head Attention | Computes attention scores for each token in the sequence. |
| LayerNorm | Stabilizes activations and accelerates training. |
| Residual Connection | Preserve information and enable deeper networks. |
| FFN/MLP | Increases the expressive power of the learned representation, often using GELU activations. |

## 1.2   Transformers are GNNs

**Summary**: Transformers are a special case of GNN

|  | GNN | Transformer |
|---|---|---|
| **Connectivity (Adjacency)** | Sparse | Full |
| **Edge Learning** | Yes | No (Implicitly) |
| **Message Computation** | $M(n_i, n_j, e_{ij})$ | $\langle n_i, n_j \rangle$ |
| **Communication per step** | # Number of Neighboring nodes | $\sim$ # Number of Heads |
| **Data requirements** | Low | High |
| **Computation** | Slow due to gather operations | Fast, Optimizable $\sim$ Matrix Multiplications |
| **Training** | Straightforward | Pre-training is needed |

## 1.3   Attention Mechanism

**Process**:
1. **Inputs:** Tokens tensor, Mask
   - **Tokens:** Inputs for Transformer/Attention Layers, which is a numerical representation of pieces of data.
   - **Mask:** A binary matrix that indicates which tokens to give attention to.
2. **Preprocessing:** Linear maps Tokens into Queries, Keys, and Values.
   - $Q = \text{Tokens} \cdot W_Q$: Represents the current token's context.
   - $K = \text{Tokens} \cdot W_K$: Represents the context of all tokens.
   - $V = \text{Tokens} \cdot W_V$: Represents the information to be passed on.
3. **Attention scores:** $\text{Scores} = \dfrac{QK^T}{\sqrt{d_k}} \cdot \text{Mask}$
   - $\text{score}_{ij} = \dfrac{(q_i \cdot k_j)m_{ij}}{\sqrt{d_k}}$
4. **Attention Normalization:** Attention Weights = softmax(scores)
   - $\text{score}_{ij}^{\text{normalized}} = \dfrac{\exp(\text{score}_{ij})}{\sum_{k=1}^{n} \exp(\text{score}_{ik})}$
5. **Value update:** New Values = Attention Weights $\cdot V$
   - $v_i^{\text{new}} = \sum_{j=1}^{n} \text{score}_{ij}^{\text{normalized}} v_j$
6. **Post Processing:** Apply LayerNorm, Residual connections, and a FFN.
7. **Outputs:** Updated tokens tensor

$$\text{Attention}(Q, K, V, M) = \text{softmax}\left(\frac{QK^T M}{\sqrt{d_k}}\right) V \tag{1}$$

### 1.3.1   Attention Maps: Visualizing Where the Model Attends

**Notes**: Softmax bias:
- Values between 0 and 1
- Categorical like
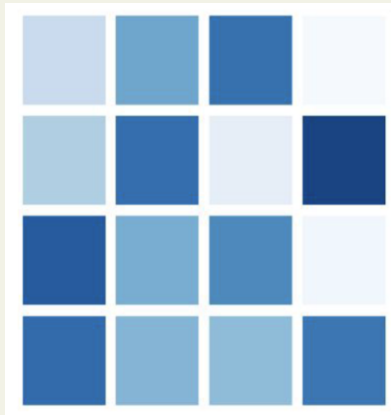- Attend to one token at a time



Figure 3

### 1.3.2   Self-Attention vs. Cross-Attention

**Notes**:
- **Self-Attention:** Attention allows connections between the same sequence without masking.

$$\text{Self-Attention}(x, \text{mask}) = \text{Attention}(\text{Linear}(x), \text{Linear}(x), \text{Linear}(x), \text{mask})$$
$$= \text{Attention}(Q(x), K(x), V(x), \text{mask})$$

- **Cross-Attention:** Attention allows connections between different sequences.

$$\text{Cross-Attention}(x, x', \text{mask}) = \text{Attention}(\text{Linear}(x), \text{Linear}(x'), \text{Linear}(x'), \text{mask})$$
$$= \text{Attention}(Q(x), K(y), V(x), \text{mask})$$

### 1.3.3   Multi-Headed Attention

**Notes**: Multiple attention mechanisms in parallel, each with different linear maps.
- Ensemble-like approach.
- Same compute and # of parameters.
- **Strategy:** Expand and contract tensors to new axis: number of heads.
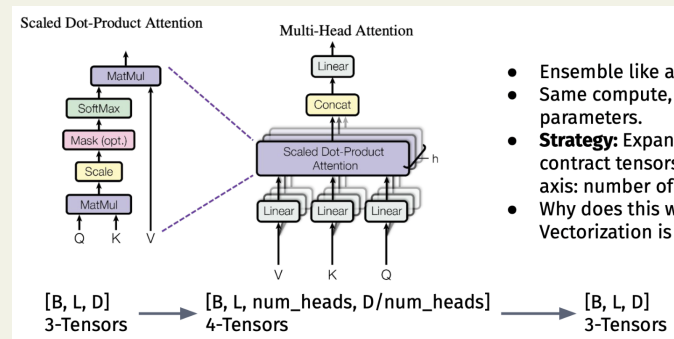


Figure 4

## 1.4 Examples
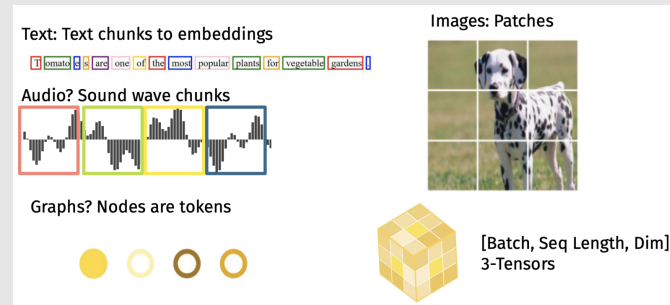
### 1.4.1 Tokens

**Example**:



Figure 5

### 1.4.2 Positional Encoding

**Example**:



$$\vec{p_t} = \begin{bmatrix} sin(w_1 t) \\ cos(w_1 t) \\ sin(w_2 t) \\ cos(w_2 t) \\ ... \\ sin(w_{d/2} t) \\ cos(w_{d/2} t) \end{bmatrix}$$
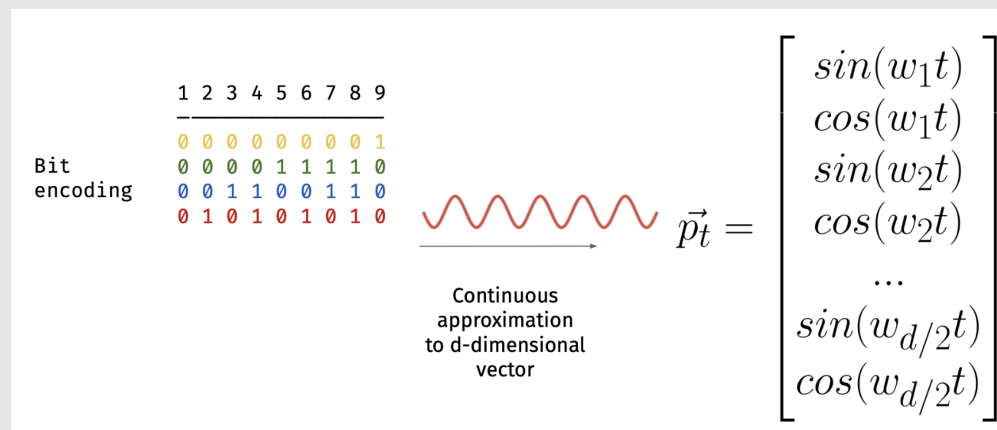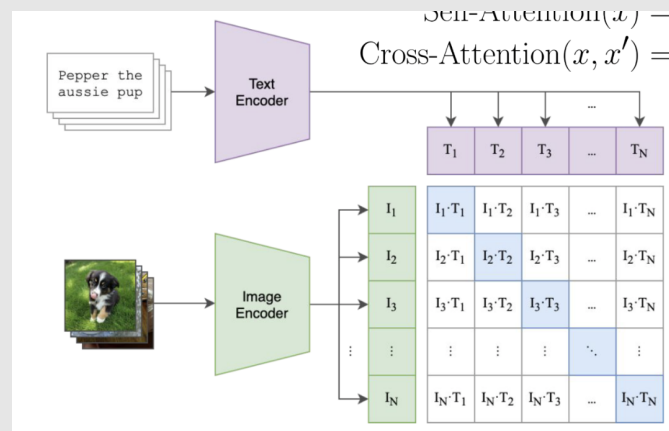
Figure 6

### 1.4.3   Cross-Attention

**Example:**



Figure 7

# 2 LLMs

**Notes**:
- Transformers on large text-like datasets.
- Transformers on "tokens" (discretized data)
- Foundational models

## 2.1 Transformers & LLMs

**Summary**:

| Concept | Description |
| --- | --- |
| Sizes of text datasets | Scale significantly impacts LLM's performance. |
| Text to Text Tasks | Many tasks can be framed as text-to-text problems. |
| Transformers & Masking: Encoders & Decoders | |
| Masked Language Modelling | Predicting masked tokens in a sequence. |
| Inputs | Tokenizing Text & Embedding Layers |
| Outputs | Auto-Regressive Decoding of Tokens<br>Decoding one token at a time, using previous outputs |

## 2.2   Scaling LLMs

**Motivation**:

### 2.2.1   Techniques

**Summary**: Table format

### 2.2.2   High-Level Impacts

**Summary**: