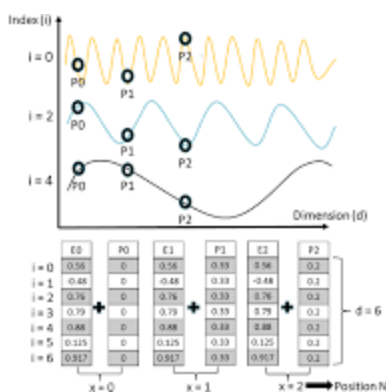# Positional Encoding methods



# Absolute Position Encoding



## 1. What is Absolute Position Encoding?

Absolute Position Encoding is a method of incorporating positional information into transformer models by assigning a unique vector to each position in the sequence. This approach allows the model to distinguish between tokens based on their absolute position within the input.

# 2. Types of Absolute Position Encoding

## 2.1 Sinusoidal Encoding (Vaswani et al., 2017)

This is the method used in the original Transformer paper.

**Mathematical Formulation:**

For position $pos$ and dimension $i$:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

**Pros:**

- Deterministic (no learning required)
- Can extrapolate to unseen sequence lengths
- Bounded values (-1 to 1)

**Cons:**

- Fixed pattern may not adapt to data-specific needs
- May not capture complex positional relationships

## 2.2 Learned Position Embeddings

This method involves learning a unique embedding for each position.

**Mathematical Formulation:**

$PE_{pos} = E_{pos}$, where $E$ is a learned embedding table

**Pros:**

- Can adapt to dataset-specific patterns
- Simple to implement

**Cons:**

- Limited to maximum sequence length seen during training
- Requires additional parameters to learn

## 2.3 Algebraic Encoding

This method uses algebraic formulations to create position encodings.

**Example Mathematical Formulation:**

$$PE_{(pos,i)} = pos/(max\_seq\_len - 1)^{i/(d_{model}-1)}$$

**Pros:**

- Deterministic
- Can handle arbitrary sequence lengths
- Smooth progression of values

**Cons:**

- May not capture complex positional relationships
- Fixed pattern may not adapt to all datasets

## 2.4 Binary Encoding

This method represents positions using their binary representations.

**Mathematical Formulation:**

$$PE_{(pos,i)} = \text{i-th bit of binary representation of pos}$$

**Pros:**

- Compact representation
- Deterministic
- Can handle long sequences with few dimensions

**Cons:**

- Discontinuous, which may affect model learning
- May require additional processing by the model

# 3. Why use Absolute Position Encoding?

- Uniqueness: Provides a distinct representation for each position
- Fixed pattern: The encoding is deterministic and doesn't require learning

- Boundedness: Values are always between -1 and 1, helping with numerical stability
- Extrapolation: Can generalize to sequences longer than those seen during training
- Efficiency: Can be pre-computed and stored for quick lookup

# 4. How does Absolute Position Encoding work?

1. Encoding generation: Generate a unique vector for each position using sinusoidal functions
2. Addition: Add the position encoding to the token embeddings
3. Model input: Feed the combined embeddings into the transformer model
4. Learning: The model learns to utilize this absolute positional information during training

# 5. Detailed Code Implementation

Here's a comprehensive implementation of Absolute Position Encoding, including visualization:

```python
import torch
import torch.nn as nn
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

class MultiMethodAbsolutePositionalEncoding(nn.Module):
    def __init__(self, d_model, max_seq_len=5000, dropout=0.1, method='sinusoidal'):
        super().__init__()
        self.d_model = d_model
        self.dropout = nn.Dropout(p=dropout)
        self.max_seq_len = max_seq_len
        self.method = method

        if method == 'sinusoidal':
            self.pe = self._sinusoidal_encoding()
        elif method == 'learned':
            self.pe = nn.Parameter(torch.randn(1, max_seq_len, d_model))
        elif method == 'algebraic':
            self.pe = self._algebraic_encoding()
        elif method == 'binary':
            self.pe = self._binary_encoding()
        else:
            raise ValueError(f"Unknown encoding method: {method}")

    def _sinusoidal_encoding(self):
        pe = torch.zeros(self.max_seq_len, self.d_model)
        position = torch.arange(0, self.max_seq_len, dtype=torch.float).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, self.d_model, 2).float() * (-np.log(10000.0) / sel
        pe[:, 0::2] = torch.sin(position * div_term)
        pe[:, 1::2] = torch.cos(position * div_term)
        return pe.unsqueeze(0)

    def _algebraic_encoding(self):
        pe = torch.zeros(self.max_seq_len, self.d_model)
        position = torch.arange(0, self.max_seq_len, dtype=torch.float).unsqueeze(1)
        for i in range(self.d_model):
            pe[:, i] = position.squeeze() / (self.max_seq_len - 1) ** (i / (self.d_model - 1))
        return pe.unsqueeze(0)

    def _binary_encoding(self):
        pe = torch.zeros(self.max_seq_len, self.d_model)
        for pos in range(self.max_seq_len):
```

```python
            for i in range(min(self.d_model, len(bin(self.max_seq_len)) - 2)):
                pe[pos, i] = (pos >> i) & 1
        return pe.unsqueeze(0)

    def forward(self, x):
        x = x + self.pe[:, :x.size(1)]
        return self.dropout(x)


class AbsolutePositionAttention(nn.Module):
    def __init__(self, d_model, num_heads, dropout=0.1, max_seq_len=5000, method='sinusoidal'):
        super().__init__()
        self.d_model = d_model
        self.num_heads = num_heads

        self.pos_encoding = MultiMethodAbsolutePositionalEncoding(d_model, max_seq_len, dropout,

        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.out_proj = nn.Linear(d_model, d_model)

        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        batch_size, seq_len, _ = x.shape

        # Add positional encoding to input
        x = self.pos_encoding(x)

        q = self.q_proj(x).view(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)
        k = self.k_proj(x).view(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)
        v = self.v_proj(x).view(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)

        attn = (q @ k.transpose(-2, -1)) / (self.d_model ** 0.5)
        attn = self.dropout(attn.softmax(dim=-1))

        out = (attn @ v).transpose(1, 2).reshape(batch_size, seq_len, self.d_model)
        return self.out_proj(out)

    def visualize_positional_encoding(self, seq_len=100):
        pos_enc = self.pos_encoding.pe[:, :seq_len, :].squeeze().detach().cpu().numpy()

        plt.figure(figsize=(12, 6))
```

```python
            sns.heatmap(pos_enc, cmap='coolwarm')
            plt.title(f"Absolute Positional Encoding ({self.pos_encoding.method})")
            plt.xlabel("Encoding Dimension")
            plt.ylabel("Position")
            plt.show()

            plt.figure(figsize=(12, 6))
            for i in range(0, pos_enc.shape[1], pos_enc.shape[1]//10):
                plt.plot(pos_enc[:, i], label=f'Dim {i}')
            plt.title(f"Absolute Positional Encoding Components ({self.pos_encoding.method})")
            plt.xlabel("Position")
            plt.ylabel("Encoding Value")
            plt.legend()
            plt.show()

if __name__ == "__main__":

    d_model = 512
    num_heads = 8
    max_seq_len = 1000

    methods = ['sinusoidal', 'learned', 'algebraic', 'binary']

    for method in methods:
        print(f"\nTesting {method.capitalize()} Encoding:")
        abs_pos_attn = AbsolutePositionAttention(d_model, num_heads, max_seq_len=max_seq_len, me

        x = torch.randn(32, 100, d_model)
        output = abs_pos_attn(x)

        print(f"Input shape: {x.shape}")
        print(f"Output shape: {output.shape}")

    abs_pos_attn.visualize_positional_encoding()
```

This implementation includes:

- An `AbsolutePositionalEncoding` module that generates the sinusoidal position encodings
- An `AbsolutePositionAttention` module that applies the encoding in the attention mechanism
- Visualization methods to display the encoding patterns

# 6. Pros of Absolute Position Encoding

1. Simplicity: Easy to implement and understand
2. Deterministic: No need for learning positional embeddings
3. Extrapolation: Can handle sequences longer than those seen during training
4. Efficiency: Can be pre-computed and stored for quick lookup
5. Stability: Bounded values help with numerical stability

# 7. Cons of Absolute Position Encoding

1. Fixed representation: May not adapt to specific patterns in the data
2. Limited expressiveness: May not capture complex positional relationships
3. No relative information: Doesn't explicitly encode relative positions between tokens
4. Potential interference: Adding position encodings might interfere with token embeddings
5. Hyperparameter sensitivity: Performance may depend on the choice of frequency scaling
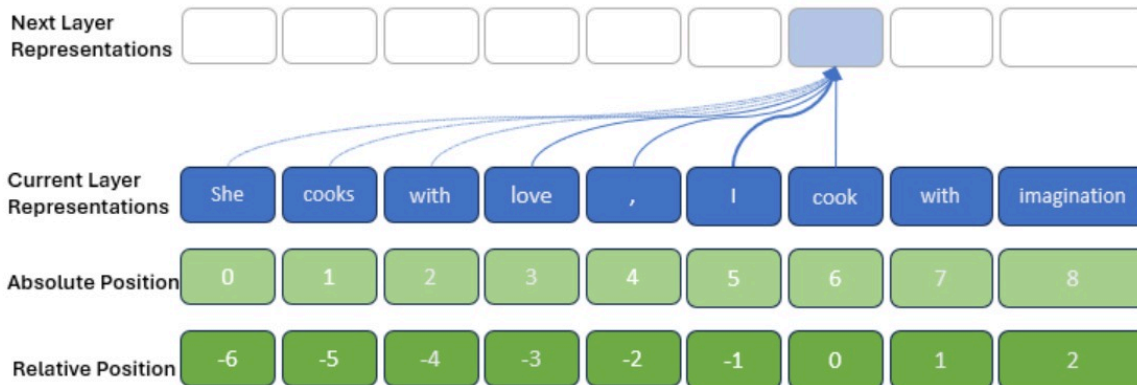
# 8. Comparison with Other Methods

- vs. Learned Position Embeddings: Absolute encoding doesn't require learning, which can be an advantage for generalization but may lack adaptability to specific datasets.
- vs. Relative Position Encoding: Absolute encoding provides unique representations for each position but doesn't directly encode relationships between positions.
- vs. RoPE: Absolute encoding is simpler but may not capture relative information as effectively as RoPE.
- vs. ALiBi: Absolute encoding provides more detailed positional information but may not scale as well to very long sequences.

The choice between these methods often depends on the specific requirements of the task, the model architecture, and the characteristics of the data being processed. Absolute Position Encoding is particularly well-suited for tasks where the absolute position of tokens is crucial and where a simple, deterministic encoding is preferred.

# Relative Position Encoding



## 1. What is Relative Position Encoding?

Relative Position Encoding (RPE) is a method of incorporating positional information into transformer models by encoding the relative distances between tokens rather than their absolute positions. This approach allows the model to focus on the relationships between tokens, which can be particularly useful for tasks that depend on local context or where the absolute position is less important.

## 2. Mathematical Formulation

There are several ways to implement Relative Position Encoding. We'll focus on two popular methods:

### 2.1 Shaw et al. (2018) Method

In this method, relative positions are clipped to a maximum distance.

For tokens at positions $i$ and $j$:

$$a_{ij} = (W_Q x_i)^T W_K^E x_j + (W_Q x_i)^T W_K^R r_{ij}$$

Where:

- $W_Q, W_K^E, W_K^R$ are learned weight matrices
- $x_i, x_j$ are token embeddings
- $r_{ij}$ is a learned embedding for the relative position $(i - j)$, clipped to a maximum absolute value

## 2.2 Transformer-XL Method (Dai et al., 2019)

This method uses sinusoidal functions to encode relative positions:

$$R_{i-j} = [R_{i-j}^{(d)}, R_{i-j}^{(d+1)}, ..., R_{i-j}^{(d+d'-1)}]$$

Where:
$$R_{i-j}^{(2k)} = \sin((i-j)/10000^{2k/d})$$
$$R_{i-j}^{(2k+1)} = \cos((i-j)/10000^{2k/d})$$

The attention score is then computed as:

$$a_{ij} = (W_Q x_i + u)^T W_K x_j + (W_Q x_i + v)^T R_{i-j}$$

Where $u$ and $v$ are learned vectors.

# 3. Why use Relative Position Encoding?

- Context focus: Emphasizes relationships between tokens rather than absolute positions
- Translation invariance: The model can learn patterns that are independent of absolute position
- Long-range dependencies: Can potentially capture long-range dependencies more effectively
- Extrapolation: Often generalizes better to sequence lengths not seen during training
- Efficiency: Can be more parameter-efficient than some absolute position encoding methods

# 4. How does Relative Position Encoding work?

1. Relative distance computation: Calculate the relative distances between all pairs of tokens
2. Encoding generation: Generate encodings for these relative distances (learned or function-based)
3. Attention computation: Incorporate the relative position information into the attention mechanism
4. Learning: The model learns to utilize this relative positional information during training

# 5. Detailed Code Implementation

Here's a comprehensive implementation of Relative Position Encoding, including visualization:

```python
import torch
import torch.nn as nn
import torch.nn.functional as F
import seaborn as sns
import matplotlib.pyplot as plt
import math

class RelativeGlobalAttention(nn.Module):
    def __init__(self, d_model, num_heads, max_len=1024, dropout=0.1):
        super().__init__()
        d_head, remainder = divmod(d_model, num_heads)
        if remainder:
            raise ValueError(
                "incompatible `d_model` and `num_heads`"
            )
        self.max_len = max_len
        self.d_model = d_model
        self.num_heads = num_heads
        self.key = nn.Linear(d_model, d_model)
        self.value = nn.Linear(d_model, d_model)
        self.query = nn.Linear(d_model, d_model)
        self.dropout = nn.Dropout(dropout)
        self.Er = nn.Parameter(torch.randn(max_len, d_head))
        self.register_buffer(
            "mask",
            torch.tril(torch.ones(max_len, max_len))
            .unsqueeze(0).unsqueeze(0)
        )
        # self.mask.shape = (1, 1, max_len, max_len)

    def forward(self, x):
        batch_size, seq_len, _ = x.shape

        if seq_len > self.max_len:
            raise ValueError(
                "sequence length exceeds model capacity"
            )

        k_t = self.key(x).reshape(batch_size, seq_len, self.num_heads, -1).permute(0, 2, 3, 1)
        v = self.value(x).reshape(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)
        q = self.query(x).reshape(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)

        start = self.max_len - seq_len
```

```python
        Er_t = self.Er[start:, :].transpose(0, 1)
        QEr = torch.matmul(q, Er_t)
        Srel = self.skew(QEr)

        QK_t = torch.matmul(q, k_t)
        attn = (QK_t + Srel) / math.sqrt(q.size(-1))
        mask = self.mask[:, :, :seq_len, :seq_len]
        attn = attn.masked_fill(mask == 0, float("-inf"))
        attn = F.softmax(attn, dim=-1)
        out = torch.matmul(attn, v)
        out = out.transpose(1, 2).reshape(batch_size, seq_len, -1)

        return self.dropout(out), attn  # Returning attention for visualization

    def skew(self, QEr):
        padded = F.pad(QEr, (1, 0))
        batch_size, num_heads, num_rows, num_cols = padded.shape
        reshaped = padded.reshape(batch_size, num_heads, num_cols, num_rows)
        Srel = reshaped[:, :, 1:, :]
        return Srel

    def visualize_attention(self, attn, seq_len):
        """
        Visualizes attention matrix for a given sequence length
        Parameters:
        - attn: Attention weights tensor with shape (batch_size, num_heads, seq_len, seq_len)
        - seq_len: Length of the sequence
        """
        # Pick the attention weights for the first head and first batch (for simplicity)
        attn_weights = attn[0, 0, :, :].detach().cpu().numpy()

        # Create a heatmap for the attention weights
        plt.figure(figsize=(10, 8))
        sns.heatmap(attn_weights, cmap='coolwarm', xticklabels=range(seq_len), yticklabels=range
        plt.title('Attention Weights Visualization')
        plt.xlabel('Key Positions')
        plt.ylabel('Query Positions')
        plt.show()

if __name__ == "__main__":

    seq_len = 10
    d_model = 64
```

```
    num_heads = 4

    model = RelativeGlobalAttention(d_model=d_model, num_heads=num_heads, max_len=1024)
    x = torch.randn(1, seq_len, d_model)

    out, attn = model(x)
    model.visualize_attention(attn, seq_len)
```

# 6. Pros of Relative Position Encoding

1. Translation invariance: Patterns learned are independent of absolute position
2. Long-range dependencies: Can potentially capture long-range relationships more effectively
3. Generalization: Often performs better on sequences longer than those seen during training
4. Context focus: Emphasizes relationships between tokens rather than their absolute positions
5. Parameter efficiency: Can be more efficient than some absolute position encoding methods
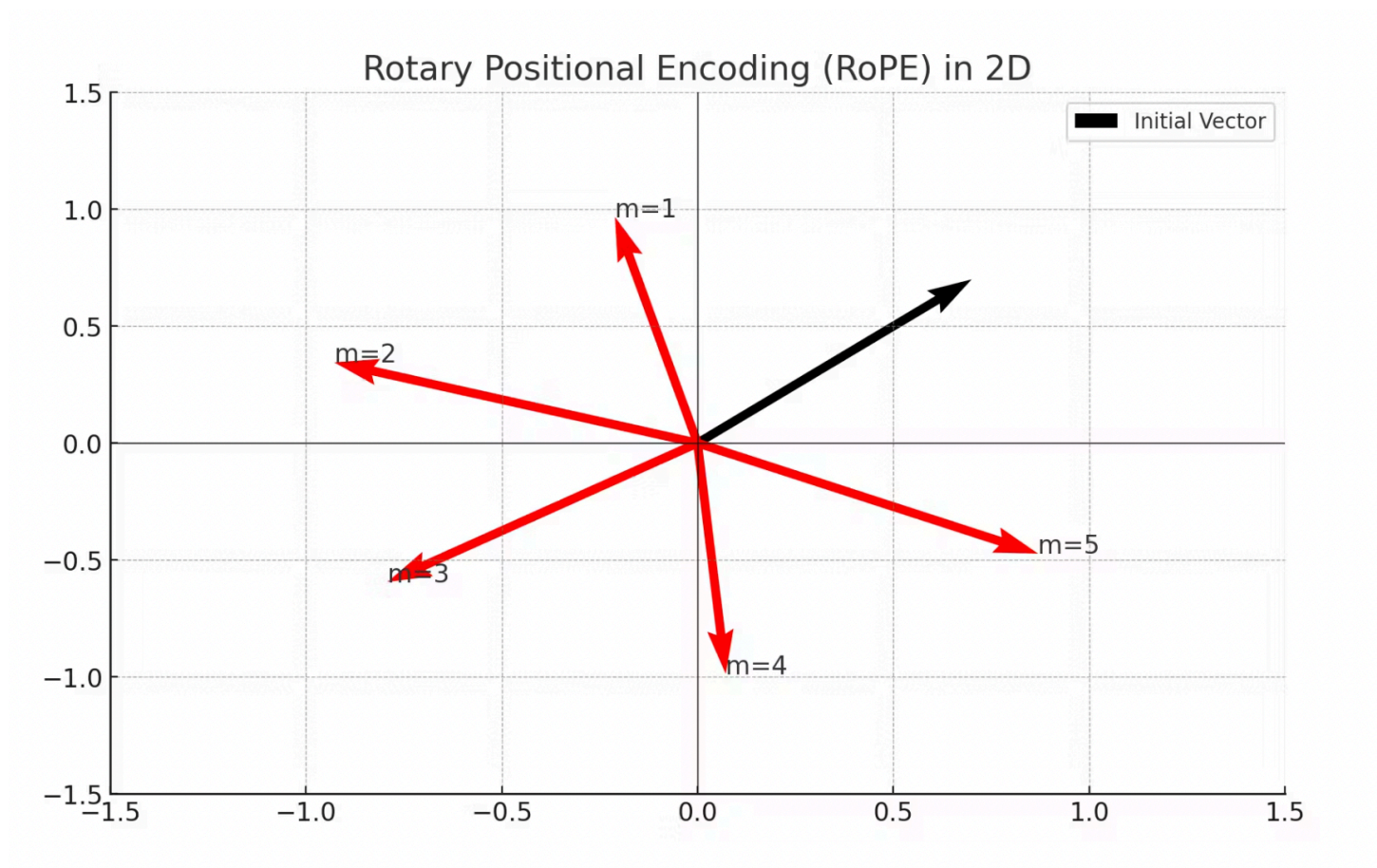
# 7. Cons of Relative Position Encoding

1. Computational complexity: Can be more computationally expensive, especially for long sequences
2. Implementation complexity: More complex to implement than absolute position encoding
3. Loss of absolute position information: May not be suitable for tasks where absolute position is crucial
4. Memory usage: Can require more memory, especially for long sequences
5. Hyperparameter sensitivity: Performance may depend on the maximum relative distance considered

# 8. Comparison with Other Methods

- vs. Absolute Position Encoding: RPE focuses on token relationships rather than absolute positions, which can be beneficial for many tasks but may lose some global position information.
- vs. Sinusoidal Encoding: RPE can capture more complex positional relationships and is often more flexible.
- vs. Learned Absolute Embeddings: RPE can generalize better to unseen sequence lengths but may be more complex to implement.
- vs. RoPE: Both methods encode relative positions, but RPE typically uses separate relative position embeddings, while RoPE integrates the information directly into token representations.

The choice between these methods depends on the specific task, model architecture, and computational constraints. Relative Position Encoding is particularly well-suited for tasks where the relationships between tokens are more important than their absolute positions, and where generalization to different sequence lengths is crucial.

# Rotary Position Embedding



Rotary Positional Encoding (RoPE) in 2D

# 1. What is RoPE?

RoPE, or Rotary Position Embedding, is a method of incorporating positional information into transformer models by applying rotation matrices to token embeddings. Introduced by Su et al. in 2021, RoPE aims to combine the benefits of absolute and relative position encoding while maintaining computational efficiency.

# 2. Mathematical Formulation

Let's define RoPE mathematically:

1. For a complex-valued vector $x = [x_1, x_2, ..., x_{d/2}]$ where each $x_i$ is a complex number, we define a rotation operation:
$R_m(x) = [x_1 e^{im\theta_1}, x_2 e^{im\theta_2}, ..., x_{d/2} e^{im\theta_{d/2}}]$
Where $m$ is the position index, and $\theta_i = w_i^{-k}$, with $w_i$ and $k$ being hyperparameters.

2. In practice, for real-valued vectors, this rotation is implemented as:
$R_m(x_{2i-1}, x_{2i}) = [x_{2i-1}\cos(m\theta_i) - x_{2i}\sin(m\theta_i), x_{2i-1}\sin(m\theta_i) + x_{2i}\cos(m\theta_i)]$
Sample Rotary Matrix:

$$R_{\Theta,m} = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix}$$

3. In the attention mechanism, for query $q$ and key $k$ at positions $m$ and $n$ respectively:
$q^T k = R_m(q)^T R_n(k) = q^T R_{m-n}(k)$

# 3. Why use RoPE?

1. **Unified absolute and relative encoding**: RoPE naturally captures both absolute and relative position information.
2. **Flexibility**: Can be easily integrated into existing transformer architectures with minimal changes.
3. **Theoretical properties**: Preserves the linear self-attention formulation while encoding position information.
4. **Efficiency**: Computationally efficient, as it can be implemented with simple matrix operations.
5. **Extrapolation**: Demonstrates good performance on sequences longer than those seen during training.

# 4. How does RoPE work?

1. **Rotation generation**: For each position, generate a rotation matrix based on the position and a set of frequencies.
2. **Embedding rotation**: Apply these rotations to the query and key vectors in the attention mechanism.
3. **Attention computation**: Compute attention scores using the rotated vectors. The resulting dot products naturally encode relative position information.
4. **Training**: The model learns to utilize this rotated space during the training process.

# 5. Detailed Code Implementation:

Here's a comprehensive implementation of RoPE, including visualization:

```python
import torch
import torch.nn as nn
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

class RotaryPositionalEmbedding(nn.Module):
    def __init__(self, d_model, max_seq_len=5000, base=10000):
        super().__init__()
        self.d_model = d_model
        inv_freq = 1. / (base ** (torch.arange(0, d_model, 2).float() / d_model))
        self.register_buffer('inv_freq', inv_freq)
        self.max_seq_len = max_seq_len

    def forward(self, positions):
        sinusoid_inp = torch.einsum("i,j->ij", positions.float(), self.inv_freq)
        emb = torch.cat((sinusoid_inp.sin(), sinusoid_inp.cos()), dim=-1)
        return emb

def rotate_half(x):
    x1, x2 = x[..., :x.shape[-1]//2], x[..., x.shape[-1]//2:]
    return torch.cat((-x2, x1), dim=-1)

class RotaryAttention(nn.Module):
    def __init__(self, d_model, num_heads, max_seq_len=5000):
        super().__init__()
        self.d_model = d_model
        self.num_heads = num_heads
        head_dim = d_model // num_heads
        self.head_dim = head_dim
        self.rotary_emb = RotaryPositionalEmbedding(head_dim, max_seq_len)

        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.out_proj = nn.Linear(d_model, d_model)

    def forward(self, x):
        batch_size, seq_len, _ = x.shape
        q = self.q_proj(x).view(batch_size, seq_len, self.num_heads, self.head_dim).transpose(1,
        k = self.k_proj(x).view(batch_size, seq_len, self.num_heads, self.head_dim).transpose(1,
        v = self.v_proj(x).view(batch_size, seq_len, self.num_heads, self.head_dim).transpose(1,
```

```python
        positions = torch.arange(seq_len, device=x.device)
        rot_emb = self.rotary_emb(positions)

        q_rot = (q * rot_emb.cos()) + (rotate_half(q) * rot_emb.sin())
        k_rot = (k * rot_emb.cos()) + (rotate_half(k) * rot_emb.sin())

        attn = (q_rot @ k_rot.transpose(-2, -1)) / (self.head_dim ** 0.5)
        attn = attn.softmax(dim=-1)

        out = (attn @ v).transpose(1, 2).reshape(batch_size, seq_len, self.d_model)
        return self.out_proj(out)

    def visualize_rotary_embeddings(self, seq_len=20):
        positions = torch.arange(seq_len)
        rot_emb = self.rotary_emb(positions).detach().cpu().numpy()

        plt.figure(figsize=(12, 6))
        plt.subplot(1, 2, 1)
        sns.heatmap(rot_emb[:, :self.head_dim//2], cmap='coolwarm')
        plt.title("Sine Component of RoPE")
        plt.xlabel("Dimension")
        plt.ylabel("Position")

        plt.subplot(1, 2, 2)
        sns.heatmap(rot_emb[:, self.head_dim//2:], cmap='coolwarm')
        plt.title("Cosine Component of RoPE")
        plt.xlabel("Dimension")
        plt.ylabel("Position")

        plt.tight_layout()
        plt.show()


if __name__ == "__main__":
    d_model = 512
    num_heads = 8
    max_seq_len = 1000
    rope_attn = RotaryAttention(d_model, num_heads, max_seq_len)

    x = torch.randn(32, 100, d_model)
    output = rope_attn(x)

    print(f"Input shape: {x.shape}")
    print(f"Output shape: {output.shape}")
```

```
rope_attn.visualize_rotary_embeddings()
```

This implementation includes:

- A `RotaryPositionalEmbedding` module that generates the rotary embeddings
- A `RotaryAttention` module that applies RoPE in the attention mechanism
- A visualization method to display the sine and cosine components of the rotary embeddings

# 6. Pros of RoPE:

1. **Unified encoding**: Captures both absolute and relative position information in a single mechanism.
2. **Efficiency**: Can be implemented with simple matrix operations, making it computationally efficient.
3. **Extrapolation**: Shows good performance on sequences longer than those seen during training.
4. **Theoretical guarantees**: Preserves the linear self-attention formulation while encoding positions.
5. **Flexibility**: Can be easily integrated into existing transformer architectures with minimal changes.

# 7. Cons of RoPE:

1. **Complexity**: The concept and implementation can be more complex to understand compared to simpler encoding methods.
2. **Hyperparameter sensitivity**: The performance may be sensitive to the choice of base frequency.
3. **Limited interpretability**: The rotated space can be less intuitive to interpret compared to direct positional encodings.
4. **Potential numerical issues**: For very long sequences, there might be numerical precision issues with the rotation angles.
5. **Task-specific performance**: While generally effective, its benefits may vary depending on the specific task and dataset.

# Mixture Position Encoding (Absolute + RPE)

## 1. What is Mixture Position Encoding?

Mixture Position Encoding is an approach that combines multiple position encoding methods to capture different aspects of positional information. In this guide, we'll focus on a mixture of Absolute Position Encoding and Relative Position Encoding (RPE). This combination aims to leverage the strengths of both methods:

- Absolute Position Encoding provides a unique representation for each position in the sequence.
- Relative Position Encoding captures the relationships between different positions.

## 2. Mathematical Formulation

Let's define the Mixture Position Encoding mathematically:

For a sequence of length $L$, we define the position encoding for position $i$ as:

$$PE_{mixture}(i) = \alpha \cdot PE_{absolute}(i) + (1 - \alpha) \cdot PE_{relative}(i)$$

Where:

- $PE_{absolute}(i)$ is the absolute position encoding for position $i$
- $PE_{relative}(i)$ is the relative position encoding for position $i$
- $\alpha$ is a learnable parameter or a fixed hyperparameter that controls the mixture

### 2.1 Absolute Position Encoding

We'll use the sinusoidal encoding from the original Transformer paper:

$$PE_{absolute}(pos, 2i) = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{absolute}(pos, 2i + 1) = \cos(pos/10000^{2i/d_{model}})$$

Where $pos$ is the position and $i$ is the dimension.

### 2.2 Relative Position Encoding

We'll use a learned relative position encoding:

$$PE_{relative}(i, j) = E_{i-j}$$

Where $E$ is a learned embedding table, and $i - j$ represents the relative distance between positions $i$ and $j$.

# 3. Why use Mixture Position Encoding?

- Comprehensive representation: Captures both absolute and relative positional information
- Flexibility: Can adjust the importance of absolute vs. relative information
- Improved performance: Often leads to better results than using either method alone
- Generalization: Can potentially handle both short and long sequences more effectively
- Task adaptability: The mixture can be optimized for different tasks

# 4. How does Mixture Position Encoding work?

1. Absolute encoding generation: Generate absolute position encodings for each position
2. Relative encoding generation: Generate or look up relative position encodings for each pair of positions
3. Mixing: Combine the absolute and relative encodings using the mixture parameter
4. Attention computation: Incorporate the mixed encoding into the attention mechanism
5. Learning: The model learns to utilize both types of positional information during training

# 5. Detailed Code Implementation

Here's a comprehensive implementation of Mixture Position Encoding, including visualization:

```python
import torch
import torch.nn as nn
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

class MixturePositionalEncoding(nn.Module):
    def __init__(self, d_model, max_seq_len=5000, dropout=0.1):
        super().__init__()
        self.d_model = d_model
        self.max_seq_len = max_seq_len
        self.dropout = nn.Dropout(p=dropout)

        # Absolute position encoding
        pe = torch.zeros(max_seq_len, d_model)
        position = torch.arange(0, max_seq_len, dtype=torch.float).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-np.log(10000.0) / d_model))
        pe[:, 0::2] = torch.sin(position * div_term)
        pe[:, 1::2] = torch.cos(position * div_term)
        pe = pe.unsqueeze(0)
        self.register_buffer('pe', pe)

        # Relative position encoding
        self.relative_position_embedding = nn.Embedding(2 * max_seq_len - 1, d_model)

        # Mixture parameter
        self.alpha = nn.Parameter(torch.tensor(0.5))

    def forward(self, x):
        seq_len = x.size(1)

        # Absolute position encoding
        absolute_pe = self.pe[:, :seq_len]

        # Relative position encoding
        relative_pos = torch.arange(-seq_len + 1, seq_len, device=x.device).unsqueeze(0)
        relative_pe = self.relative_position_embedding(relative_pos + self.max_seq_len - 1)

        # Mix absolute and relative encodings
        mixed_pe = self.alpha * absolute_pe + (1 - self.alpha) * relative_pe[:, :seq_len]

        return self.dropout(mixed_pe)
```

```python
class MixtureAttention(nn.Module):
    def __init__(self, d_model, num_heads, dropout=0.1, max_seq_len=5000):
        super().__init__()
        self.d_model = d_model
        self.num_heads = num_heads
        self.max_seq_len = max_seq_len

        self.pos_encoding = MixturePositionalEncoding(d_model, max_seq_len, dropout)

        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.out_proj = nn.Linear(d_model, d_model)

        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        batch_size, seq_len, _ = x.shape
        pos_enc = self.pos_encoding(x)

        # Add positional encoding to input
        x = x + pos_enc

        q = self.q_proj(x).view(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)
        k = self.k_proj(x).view(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)
        v = self.v_proj(x).view(batch_size, seq_len, self.num_heads, -1).transpose(1, 2)

        attn = (q @ k.transpose(-2, -1)) / (self.d_model ** 0.5)
        attn = self.dropout(attn.softmax(dim=-1))

        out = (attn @ v).transpose(1, 2).reshape(batch_size, seq_len, self.d_model)
        return self.out_proj(out)

    def visualize_positional_encoding(self, seq_len=100):
        pos_enc = self.pos_encoding(torch.zeros(1, seq_len, self.d_model)).squeeze().detach().cp

        plt.figure(figsize=(12, 6))
        sns.heatmap(pos_enc, cmap='coolwarm')
        plt.title(f"Mixture Positional Encoding (α={self.pos_encoding.alpha.item():.2f})")
        plt.xlabel("Encoding Dimension")
        plt.ylabel("Position")
        plt.show()
```

```python
        plt.figure(figsize=(12, 6))
        for i in range(0, pos_enc.shape[1], pos_enc.shape[1]//10):
            plt.plot(pos_enc[:, i], label=f'Dim {i}')
        plt.title(f"Mixture Positional Encoding Components (α={self.pos_encoding.alpha.item():.2
        plt.xlabel("Position")
        plt.ylabel("Encoding Value")
        plt.legend()
        plt.show()

# Usage example
d_model = 512
num_heads = 8
max_seq_len = 1000
mixture_attn = MixtureAttention(d_model, num_heads, max_seq_len=max_seq_len)

# Simulate input
x = torch.randn(32, 100, d_model)  # Batch of 32, sequence length 100
output = mixture_attn(x)

print(f"Input shape: {x.shape}")
print(f"Output shape: {output.shape}")
print(f"Mixture parameter α: {mixture_attn.pos_encoding.alpha.item():.4f}")

# Visualize the mixture positional encoding
mixture_attn.visualize_positional_encoding()
```

This implementation includes:

- A `MixturePositionalEncoding` module that generates the mixed absolute and relative position encodings
- A `MixtureAttention` module that applies the mixed encoding in the attention mechanism
- Visualization methods to display the encoding patterns

# 6. Pros of Mixture Position Encoding

1. Comprehensive: Captures both absolute and relative positional information
2. Flexibility: Can adjust the balance between absolute and relative information
3. Improved performance: Often leads to better results than using either method alone
4. Adaptability: Can be optimized for different tasks and sequence lengths
5. Robustness: Less likely to overfit to specific positional patterns

# 7. Cons of Mixture Position Encoding

1. Complexity: More complex to implement and understand than single encoding methods
2. Computational overhead: May require more computation, especially for the relative component
3. Hyperparameter sensitivity: Performance may depend on the initial or learned mixture parameter
4. Potential conflicts: Absolute and relative information might conflict in some cases
5. Memory usage: Requires storing both absolute and relative encodings

# 8. Comparison with Other Methods

- vs. Absolute Position Encoding: Mixture encoding captures relative information, which absolute encoding lacks
- vs. Relative Position Encoding: Mixture encoding retains absolute position information, which can be crucial for some tasks
- vs. RoPE: Mixture encoding offers more flexibility in balancing absolute and relative information, while RoPE integrates them implicitly
- vs. ALiBi: Mixture encoding can capture more complex positional relationships, while ALiBi focuses on inducing a specific positional bias

The choice between these methods often depends on the specific requirements of the task, the model architecture, and the characteristics of the data being processed.

# Attention with Linear Biases (ALiBi)



# 1. What is ALiBi?

Attention with Linear Biases (ALiBi) is a novel approach to incorporating positional information in transformer models. Unlike traditional positional encoding methods, ALiBi introduces a simple yet effective bias term in the attention mechanism. This method aims to improve the model's ability to handle longer sequences and generalize to unseen sequence lengths.

# 2. Mathematical Formulation

ALiBi modifies the standard attention mechanism by adding a linear bias term to the attention scores. The ALiBi attention is defined as:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}} + m \cdot [-(i-j)]) \cdot V$$

Where:

- $Q$, $K$, and $V$ are the query, key, and value matrices
- $d_k$ is the dimension of the keys
- $m$ is a head-specific slope
- $i$ and $j$ are the positions in the sequence
- $[-(i-j)]$ is a matrix of relative positions

The slope $m$ for each head is calculated as:

$$m_h = \frac{1}{2^{\frac{8}{H} \cdot (h+1)}}$$

Where $H$ is the total number of heads and $h$ is the head index (0-indexed).

# 3. Why use ALiBi?

- Simplicity: ALiBi is straightforward to implement and understand
- Efficiency: It doesn't require additional parameters or computations for positional encoding
- Extrapolation: ALiBi can generalize to longer sequences than seen during training
- Performance: It often outperforms traditional positional encoding methods
- Memory efficiency: ALiBi doesn't require storing position embeddings

# 4. How does ALiBi work?

1. Attention computation: Calculate standard attention scores
2. Bias addition: Add the linear bias term to the attention scores
3. Softmax: Apply softmax to get attention weights
4. Value aggregation: Multiply attention weights with values
5. Output: Produce the final output of the attention mechanism

# 5. Detailed Code Implementation

Here's a comprehensive implementation of ALiBi, including visualization:

```python
import torch
import torch.nn as nn
import matplotlib.pyplot as plt
import seaborn as sns


class ALiBiAttention(nn.Module):
    def __init__(self, d_model, num_heads, dropout=0.1):
        super().__init__()
        self.d_model = d_model
        self.num_heads = num_heads
        self.head_dim = d_model // num_heads

        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.out_proj = nn.Linear(d_model, d_model)

        self.dropout = nn.Dropout(dropout)

        # Calculate ALiBi slopes
        m = torch.arange(1, self.num_heads + 1)
        m = 1.0 / (2 ** (8 * m / self.num_heads))
        self.register_buffer("m", m.view(1, self.num_heads, 1, 1))

    def forward(self, x):
        batch_size, seq_len, _ = x.shape

        q = self.q_proj(x).view(batch_size, seq_len, self.num_heads, self.head_dim).transpose(1,
        k = self.k_proj(x).view(batch_size, seq_len, self.num_heads, self.head_dim).transpose(1,
        v = self.v_proj(x).view(batch_size, seq_len, self.num_heads, self.head_dim).transpose(1,

        # Calculate attention scores
        attn_scores = (q @ k.transpose(-2, -1)) / (self.head_dim ** 0.5)

        # Add ALiBi bias
        alibi_bias = -torch.abs(torch.arange(seq_len, device=x.device).unsqueeze(0) - torch.ara
        alibi_bias = self.m * alibi_bias.unsqueeze(0).unsqueeze(0)
        attn_scores = attn_scores + alibi_bias

        attn_weights = self.dropout(attn_scores.softmax(dim=-1))

        out = (attn_weights @ v).transpose(1, 2).reshape(batch_size, seq_len, self.d_model)
        return self.out_proj(out)
```

```python
    def visualize_alibi_bias(self, seq_len=100):
        alibi_bias = -torch.abs(torch.arange(seq_len).unsqueeze(0) - torch.arange(seq_len).unsq

        # Reshape alibi_bias to (1, seq_len, seq_len) for broadcasting
        alibi_bias = alibi_bias.unsqueeze(0)  # Shape: (1, seq_len, seq_len)

        # Expand m to match the shape for broadcasting: (num_heads, 1, 1)
        m_expanded = self.m.squeeze().unsqueeze(1).unsqueeze(2)  # Shape: (num_heads, 1, 1)

        # Now broadcast properly
        alibi_bias = m_expanded * alibi_bias  # Shape: (num_heads, seq_len, seq_len)

        fig, axes = plt.subplots(2, 4, figsize=(20, 10))
        fig.suptitle("ALiBi Bias for Different Attention Heads")

        for i, ax in enumerate(axes.flat):
            if i < self.num_heads:
                sns.heatmap(alibi_bias[i].detach().cpu().numpy(), ax=ax, cmap='coolwarm')
                ax.set_title(f"Head {i + 1}")
                ax.set_xlabel("Query Position")
                ax.set_ylabel("Key Position")
            else:
                ax.axis('off')

        plt.tight_layout()
        plt.show()

        plt.figure(figsize=(12, 6))
        for i in range(self.num_heads):
            plt.plot(alibi_bias[i, seq_len // 2], label=f'Head {i + 1}')  # Plot for the middle
        plt.title("ALiBi Bias for Different Heads (Middle Query)")
        plt.xlabel("Key Position")
        plt.ylabel("Bias Value")
        plt.legend()
        plt.show()


if __name__ == "__main__":
    d_model = 512
    num_heads = 8
    alibi_attn = ALiBiAttention(d_model, num_heads)

    x = torch.randn(32, 100, d_model)
```

```
    output = alibi_attn(x)

    print(f"Input shape: {x.shape}")
    print(f"Output shape: {output.shape}")

    alibi_attn.visualize_alibi_bias()
```

This implementation includes:

- An `ALiBiAttention` module that implements the ALiBi attention mechanism
- Visualization methods to display the ALiBi bias patterns for different attention heads

# 6. Pros of ALiBi

1. Simplicity: Easy to implement and understand
2. Efficiency: No additional parameters or computations for positional encoding
3. Extrapolation: Can handle longer sequences than seen during training
4. Performance: Often outperforms traditional positional encoding methods
5. Memory efficiency: Doesn't require storing position embeddings
6. Interpretability: The linear bias has a clear interpretation

# 7. Cons of ALiBi

1. Fixed bias: The linear bias is not learned, which might limit flexibility
2. Potential limitations: May not capture complex positional relationships as effectively as learned encodings
3. Task sensitivity: Performance gains may vary depending on the specific task
4. Hyperparameter sensitivity: The slope calculation formula might need adjustment for different model sizes or tasks
5. Limited positional information: Only captures relative positions, not absolute positions

# 8. Comparison with Other Methods

- vs. Absolute Position Encoding: ALiBi doesn't require explicit position embeddings and can generalize to longer sequences
- vs. Relative Position Encoding: ALiBi is simpler and more memory-efficient, but may capture less complex positional relationships

- vs. RoPE: ALiBi is easier to implement and understand, but RoPE might capture more nuanced positional information
- vs. Mixture Position Encoding: ALiBi is simpler and more efficient, but mixture encoding offers more flexibility in combining absolute and relative information

The choice between these methods depends on the specific requirements of the task, the model architecture, and the characteristics of the data being processed. ALiBi shines in scenarios where simplicity, efficiency, and the ability to handle variable-length sequences are crucial.