```
self.scaling_factor = scaling_factor
                     self.dim = dim
                     self.max position embeddings = max position embeddings
                                                                                \Theta = \{	heta_i = 10000^{-2i/d}, i \in [0, 1, 2, \dots d/2 - 1]\}
                     self.base = base
                     inv_freq = 1.0 / (self.base ** (torch.arange(0, self.dim, 2, dtype=torch.int64).float().to(device) / self.dim))
                     # inv freq: torch.Size([dim/2])
                     self.register buffer("inv freq", inv freq, persistent=False)
                     self.max seq len cached = max position embeddings
                     t = torch.arange(self.max_seq_len_cached, device=device, dtype=torch.int64).type_as(self.inv_freq)
                     # t: torch.Size([max position embeddings])
                     # t[pos][j]: pos\theta_j
                     t = t / self.scaling factor
                     freqs = torch.outer(t, self.inv freq) # torch.Size([max position embeddings, dim/2])
                     emb = torch.cat((freqs, freqs), dim=-1) # torch.Size([max_position_embeddings, dim])
                     self.register_buffer("_cos_cached", emb.cos().to(torch.get_default_dtype()), persistent=False)
                     self.register_buffer("_sin_cached", emb.sin().to(torch.get default dtype()), persistent=False)
                        @torch.no grad()
LlamaRotaryEmbedding
                        def forward(self, x, position ids):
                            # x: [bs, num_attention_heads, seq_len, head_size]
                            inv_freq_expanded = self.inv_freq[None, :, None].float().expand(position_ids.shape[0], -1, 1)
                           position_ids_expanded = position_ids[:, None, :].float()
                            device type = x.device.type
                           device_type = device_type if isinstance(device_type, str) and device_type != "mps" else "cpu"
                           with torch.autocast(device type=device type, enabled=False):
                               # get m theta
                                freqs = (inv freq expanded.float() @ position ids expanded.float()).transpose(1, 2)
                                emb = torch.cat((freqs, freqs), dim=-1) # torch.Size([1, len, dim])
                                cos = emb.cos()
                                sin = emb.sin()
                           return cos.to(dtype=x.dtype), sin.to(dtype=x.dtype)
```

```
def rotate_half(x):
    """Rotates half the hidden dims of the input."""
    x1 = x[..., : x.shape[-1] // 2]
    x2 = x[..., x.shape[-1] // 2 :]
    return torch.cat((-x2, x1), dim=-1)
```

```
def apply_rotary_pos_emb(q, k, cos, sin, position_ids=None, unsqueeze_dim=1):
    """Applies Rotary Position Embedding to the query and key tensors."""
    cos = cos.unsqueeze(unsqueeze_dim)
    sin = sin.unsqueeze(unsqueeze_dim)
    q_embed = (q * cos) + (rotate_half(q) * sin)
    k_embed = (k * cos) + (rotate_half(k) * sin)
    return q_embed, k_embed
```