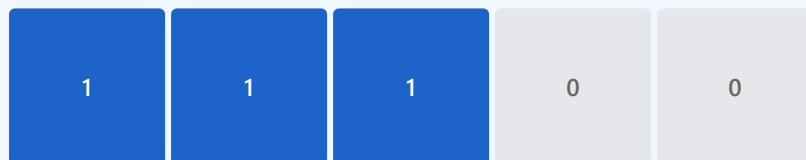# Masking Implementation

## ◆ Padding Mask

**Ignore padding tokens** in attention

◺ Shape: **(batch_size, 1, 1, seq_len)**

🎯 Marks which positions are actual tokens vs padding

| 1 | 1 | 1 | 0 | 0 |

## ⏱ Application Timing

Applied **before softmax** in attention calculation
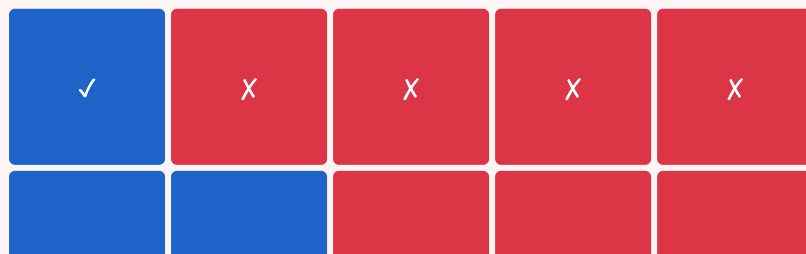
## 🔢 Mask Values

Use **-inf** for masked positions (becomes 0 after softmax)

## ▲ Look-Ahead Mask (Causal)

**Prevent attending to future tokens**

◺ **Upper triangular matrix** of -inf values

🔒 Ensures causality during decoding

| ✓ | X | X | X | X |

### 🔗 Combined Masks

Multiple masks combined using
**element-wise OR operation**

⚠️

**Proper masking is crucial for correct model behavior**