# RNN Attention vs. Self-Attention

| Aspect | RNN Attention | Self-Attention |
|---|---|---|
| 📍 **Source** | Query from **decoder** Keys/Values **encoder** from | All from the **same sequence** |
| ⚙️ **Processing** | **Sequential** processing | **Parallel** processing |
| 🔗 **Relationships** | Limited to current step | Captures relationships between **all positions** |
| ⏱️ **Complexity** | **O(n)** per step | **O(n²)** total |

⚡ **Key Advantage**

Self-attention enables parallelization and better long-range dependency modeling