# Class Notes: Attention Mechanisms in Neural Networks

## 1. Bahdanau (Additive) Attention

For each decoder hidden state $s_t$, the attention mechanism computes a context vector $c_t$ as a weighted sum of encoder hidden states $h_s$.
Keys: Encoder hidden states $h_s$
Queries: Decoder hidden state $s_t$
Values: Encoder hidden states $h_s$

$$e_{ts} = v_a^\top \tanh(W_q s_t + W_k h_s)$$

$$\alpha_{ts} = \frac{\exp(e_{ts})}{\sum_j \exp(e_{tj})}$$
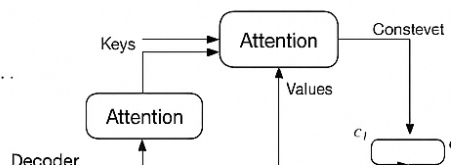


Figure 1: Bahdanau Attention flow diagram.

## 2. Luong (Multiplicative) Attention

The "global" variant uses the decoder hidden state from the previous timestep.
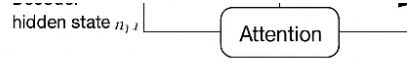Keys: Encoder hidden states $h_s$
Queries: Decoder hidden state $s_{t-1}$
Values: Encoder hidden states $h_s$

$$e_{ts} = h_t^\top h_s \quad \text{or} \quad e_{ts} = h_t^\top W h_s$$

$$\alpha_{ts} = \frac{\exp(e_{ts})}{\sum_j \exp(e_{tj})}$$

Attention weights $\alpha_{ts} = \dfrac{\dot{a_t}^{\cdot\;\cdots}}{\sum_j \exp(e_{tj})}$

Decoder
hidden state $n_{t,t}$



### 2. Luong (Multiplicative) Attention

The "global" variant of Luong attention uses the deoder hidden state $u_t$ previous timestep.

| Keys | Queries | Values |
|------|---------|--------|
| $c_t$ | query | Encoder hidden states $c$ |
| values | values | Encoder hidden states $h, \ldots$ |

Scoring function $e_{ts} = \hat{h_t}\, s_{t-1}$

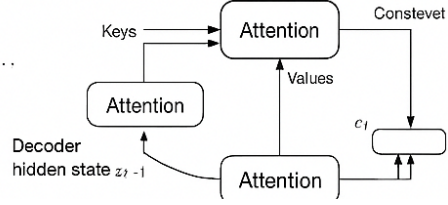Attention weights $\alpha_{ts} = \dfrac{\exp(e_{t_{oo}})}{\sum_s \exp(e_{t_s})}$



Figure 2: Luong Attention flow diagram.

# 3. Multi-Head Attention

Self-attention where $Q, K, V$ are projected into $h$ subspaces, processed in parallel, then concatenated.

$$Q_h = QW_h^Q, \quad K_h = KW_h^K, \quad V_h = VW_h^V$$

$$A_h = \text{softmax}\left(\frac{Q_h K_h^\top}{\sqrt{d_k}}\right), \quad O_h = A_h V_h$$

### 4. Multi-Head Attention

Self-attention. Aggregate sequence into single vector –using self-attention pooling.

Ket the inputs be $(x_t, \ldots x_t)$

$K_{\cdot,b} \equiv x_1\, W_{t_t} \qquad Q_{\eta-} = \dfrac{x_1 \quad W_{H}}{3_t\hat{h}_t}$

$Q_{\eta} = x_1\, W_H \qquad V_H = x_1\, W_H\,^x$

Output of head $h$ $\;\hat{h} = x_l$

$sc = \dfrac{c_{\eta_{1t}}}{g_H} \qquad \alpha_s = \dfrac{\exp(q s_4)}{\sum_{\eta} \exp \hat{h}_j}$



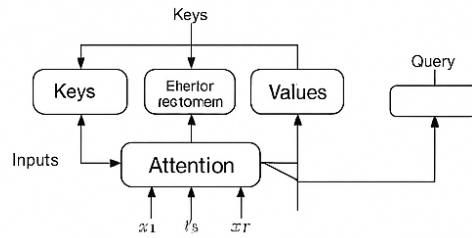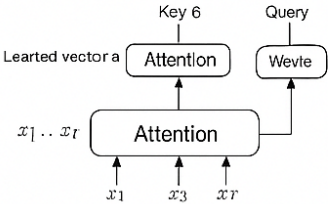Figure 3: Multi-Head Attention computation from the provided diagram.

# 4. Self-Attention

Aggregates a sequence into a single vector via scaled dot-product.

## 4. Self-Attention

Aggregate a sequence $j$ into a single vector, self-attention pooling.

| Keys | queries | Score | Values |
|------|---------|-------|--------|
| $x_1$ | $q$ | Additive | Linear $(c) + xr$ |
| Values | learned | Dot product | Scaled dot-prod |



| Feature | Keys | Score | Output | Score | Twist |
|---------|------|-------|--------|-------|-------|
| Bahdanau $(h_1 \ldots x_T)$ | $x_1, ., ., w_1 \, x_T)$ | Additive | Weighte | Weighted sum | Weighted sum |
| Luong $(x_1, \ldots x_T)$ | $x_1, \ldots, x_1, x_T$ | Dot-pred-act | Multi-head | Weighted sum | Multi-head |

Figure 4: Self-Attention flow from the provided diagram.

# 5. Summary Table

| Variant | Q Source | K Source | Scoring | Use Case |
|---------|----------|----------|---------|----------|
| Bahdanau | Decoder RNN | Encoder RNN | Additive MLP | Seq2Seq |
| Luong | Decoder RNN | Encoder RNN | Dot/General | Seq2Seq |
| Self-Attn | Same sequence | Same sequence | Scaled Dot | Transformers |
| Multi-Head | Same sequence | Same sequence | Scaled Dot (multi) | Rich relations |