**Duration: 15mins**                                                                                   **Max Marks: 10**

1. Define                                                                                                                  [3]
   - **Zero-shot learning:** The model is given zero example to refer/learn from before making predictions on unseen samples. No training happens.

   - **One-shot learning:** The model is given one example to refer/learn from before making predictions on unseen samples.

   - **Few-shot learning:** The model is given more than one example to refer/learn from before making predictions on unseen samples.

2. Why do we use an Add layer in Transformer?                                                                [1]

   The add layer in Transformer acts as a residual/skip connection, which allows gradients to propagate back to the lower layer through a bypass path.

3. Write steps (incl equations, if possible) of Bahadanu's and Vaswani's Attention mechanism. Do NOT draw the encoder-decoder architecture.                                                              [6]

| Bahadanu's (Standard attention ) | Vaswani's (Self-attention) |
|---|---|
| 1. Compute hidden representations $h_j$ <br> 2. Compute attention model <br><br> $$e_{ij} = a(s_{i-1}, h_j)$$ <br><br> 3. Compute softmax (attention weights) <br><br> $$a_{ij} = \frac{\exp(e_{ij})}{\sum \exp(e_{ik})}$$ <br><br> 4. Compute context vector <br><br> $$c_i = \sum a_{ij} h_j$$ | 1. Compute query ($Q$), key ($K$) and value ($V$) through projection. <br><br> $$Q = W_Q \cdot X$$ $$K = W_K \cdot X$$ $$V = W_V \cdot X$$ <br> 2. Compute dot product between a query vector and all other key vectors. <br><br> $$\forall_j \, s_{ij} = q_i \, k_j$$ <br> 3. Scale by sqrt($d_k$) <br><br> $$s'_{ij} = s_{ij} / \text{sqrt}(d_k)$$ <br><br> 4. Compute softmax (attention weights) <br><br> $$a_{ij} = Softmax(s'_{ij})$$ <br><br> 5. Compute attended value vectors <br><br> $$v'_{ij} = a_{ij} \cdot v_{ij}$$ <br><br> 6. Sum all value vectors <br><br> $$\sum_j v'_{ij}$$ |