# Fundamentals of Machine Learning
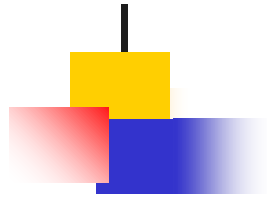
## NETWORK ARCHITECTURES, ATTENTION
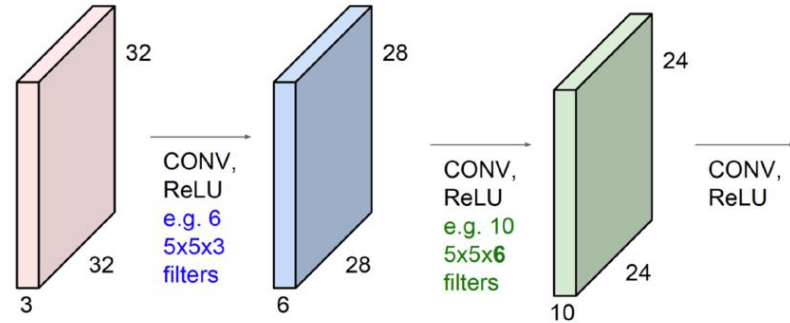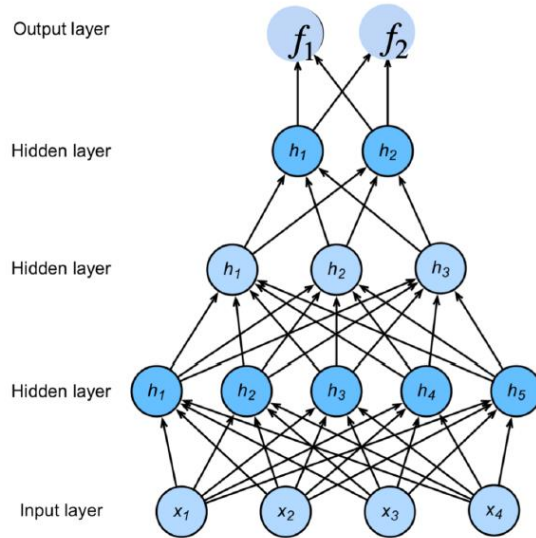
Amit K Roy-Chowdhury

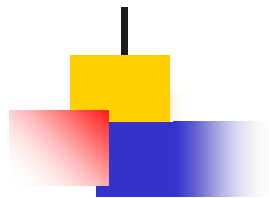# Feedforward Networks



MLP

CNN

# Recurrent Neural Networks



Recurrent network

input layer

hidden layers

output layer

$$h_t = f_W(h_{t-1}, x_t)$$

new state      old state    input at some time step

y

RNN

x

# Autoencoders

# Transformers

# Key, Query, Value

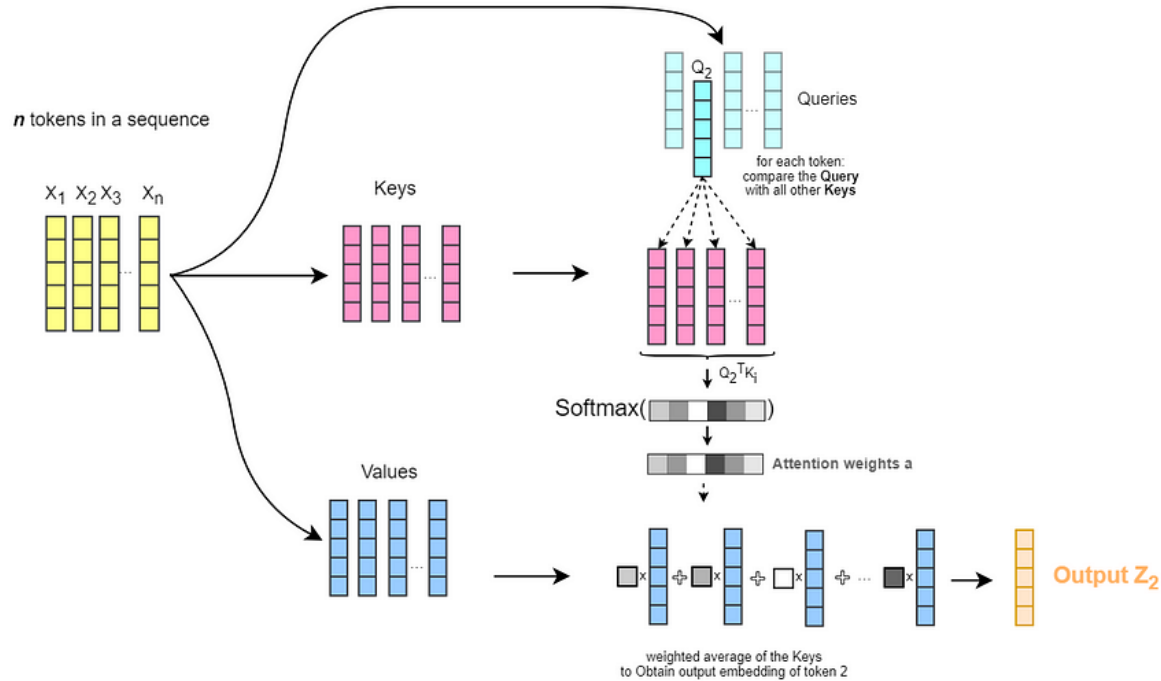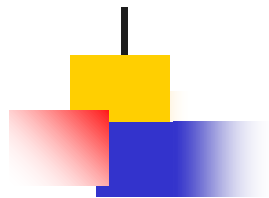Value: $\mathbf{V} \in \mathbb{R}^{m \times v}$   set of $m$ feature vectors
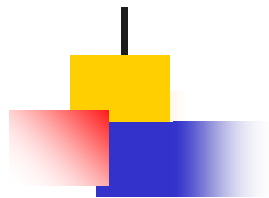
Query: $q \in \mathbb{R}^q$

Keys: $\mathbf{K} \in \mathbb{R}^{m \times k}$   set of $m$ **keys**

Find the query that is most similar to a key and use the corresponding value.

$$\text{Attn}(q, (k_1, v_1), \ldots, (k_m, v_m)) = \text{Attn}(q, (k_{1:m}, v_{1:m})) = \sum_{i=1}^{m} \alpha_i(q, k_{1:m}) v_i$$

$$0 \leq \alpha_i(q, k_{1:m}) \leq 1 \qquad \sum_i \alpha_i(q, k_{1:m}) = 1$$
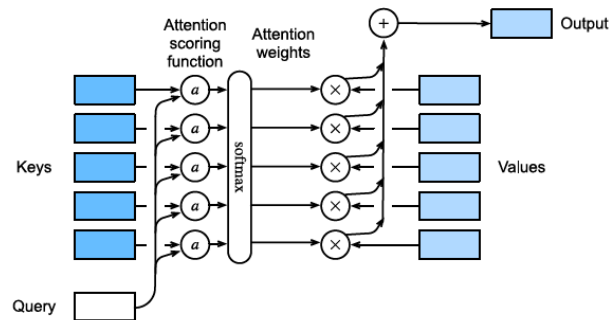
# Computing Attention Weights

Find the query that is most similar to a key and use the corresponding value.
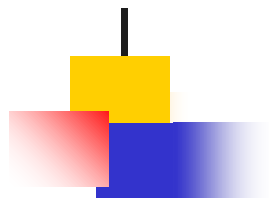
$$\text{Attn}(q, (k_1, v_1), \ldots, (k_m, v_m)) = \text{Attn}(q, (k_{1:m}, v_{1:m})) = \sum_{i=1}^{m} \alpha_i(q, k_{1:m}) v_i$$

$$0 \leq \alpha_i(q, k_{1:m}) \leq 1 \qquad \sum_i \alpha_i(q, k_{1:m}) = 1$$

Attention similarity/score: $a(q, k_i) \in \mathbb{R}$

$$\alpha_i(q, k_{1:m}) = \text{softmax}_i([a(q, k_1), \ldots, a(q, k_m)]) = \frac{\exp(a(q, k_i))}{\sum_{j=1}^{m} \exp(a(q, k_j))}$$
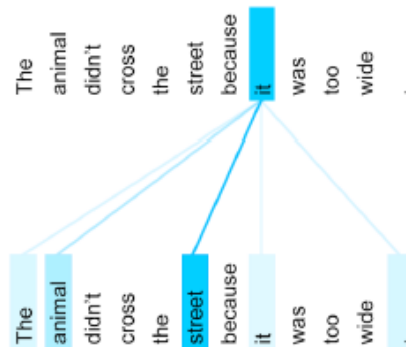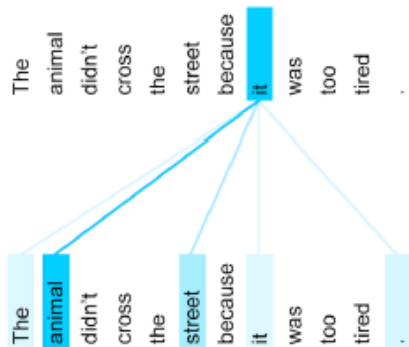
# Self Attention

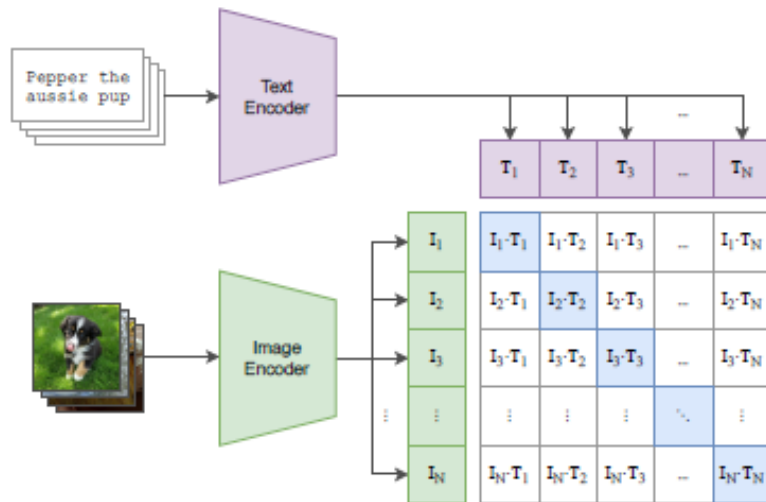$$y_i = \text{Attn}(x_i, (x_1, x_1), \ldots, (x_n, x_n))$$

input tokens $x_1, \ldots, x_n$, where $x_i \in \mathbb{R}^d$

query is $x_i$, and the keys and values are all the (valid) inputs $x_1, \ldots, x_n$

# Multimodal Models



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction