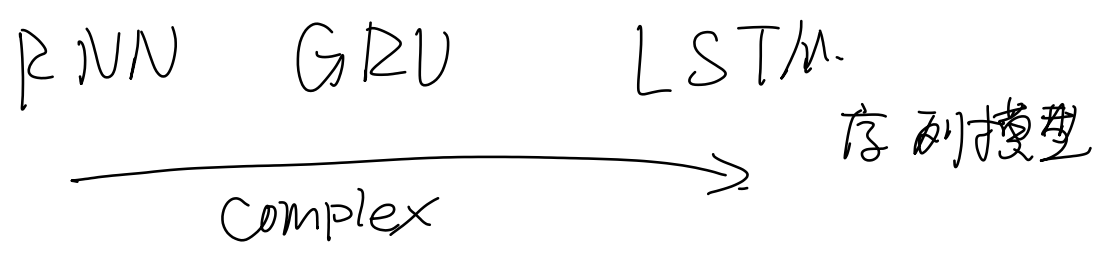


Transformer Network Intuition



Self - Attention

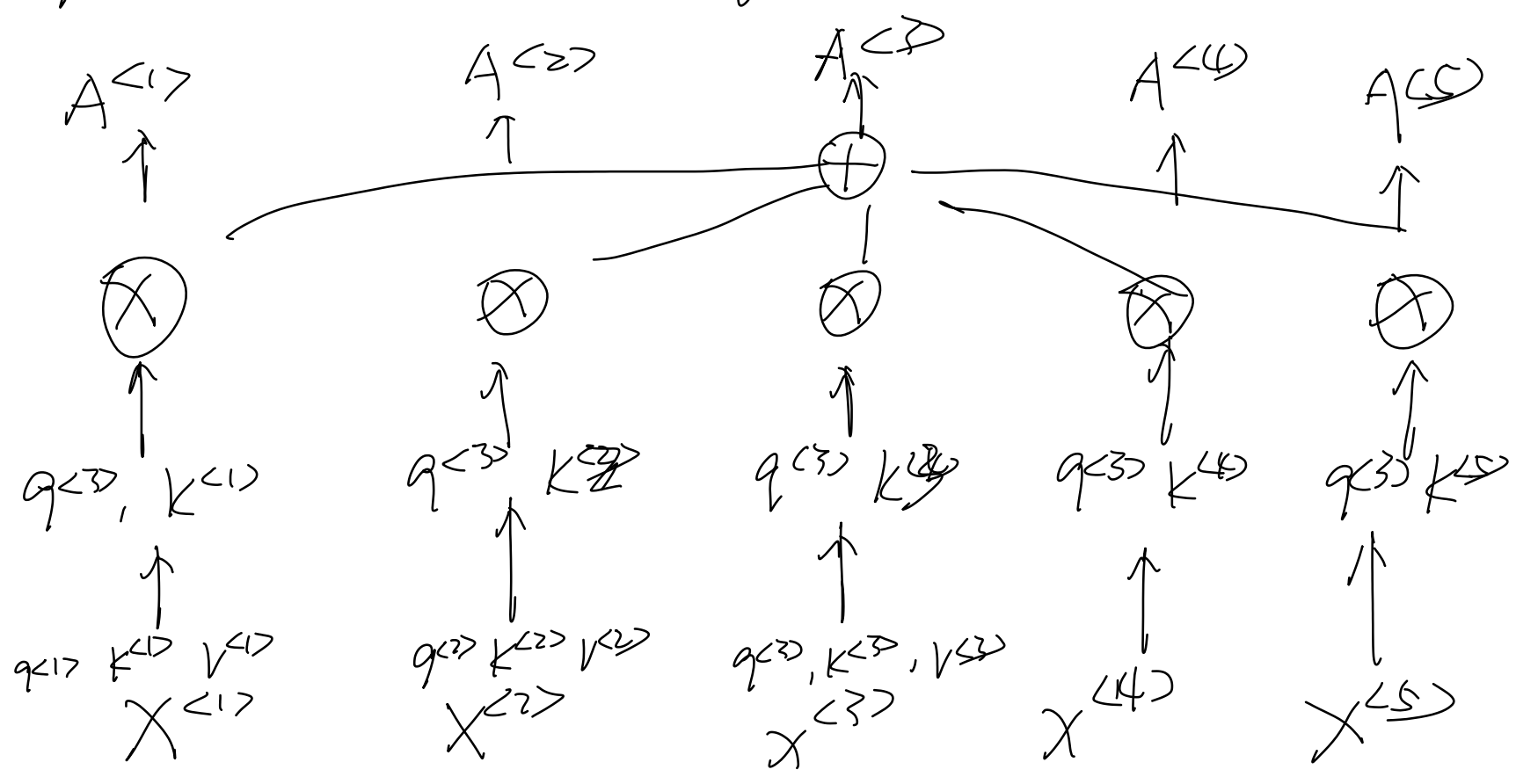
$A(q, k, V)$ = attention-based vector representation of a word.

RNN:

$$a^{<t, t'>} = \frac{\exp(e^{<t, t'>})}{\sum_{t'=1}^T \exp(e^{<t, t'>})}$$
$$A(q, k, V) = \sum_i \frac{\exp(q \cdot k^{<i>})}{\sum_j \exp(q \cdot k^{<j>})} V^{<i>}$$

Self - Attention

Query (Q)	Key (K)	Value (V)
$q^{<1>}$	$k^{<1>}$	$v^{<1>}$
$q^{<2>}$	$k^{<2>}$	$v^{<2>}$
$q^{<3>}$	$k^{<3>}$	$v^{<3>}$
$q^{<4>}$	$k^{<4>}$	$v^{<4>}$
$q^{<5>}$	$k^{<5>}$	$v^{<5>}$



Multi - Head Attention & Transformer Network