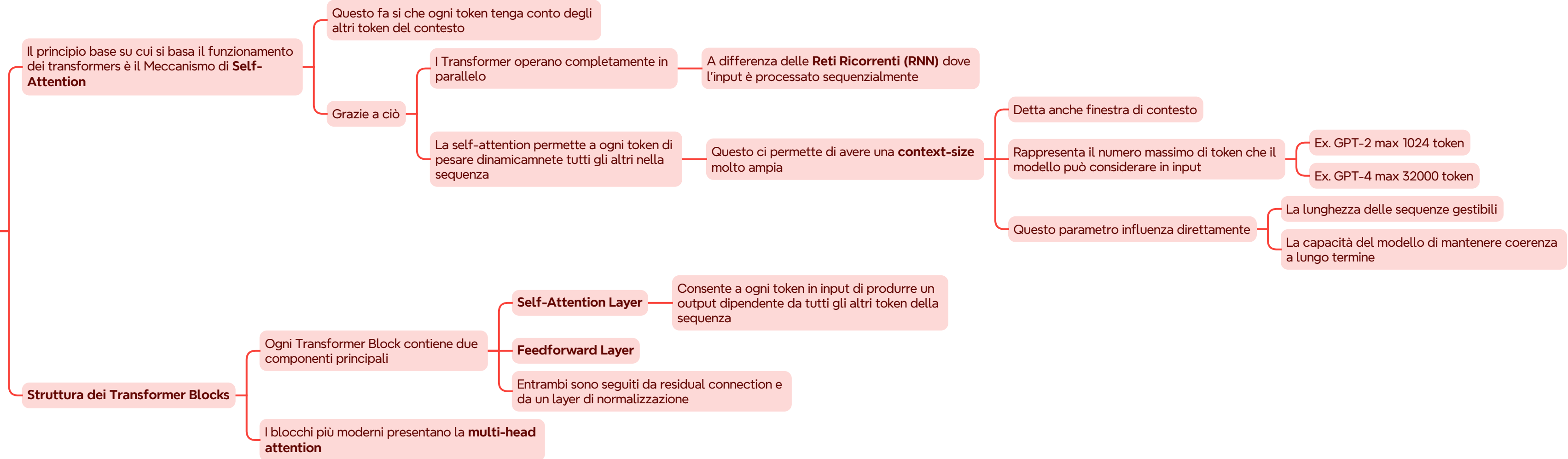
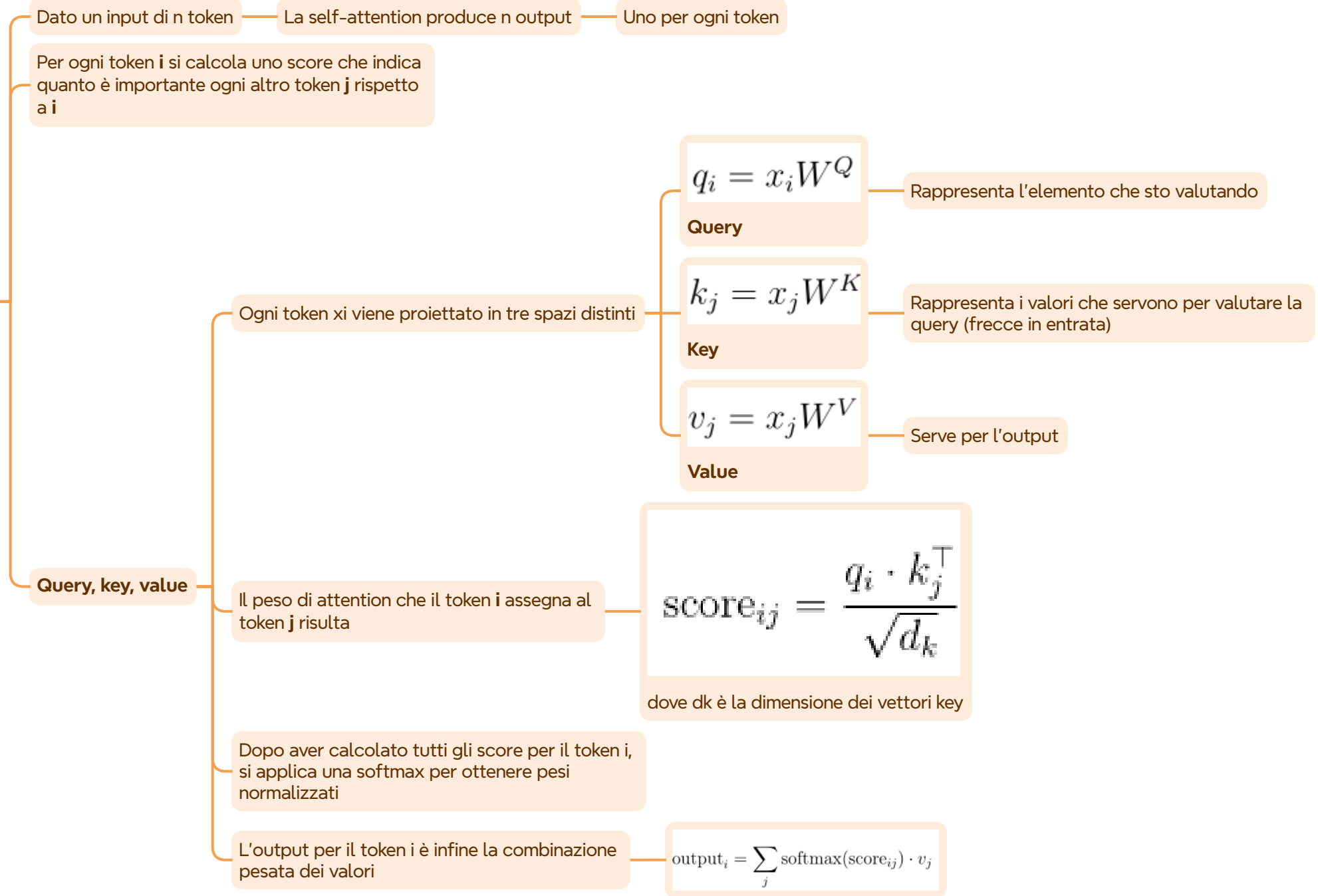


Transformers

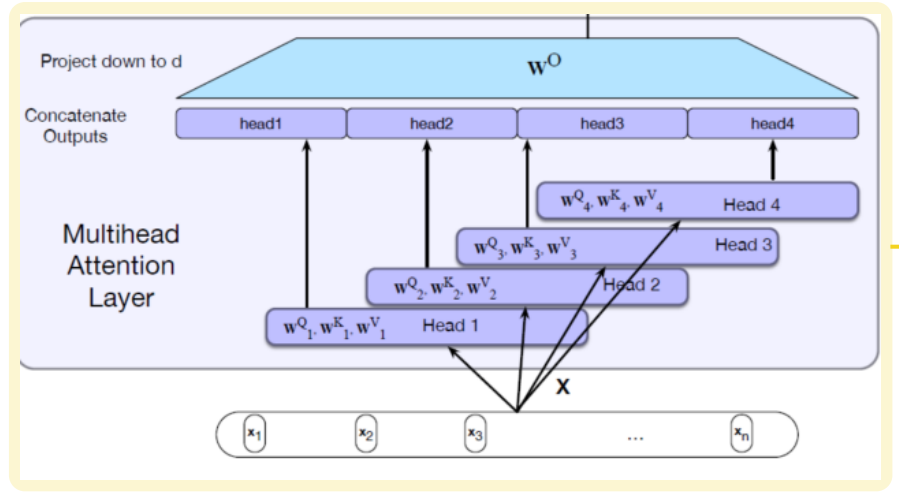
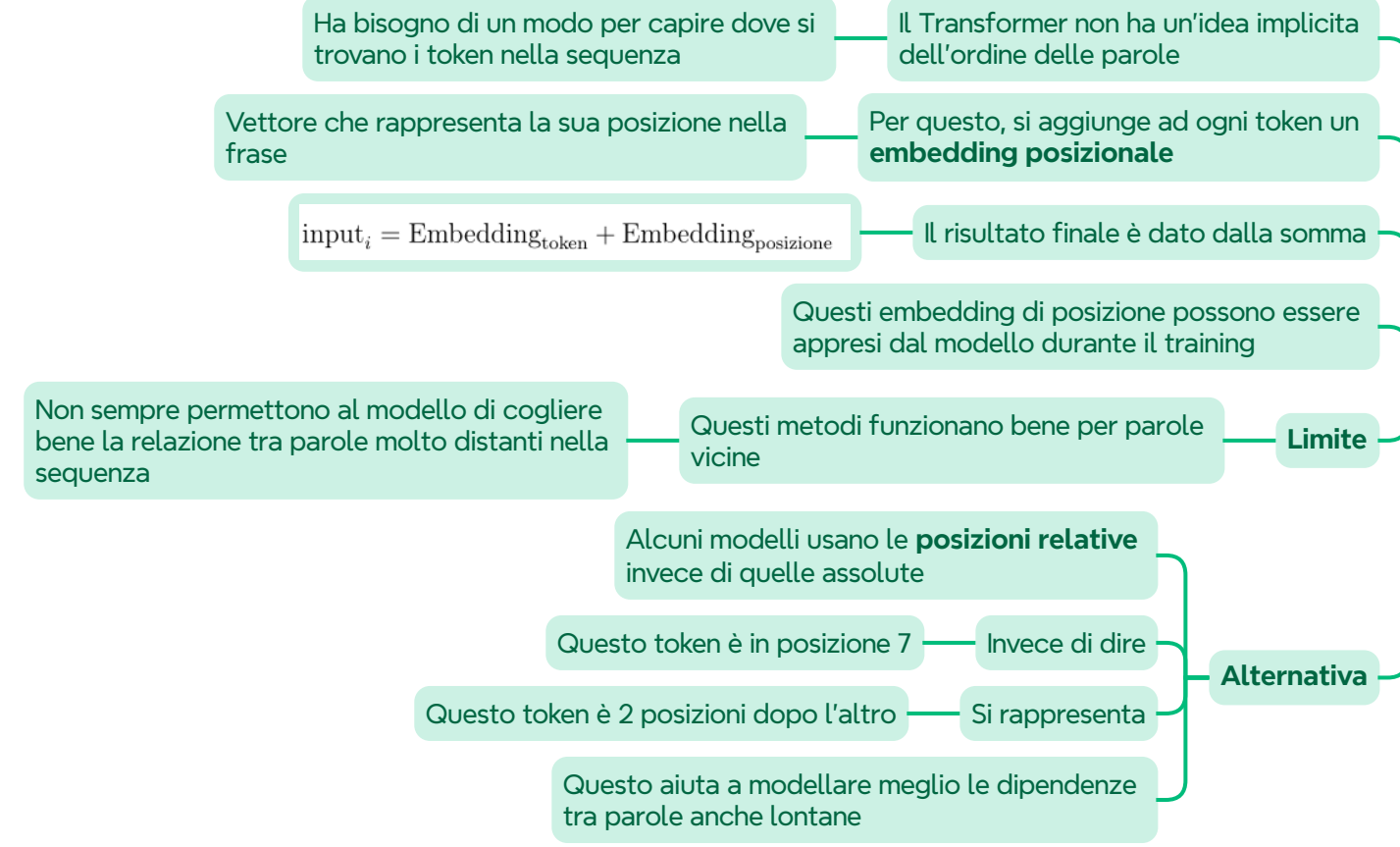
1. Introduzione



2. Funzionamento del blocco di self-attention



4. Codifica Posizionale dei Token



Consente al modello di apprendere più rappresentazioni dell'input da differenti spazi di attenzione

La Multi-Head Attention è un'estensione del meccanismo di Self-Attention

Ogni testa di attenzione (head) lavora con una proiezione diversa dell'input originale

3. Multi-head Attention Layer

