

## **Riassunto: "Attention Is All You Need"**

Il paper introduce il modello Transformer, un'architettura di rete neurale basata solo su meccanismi di attenzione, senza ricorrenze o convoluzioni. Questo modello si collega all'encoder e al decoder attraverso un meccanismo di attenzione, ottenendo risultati superiori in termini di qualità, parallelizzazione e tempo di addestramento rispetto ai modelli tradizionali. Sperimentando su due compiti di traduzione automatica, il Transformer raggiunge risultati migliori rispetto ai modelli esistenti, inclusi ensemble. Inoltre, dimostra di generalizzare bene ad altri compiti, come il parsing di frase in lingua inglese.

## **Riassunto: "Transformer: A Model Architecture for Language Understanding"**

Il Transformer è un'architettura di modello che usa esclusivamente meccanismi di attenzione per modellare dipendenze globali tra input e output, evitando l'uso di reti ricorrenti. Questo permette una maggiore parallelizzazione e ha dimostrato di raggiungere risultati di qualità superiore nella traduzione con poche ore di addestramento. L'architettura del Transformer consiste in encoder e decoder composti da stack di livelli con meccanismi di auto-attenzione multi-testa e reti feed-forward completamente connesse. Il modello utilizza connessioni residue e normalizzazione dei livelli per facilitare l'apprendimento delle dipendenze.

## **Riassunto di "Decoder"**

Il decoder è composto da 6 livelli identici e utilizza un'attenzione multi-head sullo stack dell'encoder. Ogni livello del decoder ha tre sub-livelli: due di auto-attenzione e uno di attenzione multi-head sull'output dello stack dell'encoder. Si utilizzano connessioni residue e normalizzazione dei livelli. Viene modificato il sub-livello di auto-attenzione per impedire a posizioni successive di attendere posizioni successive. Viene introdotta una mascheratura per assicurare che le previsioni dipendano solo dagli output conosciuti delle posizioni precedenti. L'attenzione è descritta come una funzione che mappa una query e un insieme di coppie chiave-valore in un output. Viene introdotta l'attenzione Scaled Dot-Product, in cui i prodotti scalari delle query con le chiavi vengono divisi per la radice quadrata della dimensione delle chiavi. Si utilizzano diverse proiezioni lineari per le query, le chiavi e i valori, seguite dall'attenzione multi-head.

## **Riassunto: "Multi-Head Attention in Transformer Model"**

Il modello Transformer utilizza l'attenzione multi-head per processare le informazioni in modi diversi: - Nelle "encoder-decoder attention" le query provengono dal layer precedente del decoder, mentre le chiavi e i valori di memoria provengono dall'output dell'encoder. Questo permette a ogni posizione nel decoder di considerare tutte le posizioni nella sequenza di input. - L'encoder include strati di self-attention. Qui chiavi, valori e query provengono dallo stesso posto, cioè dall'output del layer precedente nell'encoder. Ogni posizione nell'encoder può considerare tutte le posizioni nel layer precedente dell'encoder.

# Riassunto del libro "Self-Attention in Transformer Models"

Il libro descrive come i livelli di self-attention nel decoder permettano a ogni posizione di "guardare" tutte le posizioni precedenti per mantenere la proprietà auto-regressiva. Viene spiegato come venga impedito il flusso di informazioni verso sinistra all'interno del decoder. Vengono utilizzati livelli di feed-forward e embedding per convertire i token in vettori e calcolare le probabilità del token successivo. Viene introdotto anche il concetto di "positional encoding" per dare informazioni sulla posizione relativa dei token nella sequenza. Infine, vengono forniti dettagli sulle complessità e le operazioni sequenziali per diversi tipi di layer nei modelli di trasformatori.

## Riassunto: Positional Encoding and Self-Attention

- **Positional Encoding:**
  - Usiamo funzioni seno e coseno di diverse frequenze per codificare la posizione di ciascuna dimensione.
  - Le lunghezze d'onda formano una progressione geometrica da  $2\pi$  a  $10000 \cdot 2\pi$ .
  - Questa codifica aiuta il modello a imparare facilmente a prestare attenzione alle posizioni relative.
- **Self-Attention:**
  - Confrontiamo i layer di self-attention con i layer ricorrenti e convoluzionali.
  - Consideriamo la complessità computazionale, la parallelizzazione e la lunghezza del percorso tra dipendenze a lungo raggio.
  - I percorsi più brevi tra qualsiasi combinazione di posizioni nelle sequenze di input e output facilitano l'apprendimento delle dipendenze a lungo raggio.

## Riassunto del libro "Self-Attention Layers vs Recurrent Layers"

- **Self-Attention vs Recurrent Layers:**
  - Le self-attention layers sono più veloci delle recurrent layers quando la lunghezza della sequenza è minore della dimensionalità della rappresentazione.
  - Limitare self-attention a un vicinato di dimensione  $r$  può migliorare le prestazioni computazionali per sequenze molto lunghe.
- **Convoluzione vs Self-Attention:**
  - Le convoluzioni richiedono un numero maggiore di operazioni rispetto alle self-attention layers.
  - Le convoluzioni separate riducono significativamente la complessità computazionale.
- **Addestramento:**

- I modelli sono stati addestrati su dataset standard come WMT 2014 English-German e English-French.
- Le coppie di frasi sono state raggruppate in batch in base alla lunghezza approssimativa della sequenza.
- I modelli sono stati addestrati su una macchina con 8 GPU NVIDIA P100, con tempi di addestramento di circa 0.4 secondi per i modelli base e 1 secondo per i modelli più grandi.

## Riassunto: "Optimizer e Regularizzazione nel Modello Transformer"

- **Optimizer:** Nel modello Transformer è stato utilizzato l'ottimizzatore Adam con specifici valori per i parametri  $\beta_1$ ,  $\beta_2$  e  $\epsilon$ . La variazione del tasso di apprendimento durante l'addestramento è stata regolata da una formula che aumenta linearmente il tasso di apprendimento per i primi passaggi di addestramento e successivamente lo riduce proporzionalmente all'inverso della radice quadrata del numero di passaggi.
- **Regularizzazione:** Durante l'addestramento del modello Transformer sono state impiegate tre tipologie di regularizzazione per migliorare le prestazioni. Il modello Transformer ha ottenuto punteggi BLEU migliori rispetto ai modelli precedenti nei test da inglese a tedesco e da inglese a francese, con un costo di addestramento inferiore.

## Riassunto del libro "Transformer Neural Network Models"

Il libro "Transformer Neural Network Models" tratta dell'utilizzo di diversi modelli neurali per il machine translation. Vengono confrontati vari modelli come Deep-Att + PosUnk Ensemble, GNMT + RL Ensemble, ConvS2S Ensemble, Transformer (base model) e Transformer (big). Viene anche descritto l'utilizzo di tecniche come Residual Dropout e Label Smoothing per migliorare le prestazioni dei modelli. Residual Dropout consiste nell'applicare dropout alle output di ogni sub-layer prima di essere sommate agli input e normalizzate. Label Smoothing, invece, viene utilizzato durante il training per migliorare l'accuratezza e il punteggio BLEU, nonostante possa aumentare la perplexity.

## Riassunto: "Machine Translation with Transformer Models"

Il modello Transformer (big) ha ottenuto ottimi risultati nella traduzione da inglese a tedesco e da inglese a francese nel compito WMT 2014. Ha superato i modelli precedenti di oltre 2.0 BLEU, stabilendo un nuovo punteggio di 28.4 BLEU per l'inglese-tedesco e 41.0 BLEU per l'inglese-francese. Il modello è stato allenato su 8 GPU P100 per 3.5 giorni. Anche il modello base ha superato tutti i modelli pubblicati in precedenza, a un costo di allenamento inferiore. Sono state sperimentate diverse varianti del modello base per valutarne l'importanza nei risultati di traduzione.

## Riassunto: Il modello Transformer

Il modello Transformer è basato sull'attenzione e non utilizza più i tradizionali strati ricorrenti nei modelli encoder-decoder. Durante gli esperimenti, si è osservato che variare il numero di teste di attenzione e le dimensioni chiave e valore di attenzione può influenzare la qualità del modello. Ridurre le dimensioni della chiave di attenzione può peggiorare la qualità del modello. Modelli più grandi tendono ad essere migliori, e l'utilizzo del dropout è utile per evitare l'overfitting. Inoltre, sostituire l'encoding posizionale sinusoidale con embedding posizionali appresi non ha

influenzato significativamente i risultati. Il modello Transformer ha ottenuto ottimi risultati nel parsing di costituenza inglese, superando altri modelli, anche in condizioni di dati limitati.

## **Riassunto "Transformer"**

Il Transformer può essere addestrato molto più velocemente di altre architetture come quelle basate su livelli ricorrenti o convoluzionali. Ha ottenuto risultati eccellenti nelle traduzioni dall'inglese al tedesco e francese nel 2014. L'obiettivo futuro è estendere il Transformer ad altri tipi di dati oltre al testo e migliorare la generazione non sequenziale. Il codice utilizzato è disponibile su GitHub.

## **Riassunto delle principali ricerche in Machine Learning e NLP**

Il testo contiene una serie di ricerche nel campo del Machine Learning e del Natural Language Processing, tra cui modelli di traduzione neurale veloci, reti con attenzione strutturata, ottimizzazione stocastica, trucchi di fattorizzazione per reti LSTM, embedding di frasi auto-attentive strutturate, apprendimento multitask con sequenze, approcci efficaci alla traduzione neurale basata sull'attenzione, modelli di sommarizzazione astratti rinforzati, reti neurali end-to-end con memoria, e altro ancora. Le ricerche riguardano anche l'uso di embedding di output per migliorare i modelli linguistici, la traduzione neurale di parole rare con unità di sotto-parole, reti neurali estremamente grandi con layer di esperti a varie porte, il dropout come metodo per prevenire l'overfitting nelle reti neurali, e nuove architetture per la visione artificiale.

## **Riassunto: "Grammar as a foreign language"**

Il libro parla di diversi studi e sistemi di traduzione automatica basati su reti neurali. Viene evidenziata l'importanza delle leggi che regolano il processo di registrazione e voto negli Stati Uniti. Inoltre, viene menzionata l'attenzione dell'encoder su dipendenze a lunga distanza durante il processo di traduzione.

## **Riassunto: *The Law will never be perfect***

Il libro discute l'importanza di applicare la legge in modo giusto, nonostante essa non possa mai essere perfetta. Viene evidenziata la mancanza di una corretta applicazione della legge secondo l'autore. Viene inoltre menzionato un esempio di risoluzione di anafora in un'immagine.

## **Riassunto: "The Law will never be perfect"**

Il libro discute dell'importanza che la legge sia applicata in modo giusto, nonostante non possa essere perfetta. Vengono forniti esempi di come diverse parti del processo legale svolgano compiti diversi per garantire una corretta applicazione della legge.