

2. Text to Image

(א) תאביר יהיו כאלה:

(1) Text Encoder: מקבל טקסט ומוציא וקטור ארוך (אנצ'ורי) ממייצג
הטקסט. הוויזואליזציה של RNN או
Transformers

(2) Image Decoder: מקבל וקטור ממייצג טקסט כמייצג של תמונה ומוציא
וקטור (ע"י רשת CNN או אחרת)

ע"י שימוש ב-Transformers או רשתות RNN, מודל יכול ליצור
אינטקטואל (אנצ'ורי) שמתקדם בהדרגה, ומקבל גישה יכול ליצור
אינטקטואל (אנצ'ורי) יחד עם קול אחרות. בדרך כלל, ישנה גישה
אחרת על מנת ליצור תמונה מלאה.

(ב) גישה נוספת היא שימוש ב-attention מודל. אלו מודלים אחרים
אשר מקבלים באופן ישיר את הוויזואליזציה של טקסט. מודלים אחרים
הם: מודל attention ב-Text Encoder, אלו מודלים אחרים
אשר מקבלים באופן ישיר את הוויזואליזציה של טקסט.

ב. נראה כי attention מודל אחר, אשר מקבל את הוויזואליזציה של טקסט
(הוויזואליזציה של ה-Encoder) ליצור תמונה. הוויזואליזציה של טקסט
אשר מקבלים באופן ישיר את הוויזואליזציה של טקסט.
בדרך כלל, ישנה גישה אחרת ליצור תמונה מלאה.

Transformers networks 4

באשר מבצעים סקור attention, סקור ה tokens לא נשנה אך הניקוד
 לכל token. סקור attention נעשה יחד עם
 כל ה tokens כדי לסקור את היחסים ביניהם אינם על גלגל ולא
 אומרים חלקי אחד וסדר ביניהם.

בשלב הבא נעצירה סקור FC היא נחלת על כל ורכי
 נאמן לפרגמטיות איחוד של וברור לכל חלוקה אספק

אכן הרבה של סקור attention ! FC גייזר
 Permutation invariant operator

אם נביא positional encoding, לעומת זאת הניקוד של מקום
 ה token ביחס ל tokens האחרים וכל המקום האבסולוטי
 של הניקוד, הניקוד,

נקרא להיות ה Permutation invariant
 כי השכבה שנינו לביא לפרגמטיות. אף המקום של
 ה tokens וכל המקום לא קודם.