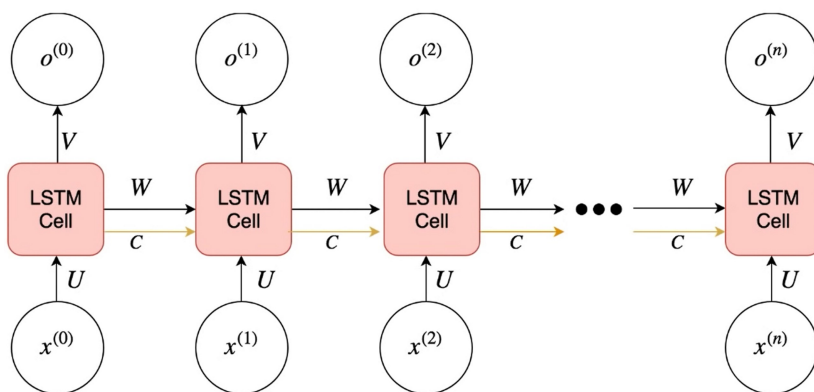
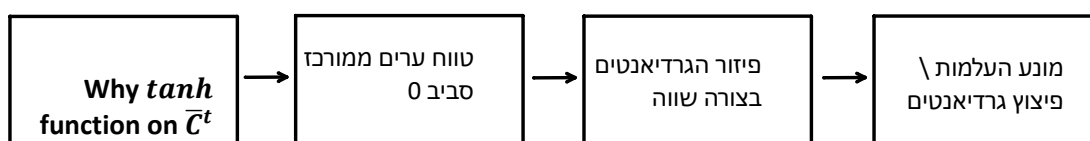


מבט מבחוץ, כאשר אנחנו משתמשים ביחידת LSTM בודדה:

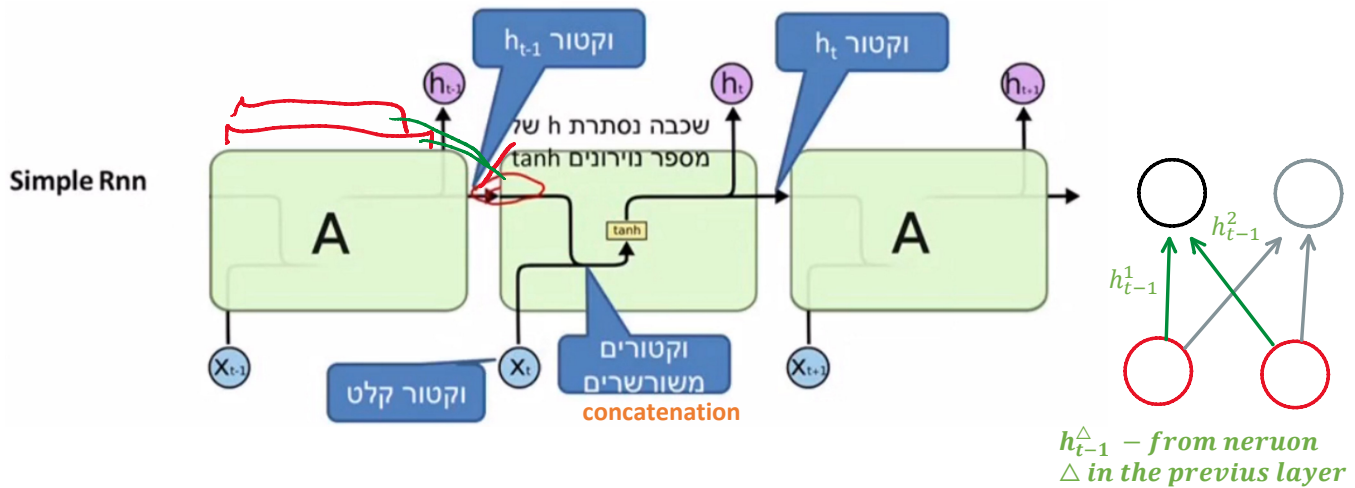


בנוסף לשערים הללו, יש לנו עוד וקטור \tilde{C}^t שתפקידו הוא לשנות את מצב התא C_t פונקציית האקטיבציה שלו היא Tanh, מכיוון שאנחנו יוצרים טווח ערכים ממורכז סביב 0, זה יאפשר לנו "לפזר" את הגרדיאנטים בצורה שווה,

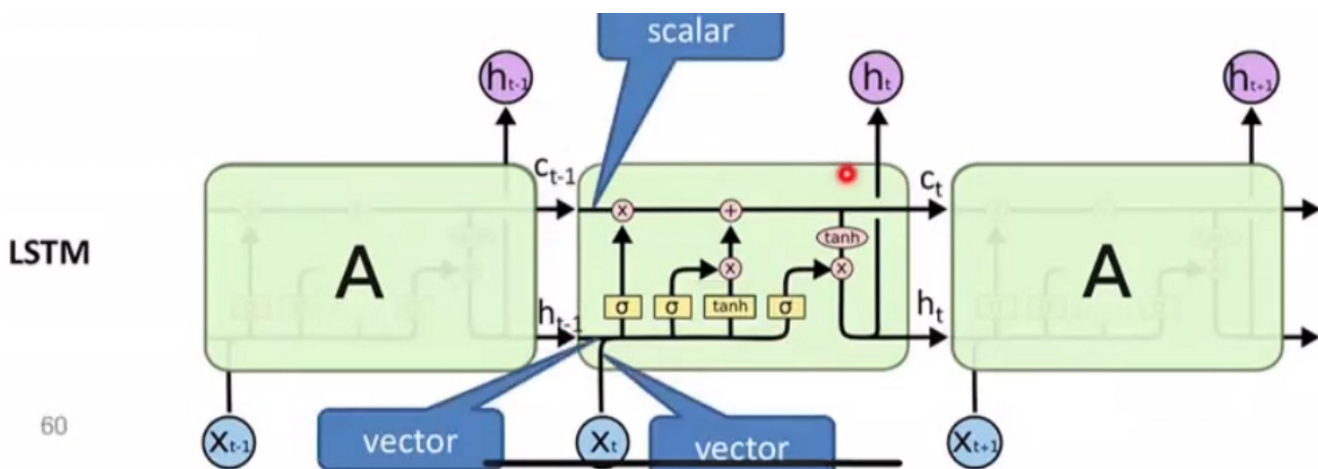


LSTM Input / Output

ברשת RNN רגילה, אנחנו מקבלים את h_{t-1} כווקטור שמגיע מכל הניורונים בשכבה הקודמת ב- $t-1$, כמו ברשת רגילה



ב LSTM אנחנו מקבלים את הווקטור h כמו ברשת RNN רגילה (מכל אחד מהניורונים בשכבה הקודמת) אבל ב-LSTM בנוסף נקבל את מצב הנירון C , שאותו הוא מקבל רק מעצמו בפרק הזמן הקודם, ולכן C סקלר



כניסות ל-LSTM

h_{t-1} ווקטור שמגיע מכל יחידות ה-LSTM שפעלו בזמן $t-1$

C_{t-1} סקלר שמגיע מאותה יחידת LSTM שפעלה בזמן $t-1$

יציאות מ-LSTM

Cell State - C_t : סקלר נושא איתו זיכרון ארוך טווח וקצר טווח (נשלט על ידי מספר שערים)

h_t סקלר $[-1, 1]$ המחושב על פי המצב האחרון של C_t (נשלט בעזרת שער בודד)

רביבי LSTM

יחידת LSTM מכילה בתוכה 4 נוירונים נסתרים, לכל אחד מהם מטריצת משקולות וביאסים
כל אחד מהנוירונים מקבל את שרשרת הוקטורים h_{t-1}, x_t

לדוגמה:

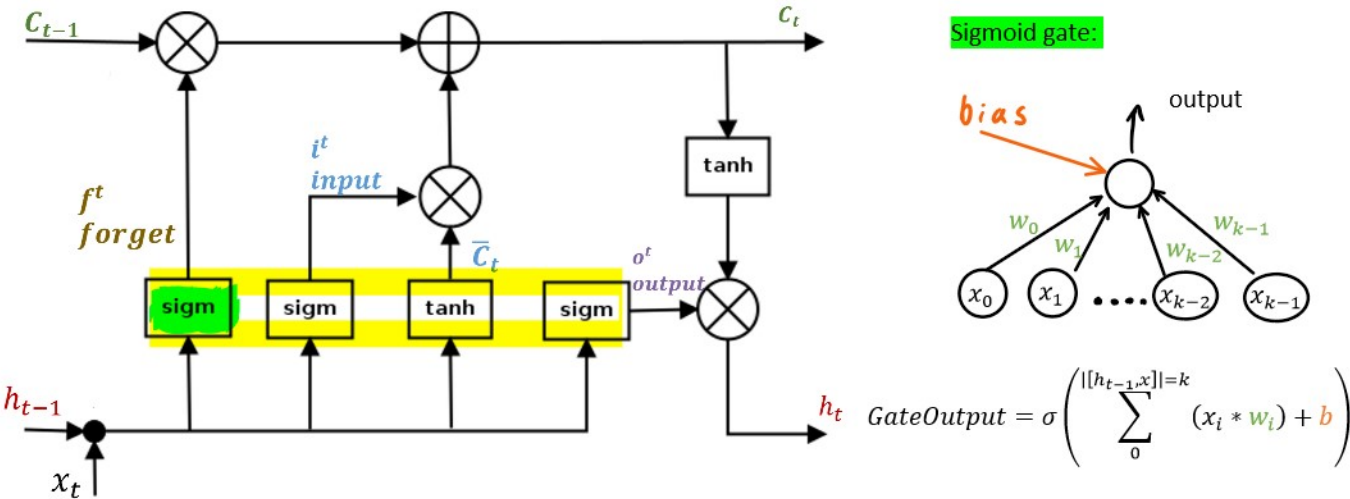
אם יש לנו רשת של 100 יחידות LSTM וקלט x מממד 50

אז מטריצת המשקולות לשער סיגמואיד יחיד, באחת מיחידות הLSTM תהיה (100 + 50) כלומר וקטור באורך 150, ועוד ביאס בודד לעצמו

אם מסתכלים על כל 4 הנוירונים הנסתרים, אז לכל אחד מהם יש אותו גודל משקולות ולכן לכל יחידת LSTM יהיו $(100 + 50 + 1_{bias}) * 4 = 604$ משקולות

אם מסתכלים על כל 100 היחידות שיש ברשת שלנו אז יש $604 * 100 = 60,400$ משקולות לכל הרשת

מספר הנוירונים ברשת הם מספר השערים, סה"כ יהיה מספר יחידות ה LSTM כפול מספר השערים בכל אחת (4) ולכן $100 * 4 = 400$

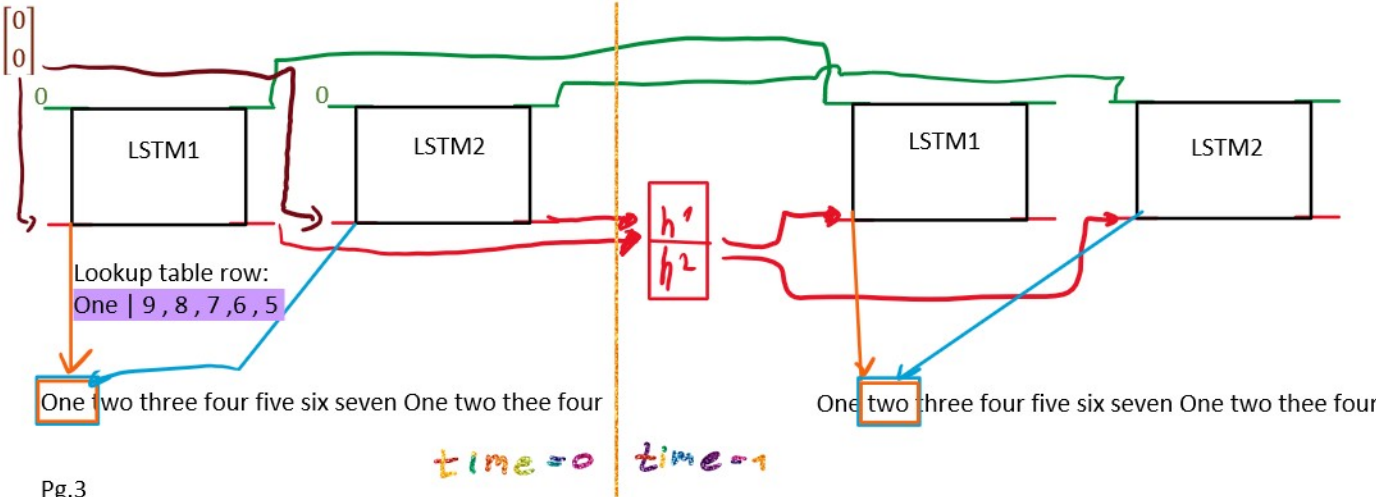


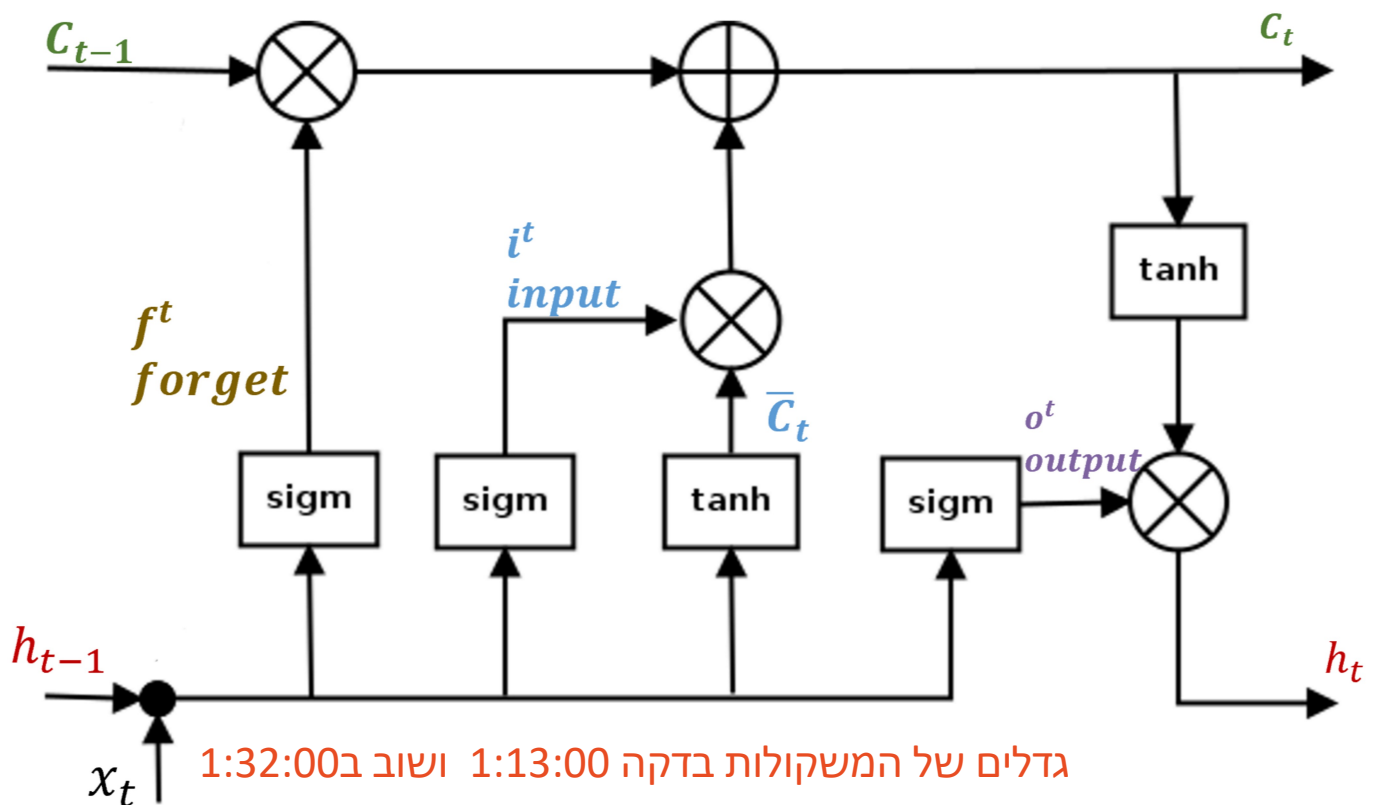
דוגמה מעשית:

אם יש לנו רשת של 2 יחידות LSTM, הקלט של כל אחת מהיחידות הוא ייצוג המילה (מהLookup table)

אז בהנחה שכל מילה מקודדת לווקטור באורך 5 אז אורך קלט ואורך ה h_{t-1} יהיה 2 [כמספר יחידות ה LSTM]

ואז השרשור $[h_{t-1}, x_t]$ לדוגמה בודדת $[h^1, h^2, x^1, x^2, x^3, x^4, x^5]$ ----- לאיטרציה הראשונה $t=0$: $[0, 0, 9, 8, 7, 6, 5]$





$$C_t = f^t \times C_{t-1} + i^t \times \bar{C}_t$$

Input Gate

Controls weather the memory cell is updated based on the new input x_t

$$\bar{C}_t = \tanh(W^c[h_{t-1}, x_t] + b^c)$$

$$i^t = \sigma(W^i[h_{t-1}, x_t] + b^i)$$

forget Gate $f^t \in [0,1]$

Controls weather the memory cell is reset to 0.// controls how much of the old state should be forgotten.
 $f_t = 1$ we preserve all the history
 $f_t = 0$, we forget all history

$$f^t = \sigma(W^f[h_{t-1}, x_t] + b^f)$$

output gate

Weather the information of the current cell state is made visible

$$o_t = \sigma(W^o[h_{t-1}, x_t] + b^o)$$

W^N Weights of N
 b^N Biases of N
 (Those Part Learn through back-prop)

$[h_{t-1}, x_t] \leftrightarrow h_{t-1}$ Concatinated with x_t

$h_{t-1} = [1,2,3], x_t = [4,5,6] \Rightarrow [h_{t-1}, x_t] = [1,2,3,4,5,6]$