

SiMBA: Simplified Mamba-based Architecture for Vision and Multivariate Time series

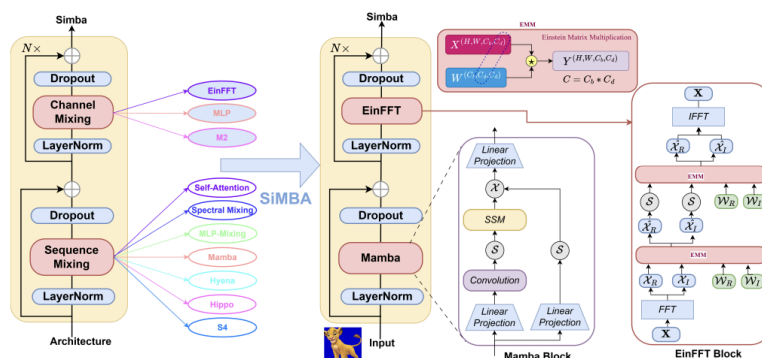


Fig. 1: Simplified Mamba Based Architecture.

המאמר הזה משך את תשומת ליבי כי מצד אחד יש בו שימוש נרחב בהתמרת פורייה ובייצוגים של דאטה בתחום התדר. החולשה שלי לתחום התדר נובעת מכך שביליתי כמה מהשנים הראשונות של הקריירה בתחום עיבוד של אותות אלחוטיים. מצד שני המאמר גם משתמש בארכיטקטורת ממבה שסקרתי בהרחבה בחודשים האחרונים (וכנראה אמשך עם זה כי מאמרים מעניינים בנושא מרתק זה לא מפסיקים להגיע).

אוקיי, אז מה יש לנו במאמר הזה? המאמר מציע שדרוג יפה לארכיטקטורה של ממבה המערב כאמור התמרות פורייה וקצת משחקים בתחום התדר. הארכיטקטורה המוצעת מתאימה גם לדאטה ויזואלי וגם לסדרות זמן multivariate. המאמר כתוב בצורה די מסורבלת והיה לי לא טריוויאלי לגלות מה הם באמת עשו עקב הסברים וסימונים לא ברורים. אבל כאמור הרעיון מאחורי המאמר הוא די חמוד.

המחברים מנסים לשפר את ממבה על ידי הוספת שכבה שבגדול לוקחת את הייצוגים המופקים על ידי ממבה ו"מחזקת" אותם על ידי פלטור תדרים מסוימים מהם (הייצוגים). קודם כל נציין שמפעילים את המנגנון המוצע, שקיבל שם EinFFT, על כל איבר סדרה בנפרד (פאץ' של תמונה) בצורה ממוקבלת. כאמור הסיפור מתחיל מהפעלת התמרת פורייה על הפלט (=ייצוג פיסת דאטה) של שכבת ממבה. ואז המאמר הופך להיות די לא ברור והדבר הזה גזל ממני בערך שעותיים כדי להבין שלא אני מפספס משהו אלא המאמר עצמו קצת לא מדויק (בתקווה עמדתי במשימה זו).

כאמור הרעיון הוא מפלטור תדרים (תלויות) הלא נחוצים (לביצוע המשימה) בייצוגי איברי הסדרה. הפלטור מתבצע במחבר הייצוג של הדאטה (כלומר אמבדינג) ונקרא channel-mixing. כלומר שכבה זו היא משמשת בתור תוספת/החלפה ל-MLP שלפעמים משמשת לאותה המטרה.

אבל איך הוא עושה את זה משתנה בין נוסחה לנוסחה במאמר. במאמר עצמו (נוסחה 4) קודם כל מפעילים שכבה לינארית במישור המרוכב ולאחריה סיגמויד (גם במישור המרוכב). ב-[appendix](#) (בתחילת עמוד 22) זה כבר מופיעה שכבה אחת של ReLU, לאחר מכן עוד שכבה לינארית, לאחר מכן מפעילים פונקציית [softshrink](#) שמטרתה היא לאפס תדרים סביב אפס באינטרבל באורך $\lambda \cdot 2$ ולהזיז כאלו מעבר לזה ב- λ . כלומר איזה stop-band filter מוזר.

הגרסה השלישית מגיע מהדף האחרון של ה-appendix ששם יש רק ReLUs. לא הסתכלתי בקוד אז לא ברור מה באמת קורה שם. כל הפעולות האלו מתבצעות בצורה נפרדת במישור הממשי ובמישור המדומה שלאחר מכן משלבים אותם. בשלב האחרון מבצעים התמרת פורייה הפוכה (IFFT).

אוקיי, אז בואו נחזור לעיקר. המנגנון שבא אחרי שכבת ממבה נקרא EinFFT וכבר הבנו ש-FFT מתאים להתמרת פורייה. אבל מה זה Ein? באופן לא מפתיע אלו 3 האותיות הראשונות נלקחו מאינשטיין. אז מה בעצם אינשטיין עושה כאן?

למעשה המאמר משתמש בסכימת אינשטיין שהיא דרך לרשום מכפלות הטנזורים או המטריצות במקרה פרטי. למעשה במקום לרשום כל איבר i, j של המכפלת מטריצות A ו- B בתור מכפלה פנימית של שורה i ועמודה j סכימת אינשטיין כותבת אותו ללא סימן של סכום (=סיגמה) אלא על ידי ציון של מספר שורה i , מספר עמודה j ואינדקס סכימה k .

אז איך המאמר משתמש בסכימה הזו? הרי אמרתי שהסכימה הזו מוגדרת גם לטנזורים ומתברר שלחבילות תוכנה כמו pytorch יש חבילות שיודעות לבצע מכפלת טנזורים רב מימדיים המבוטאים דרך סכימת אינשטיין בצורה די יעילה. וזה בדיוק מה שעושים במאמר. המאמר מפרק את המטריצות מהשכבות הלינאריות של EinFFT לכמה מטריצות במימד נמוך יותר ובונה מזה טנזור רב מימדי הבנוי ממטריצות בלוקיות (אפסים מחוץ לבלוקים). הטענה במאמר שזה מאפשר לבצע את המכפלות (באימון אבל כנראה גם באינפרנס) בצורה מהירה יותר על ידי ניצול טוב יותר של משאבי החומרה.

הביצועים כרגיל מפתיעים לטובה....

$$\begin{aligned}
 X &= \text{FFT}(x) \\
 X_{\text{real}_1} &= \max(\text{Re}(X) \cdot W_{1,\text{real}} - \text{Im}(X) \cdot W_{1,\text{imag}} + B_{1,\text{real}}, 0) \\
 X_{\text{imag}_1} &= \max(\text{Re}(X) \cdot W_{1,\text{imag}} + \text{Im}(X) \cdot W_{1,\text{real}} + B_{1,\text{imag}}, 0) \\
 X_{\text{real}_2} &= \text{Re}(X_{\text{real}_1} \cdot W_{2,\text{real}} - X_{\text{imag}_1} \cdot W_{2,\text{imag}} + B_{2,\text{real}}) \\
 X_{\text{imag}_2} &= \text{Im}(X_{\text{real}_1} \cdot W_{2,\text{imag}} + X_{\text{imag}_1} \cdot W_{2,\text{real}} + B_{2,\text{imag}}) \\
 X_{\text{shrink}} &= \text{softshrink}(X_{\text{real}_2} \& X_{\text{imag}_2}, \lambda = \text{sparsity_threshold}) \\
 x_{\text{ifft}} &= \text{IFFT}(X_{\text{shrink}}, \text{dim} = 1, \text{norm} = 'ortho') \\
 x_{\text{reshaped}} &= \text{reshape}(x_{\text{ifft}}, (B, N, C))
 \end{aligned}$$