

# שִׁכְחֹן דיגיטלי: הקרב על הזיכרון החוץ-גופי של המכונה

מאז ומעולם, ההיסטוריה האנושית היא רצף של טכניקות שנועדו להילחם בכוח ההרסני של השכחה. מרגע המעבר שלנו מתרבות בעל-פה ל"תרבות ארכיונית"<sup>1</sup>, המצאת הכתב אפשרה לנו לייצר "זיכרון חוץ-גופי"<sup>2</sup> – כספת חיצונית לידע, לחוקים ולסיפורים, ששחררה את המוח האנושי מהצורך לשנן הכול.

כעת, בעידן הבינה המלאכותית (AI), אנו עומדים מול פרדוקס קיומי-טכנולוגי: **מודלי השפה הגדולים (LLMs)**, הישגיות האינטלקטואליות החזקות ביותר שיצרנו, סובלים ממה שמוגדר כ"אמנזיה קונטקסטואלית"<sup>3</sup>, או בשפה העברית הרשמית: **שִׁכְחֹן**<sup>4</sup>.

## שִׁכְחֹן והבל הזיכרון המכני

השִׁכְחֹן הזה אינו תוצאה של כשל, אלא של **מגבלה מבנית**: מודלי LLM הם מטבעם "**חסרי מצב (Stateless)**"<sup>5</sup>. משמעות הדבר היא שכל הפעלה של המודל היא אירוע חד-פעמי, חדש לחלוטין, והוא אינו זוכר באופן אינהרנטי את ההיסטוריה, הכללים או התוצאות של הפעלה קודמת<sup>6</sup>. כאשר "חלון הקונטקסט" של המודל מתמלא, או כאשר סשן העבודה מסתיים, המידע הקריטי **נמחק**, והקונטקסט נעלם<sup>7</sup>. זוהי "אמנזיה של המכונה"<sup>8</sup>, המאלצת את הסוכן להתחיל "מאפס" בכל פעם, והאדם נדרש להסביר שוב ושוב מהם הכללים ומהו הפרויקט<sup>9</sup>.

בכך, המכונה משקפת בפנינו את ההבל הקהלתי: כל עמל, כל מאמץ חישובי, עלול להפוך לחסר ערך ברגע שהזיכרון נדחק החוצה. הפתרון ההנדסי לבעיה זו הוא יצירת **זיכרון לטווח ארוך (Persistent Memory)** מלאכותי עבור סוכני ה-AI<sup>10</sup>, המחייב את המודל לבצע "הקצאת משאבים קוגניטיבית" כדי להתגבר על השכחה.

## המנגנון הקוגניטיבי: תקציב Tokens ואכיפה

כדי שמערכת ארבעת קובצי הזיכרון (**Code Memory**) תפעל, על מודל ה-LLM לקרוא את תוכן הקבצים הללו בתחילת כל סשן עבודה<sup>11</sup>. דחיפת התוכן של קובצי הזיכרון לתוך חלון הקונטקסט היא למעשה הדרך שבה הסוכן "**נזכר**" באופן יזום בכללים, במבנה ובמשימות<sup>12</sup>. זהו המנגנון הנדרש כדי להתמודד עם המגבלה המובנית של השִׁכְחֹן הקונטקסטואלי<sup>13</sup>. המנגנון הזה יוצר "לולאת משוב קוגניטיבית" 14 החיונית לרציפות:

הסוכן קורא  $\rightarrow$  מבין  $\rightarrow$  מבצע  $\rightarrow$  מעדכן  $\rightarrow$  הסשן הבא קורא את העדכון  $\rightarrow$  ממשיך מהנקודה המדויקת שבה הסשן הקודם הפסיק<sup>15</sup>. זוהי הדמיה הנדסית של זיכרון ארוך-טווח<sup>16</sup>.

## תהליך עבודה חובה בתחילת כל סשן:

1. קרא את **PLANNING.md** – הבן את הארכיטקטורה ואת השלבים<sup>17</sup>.
2. בדוק את **TASKS.md** – ראה מה הושלם, מה הבא בתור, מה התלויות<sup>18</sup>.
3. עבוד על המשימה הבאה בתור<sup>19</sup>.
4. סמן משימות כהושלמו מיד עם תאריך (**Mark Immediately**)<sup>20</sup>.
5. הוסף משימות חדשות שהתגלו תוך כדי עבודה (תפקידו של **TASKS.md** להיות "**מסמך חיי**")<sup>21</sup>.

## הקצאת תקציב Tokens אופיינית: היררכיית הזיכרון הדיגיטלי

"הקצאת תקציב Tokens אופיינית"<sup>22</sup> היא חלוקה מומלצת של תקציב האסימונים (Tokens) הכולל של חלון הקונטקסט של המודל, בין ארבעת עמודי הזיכרון: **PRD.md**, **CLAUDE.md**, **TASKS.md** ו-**PLANNING.md**<sup>23</sup>. חלוקה זו אינה אקראית, אלא משקפת היררכיה קיומית שבה אכיפת חוקים קריטית יותר מחזון<sup>24</sup>.

חלוקת התקציב האופיינית המוצעת:

קובץ	ייעוד	אחוז אופייני מתקציב הקונטקסט
CLAUDE.md	ספר החוקים הקנוני (אכיפת כללים)	30%–25% <sup>25</sup> (הכי קריטי)
PLANNING.md	הארכיטקטורה הטכנית (המפה)	25%–20% <sup>26</sup>
TASKS.md	סטטוס ביצוע (יומן המשימות החי)	25%–20% <sup>27</sup>
PRD.md	חזון ואסטרטגיה (המוטיבציה)	20%–15% <sup>28</sup>
היתרה	דיאלוג, תוצאות כלים, קריאת קוד	20%–10% <sup>29</sup>

#### עקרונות תכנון המשקלים:

1. **CLAUDE.md (המשקל הגבוה ביותר, 30%–25%)**: זהו ספר החוקים של הפרויקט<sup>30</sup>. הוא מקבל את המשקל הגבוה ביותר משום שתפקידו **לאכוף מגבלות טכניות** (כגון הנחיות פורמט ספציפיות, למשל "השתמש ב-`\textthebrew` לטקסט בעברית")<sup>31</sup>. אכיפה אוטומטית של כללים אלה בתחילת כל ששן מונעת טעויות חוזרות<sup>32</sup>.
2. **PLANNING.md (25%–20%)**: קובץ זה מתאר את הארכיטקטורה והאסטרטגיה הטכנית (איך הדברים נבנים)<sup>33</sup>. קריאה שלו מאפשרת לסוכן להבין היכן הוא נמצא ב"מפה" של תהליך העבודה<sup>34</sup>.
3. **TASKS.md (25%–20%)**: זהו "פנקס הביצוע" או הסטטוס העדכני של המשימות שהושלמו והמשימות שנותרו<sup>35</sup>. קריאתו מאפשרת לסוכן לדעת מהי הפעולה הבאה שיש לבצע<sup>36</sup>.
4. **PRD.md (המשקל הנמוך ביותר, 20%–15%)**: קובץ זה מספק את החזון והדרישות העסקיות (מה עושים)<sup>37</sup>. הוא חשוב כקונטקסט רקע ("המוטיבציה"), אך הוא פחות קריטי לביצוע המידי מאשר הכללים (CLAUDE.md) או המצב הנוכחי (TASKS.md)<sup>38</sup>. התכנון נועד להבטיח שבכל ששן עבודה, גם כאשר הסוכן הוא "**חסר מצב**", הוא יקבל באופן מיידי את **הקונטקסט הקוגניטיבי** הנדרש כדי להמשיך את העבודה באופן עקבי, כאילו הוא שותף קוגניטיבי מתמשך<sup>39</sup>.

#### המשמעות המעשית של הקצאת הזיכרון

המשמעות המעשית של הקצאת תקציב קבועה מראש היא יצירת **זיכרון לטווח ארוך (Persistent Memory)** עבור סוכני ה-AI<sup>40</sup>:

- **אכיפת כללים אוטומטית**: משקל גבוה ל-CLAUDE.md מבטיח שהסוכן אוסף אוטומטית את "חוקי המשחק" בכל אינטראקציה, מה שמונע טעויות חוזרות<sup>41</sup>.
- **עקביות רציפה**: הסוכן אינו מתחיל "מאפס"; הוא מתחיל מהמצב המדויק שבו הפסיק בסשן הקודם<sup>42</sup>.
- **מניעת שִכחון קונטקסטואלי**: המנגנון מאלץ את הטעינה מחדש של הזיכרון הקריטי<sup>43</sup>.
- **שיפור הפרודוקטיביות**: נחסך הצורך ב"הסבר חוזר" למשתמש, מה שהוכח כמשפר את ביצועי הסוכן במשימות מורכבות בעד 39% בהשוואה למערכות ללא ניהול קונטקסט<sup>44</sup>.

#### דוגמה בסיסית לתכנון משקלים:

נניח שלמודל ה-LLM (כגון Claude) יש חלון קונטקסט כולל של **40,000 אסימונים (Tokens)**<sup>45</sup>, ואנו משתמשים במשקלים הטיפוסיים העליונים:

קובץ	משקל ל אחוז (ז)	חישוב Tokens ( $40,000 \times$ שקל)	משמעות מעשית (דוגמה)
CLAUDE. md		<sup>46</sup> 30%	$\mathbf{40,000 \times 0.30 = 12,000 \text{ Tokens}}$ <sup>47</sup>
PLANNING. md	25 <sup>49</sup> %		"השלב הנוכחי הוא 7 מתוך 10" (היכן אנחנו בתהליך) <sup>51</sup> .
TASKS.md	25 <sup>52</sup> %		"משימה 12 הושלמה. המשימה הבאה היא 13" (סטטוס עבודה) <sup>54</sup> .
PRD.md	20 <sup>55</sup> %		"המטרה היא ייצוא נתונים ל-CSV" (החזון הכללי) <sup>57</sup> .

**אופן התכנון:** התכנון מחייב שהקובץ CLAUDE.md יהיה קצר מ-12,000 אסימונים, וששאר הקבצים יישארו בטווח שהוקצה להם, או שיחולקו לקבצי משנה ויקראו רק חלקית (קריאה מדורגת)<sup>58</sup>.

#### נוסחת החישוב ומקור הנתונים

הנוסחה לחישוב מגבלת האסימונים המרבית  $T_i$  עבור קובץ  $i$  (מתוך תקציב כולל  $T_{\text{Total}}$ ) היא נוסחת הקצאה לינארית פשוטה<sup>59</sup>:

$$T_i = T_{\text{Total}} \times P_i$$

כאשר:

- $T_i$ : מספר האסימונים המרבי המוקצה לקובץ  $i$  (למשל, CLAUDE.md)<sup>60</sup>.
- $T_{\text{Total}}$ : מספר האסימונים הכולל של חלון הקונטקסט של המודל (לדוגמה, 200,000 אסימונים עבור Claude 3.5 Sonnet)<sup>61</sup>.
- $P_i$ : המשקל (האחוז) המוקצה לקובץ  $i$  (למשל, 0.30 עבור CLAUDE.md)<sup>62</sup>.

**מקור המספרים והאסמכתאות:** המספרים והאחוזים המוצגים אינם תיאורטיים מופשטים, אלא נובעים מתוך הקיטקטורה ההנדסית (**Code Memory**) שפותחה כפתרון מעשי לבעיות LLMs כמו שֶׁחֶן<sup>63</sup>. חלוקה זו היא "שימוש הנדסי מעשי" (**Engineering Implementation**), שנבחנה והשתכללה דרך הניסיון של Anthropic בניהול קונטקסט לטווח ארוך<sup>64</sup>. הוכחה פרקטית מראה כי מערכת זו שימשה לבניית הספר עצמו (גרסה 4.0), והשיגה תוצאות איכותיות כמו **100% תאימות לכלים**<sup>65</sup>.

**אופטימיזציה:** הסכום של כל המשקלים חייב להיות נמוך או שווה ל-100%<sup>66</sup>. אם חלון הקונטקסט קטן, נדרשת אופטימיזציה נוספת, כמו שימוש בטכניקות **Prompt Caching** הנתמכות על ידי מודלים מודרניים (כגון Claude 3.5 Sonnet) ומפחיתות את עלות האסימונים של קבצים חוזרים (כמו CLAUDE.md) **בעד**<sup>67</sup> **90%**.