

פרק 9

הפרישה - מעבדה לייצור בעולם האמיתי

תקציר

המעבר מפרויקט ניסיוני מוצלח למערכת ייצורית פعليה הוא אחד האתגרים המורכבים ביותר בהטמעת בינה מלאכותית. פרק זה בוחן את אפשרויות הפרישה השונות - מתשתיות מקומיות דרכ' פתרונות ענן ועד ארכיטקטורות היברידיות - ומספק למנהלים את הכלים לקבלת החלטות מושכלות. נלמד לחשב עלויות אמיתיות, לתכנן סקלబיליות, ולנהל את המורכבות הטכנולוגית של מערכות AI ביצור. זהו הפרק שבו תיאוריה הופכת למעשה, והחלומות הופכים למציאות תפעולית.

1.9 מטרות הלמידה

בתום פרק זה תוכלו:

- להבין את אפשרויות הפרישה השונות: O-nH, duolC, sesimerP
- לקבל החלטות פרישה מושכלות על בסיס OCT ושיקולים עסקיים
- לנוהל C-sreniatno ותשתיות לפרישת מערכות AI
- לתכנן ולהוציא לפועל מעבר מסביבת פיתוח לייצור
- להבין עקרונות S-gnilacS וניהול גדילה
- לבנות תוכנית RD ו-oG-L-evi מקיפה

2.9 רגע האמת: מעבדה לייצור

כשאדם ראשון מכניס מפתח לדלת ביתו ברגע סמלי של מעבר דירה, רגש ההתרגשות מתמזג עם חרדה הלא-נון. בדומה לכך, כשמנהלת טכנולוגיה מפעילה לראשונה מערכת בינה מלאכותית בסביבת ייצור, היא חווה את אותו תערובת רגשות: התרגשות מפוטנציאלי הטכנולוגיה, חרדה מפני כשלים אפשריים, ותקוות שהשकעת החודשים האחרוניים תניב פירות.

בעולם הפיתוח, ההבדל בין סביבת מעבדה לבין ייצור הוא כמו ההבדל בין דגם ארכיטקטוני מושלם לבין בניית אמיתי שבו אנשים חיים [7]. במעבדה, הכל שולט וمبוקר: כמות הנתונים מוגבלת, מספר המשתמשים קטן, והטעויות הן הזדמנויות למידה. בייצור, המיציאות אוצרית יותר: אלף משתמשים בו-זמנית, נתונים לא צפויים, דרישות זמינות 7/42, וכל שגיאה עלולה לגרום במוניטין החברה או בחווית הלקוחות.

פרק זה עוסק במסע המורכב הזה - המעבר מפרויקט פילוט מוצלח, שאולי רץ על המחשב הניד של מנהל המוצר, למערכת ארגונית שמשרתת אלף משתמשים ביום. נבחן את שלוש הדרכים העיקריות לפritis מערכות בינה מלאכותית, נלמד לחשב את העליות האמיתיות של כל גישה, ובין כיצד לתכנן תשתיות שתוכל לצמוח עם הצרכים העסקיים.

3.9 שלוש דרכי לפritis: O-nP-merP-nO, duolC, sesimerP-nO

1.3.9 O-nP-merP-nO: השיטה המוחלטת

דמיינו בנק גדול בתל-אביב, שהוק הבנקאות מחייב אותו לשמר את כל נתוני הלகוחות בתוך גבולות המדינה, תחת שליטה פיזית מלאה. עבור ארנון זהה, פריסה מקומית (sesimerP-nO) אינה רק העדפה - היא הכרה. המשמעות היא שכ התשתיות - שירותי, אחסון, רשת, ומודלי הבינה המלאכותית - נמצאים בחדרי השירותים של הארגון עצמו. היתרונות של גישה זו מרשים: שליטה מוחלטת על הנתונים, אפשרות להתאים כל פרט לצרכים הספציפיים, אי-יתלוות בספק ענן חיצוניים, ואפשרות לענות על דרישות רגולטוריות חמורות. בשנתוני המטופלים במערכת בריאות או סודות מסחריים רבים מעורבים, הידעם שה מידע מעולם לא עוזב את חדר השירותים מספקת שקט נפשי לא מבוטל. אך יחד עם היתרונות באים אתגרים משמעותיים. הארגון צריך לרכוש ציוד חומרה יקר - כרטיסי מסך UPG לריצת מודלים גדולים עולמים عشرות אלפי דולרים כל אחד. יש לשכור או להכשיר צוות טכני מיומן שיתחזק את התשתיות 7/42. עלויות החשמל והקיורו של שירותי חזקים יכולות להגיע לאלפי שקלים בחודש. והחשוב מכל - קשה מאוד לצמוח במהירות: רכישת שרת חדש יכולה לקחת שבועות או חודשים.

דרישות טכניות לפritis מקומית

1. חומרה מתקדמת:

- שירותי עם UPG חזקים (לדוגמה: H001A, AIDIVN 001)
- זיכרון MAR נרחב (BT1-BG652 למודלים גדולים)
- אחסון מהיר eMVN/DSS (טריה-בייטים)
- רשת פנימית מהירה (spbG01 ומעלה)

2. תשתיות תומכת:

- מערכת קיורו יעילה לשירותים
- מקור כוח עצמאי (SPU) ומערכות גיבוי
- אבטחה פיזית לחדר השירותים
- מערכות גיבוי וניהול אסונות

3. כוח אדם מיומן:

- מהנדסי spOveD לניהול התשתיות
- מנהלי מערכת swodniW/xuniL/W
- מומחי אבטחת מידע
- מהנדסי IA/LM לתחזוקה ושיפור

2.3.9 desaB-duolC: הגמישות האינסופית

בניגוד לבנק, סטארטאפ טכנולוגי בהרצליה שרצה לבנות כל IA לנитוח טקסט אינו רוצה להשקיע מיליוני שקלים בתשתיות לפני שהוכח את עצמו בשוק. עבורו, פתרונות הען של A beW nozamA mroftalP duolC elgooG tfosorciM secivreS שימוש, גמישות אינסופית, והתחלה מהירה [5].
מודל הען מבוסס על רעיון פשוט אך מהפכני - במקומות קבועים ולתזק תשתית, הארגון משכיר כוח חישוב לפי הצורך. צרכים UPG חזק לשעתים כדי לאמן מודל? משלמים רק על השעתים. השימוש גדול פתאום פי עשרה? התשתיות מתרחבת אוטומטית. הפרויקט נכשל? מפסיקים לשלם ואין הוצאות קבועות.
השחקנים הגדולים בתחום מציעים היום שירותים IA מותחכמים מאוד:

(SWA) secivreS beW nozamA

rekaMegaS nozamA -- LM: פלטפורמה מלאה לבניה, אימון ופרישה של מודלי

kcordeB nozamA -- (natiT ,amall ,edualC) דרך IPA אחד מוביילים sMLL

dneherpmoC nozamA -- PLN מנהלים לניטוח טקסט

tcartxeT nozamA -- חילוץ טקסט ומידע ממסמכים

secnatsnI UPG 2CE -- שירותים וירטואליים עם UPG לצרכים מותאמים

eruzA tfosorciM

ecivreS IAnepO eruzA -- ALS 1o-TPG ,o4-TPG ,4-TPG עם גישה ארגונית ל-

spOLM: פלטפורמה מלאה ל- gninraeL enihcaM eruzA --

secivreS evitingoC eruzA -- (egaugnaL ,hceepS ,noisiV) מגוון שירותים IA מוכנים

oidutS IA eruzA -- סביבה מאוחדת לפיתוח אפליקציות IA

(PCG) mroftalP duolC elgooG

xetreV -- IA: פלטפורמה מלאה לכל מחזור החיים של LM

elgooG/inimeG artlU/orP inimeG -- IPA: גישה למודלי

duolC -- UPT: מעבדי rosneT מיוחדים לאימון מהיר

larutaN -- IA egaugeN PLN מתקדמים כל

אך הען אינו ללא מחיר. עלויות יכולות גדול במהירות - שימוש אינטנסיבי ב-UPG יכול להגיע לאלפי דולרים ביום. העברת נתונים החוצה מהען (ssergE) עולה כסף. והטלות בספק יכולה להפוך לכלא זהב - מעבר ספק מצב אTEGRIS טכניים ועסקיים לא פשוטים. במקרה, שאלות של פרטיות ורגולציה יכולות להיות מורכבות יותר כשהנתונים נמצאים בשירותים זרים.

3.3.9 dirbyH: הטוב משני העולמות

ארגון פיננסי בינלאומי מצא את עצמו בדיילה: מצד אחד, נתוני ל��וחות רגושים שחביבים להישאר merP-n. מצד שני, נדרש ככוח חישוב אדריך לפרויקטים ניסיוניים ואמניים. הפתרון? ארכיטקטורה היברידית - שילוב חכם של תשתיית מקומית וширוטי ענן.

בגישה היברידית, הארגון שומר את הנתונים הרגושים והמודלים המרכזיים merP-n, אך מנצל את הענן לצרכים משתנים: סביבות פיתוח וניסוי, אימון מודלים מזדמן שדורש UPG חזק, gnitsruB-1 - היכולת להתרחב זמן קצר בעומס שי. זהו הניסיון להשיג גמישות ועלות-תועלות מיטביים.

למשל, חברת סייבר ישראלית יכולה לróż את מודל זהה הייחודי שלה על שירותי מקומיים, שם זורמים נתונים הלכו-חזרו הרגושים, אך כשהיא רוצה לנשות ארכיטקטורת מודל חדשה או לאמן על דאטאסט ענק, היא מעלה סביבת ענן זמן קצר ב-SWA, עובדת שם כמה ימים, ואז מורידה הכל - ומשלמת רק על מה שהשתמשה.

סטרטגיות dirbyH נפוצות

.1 :dirbyH evitisneS-ataD

- נתונים רגושים נשאים merP-n
- שימושים קבועים (אימון, ecnerefni, בקנה מידה גדול) בענן
- תקשורת מאובטחת בין הסביבות

.2 :tilpS dorP-veD

- יוצר merP-n ליציבות ושליטה
- פיתוח ובדיקות בענן לגמישות
- תהליכי DC/IC מתואימים

.3 :ygetartS gnitsruB

- קיבולת בסיס merP-n
- התרחבות אוטומטית לענן בעומס שי
- חזרה לשתייה מקומית כשהעומס יורך

.4 :dirbyH duolC-itluM

- ליבה merP-n
- שימוש במספר ספקי ענן לפי יתרונות (SWA לכוח חישוב, eruzA ל-Opc, IAnep fo ecivreS)
- מניעת rodneV ni-kcoL

4.9 נוסחאות מנהליות: חישוב OCT

מנהלים מנוסים יודעים שהעלות האמיתית של טכנולוגיה לעולם אינה רק מחיר התג [6]. pihsrenwO fo tsoC OCT (OCT) - עלות הבעלות הכוללת - היא המدد האמיתית שלוקח בחשבון את כל ההוצאות לאורך זמן. הבה נגדיר את הנוסחאות המנהליות לחישוב OCT עבור כל גישת פרישה.

1.4.9 duolC OCT - עלות בעלות בענן

$$TCO_{Cloud} = (\text{Compute} \times \text{Hours}) + \text{Storage} + \text{Egress} + \text{Support}$$

הסבר המרכיבים:

(MAR ,UPG ,UPC) -- etupmoC : עלות שעתי של משאבי חישוב --

sruoH : מספר השעות שהשתטים פעילים --

egarotS : עלות אחסון נתונים ומודלים --

ssergE : עלות העברת נתונים החוצה מהענן --

troppuS : חבילות תמיכה טכנית --

דוגמה מסכנית:

חברת SaaS מרכiza מודל IA על SWA :

Compute: g5.2xlarge (GPU instance) = \$1.21/hour

Hours: 730 hours/month (24/7)

Storage: 500GB at \$0.10/GB = \$50

Egress: 1TB at \$0.09/GB = \$90

Support: Business support = \$100/month

$$\begin{aligned} TCO_{Cloud_monthly} &= (1.21 \times 730) + 50 + 90 + 100 \\ &= 883.3 + 50 + 90 + 100 \\ &= \$1,123.3 \text{ per month} \\ &= \$13,479.6 \text{ per year} \end{aligned}$$

2.4.9 merP-nO OCT - עלות בעלות מקומית

$$TCO_{On-Prem} = \frac{\text{Hardware}}{5} + \text{Electricity} + \text{Cooling} + \text{HR} + \text{Facilities}$$

הסבר המרכיבים:

(5/eraudraH) -- 5 : עלות החומרה מחולקת ב-5 שנים (הנחה פחות)

yticirtcelE : עלות חשמל לתפעול השירותים --

gnilooC : עלות קירור חדר השירותים --

RH : עלויות כוח אדם (מנהל מערכות, spOveD , אבטחה)

seitilicaF : שכירות/תחזוקת מבנה, אבטחה פיזית --

דוגמה מסכנית:

אותה חברת שוקלת לעבור ל-O-n-merP :

Hardware:

- Server with NVIDIA A100 GPU: \$15,000
 - Storage: \$5,000
 - Network equipment: \$3,000
- Total: \$23,000

$$\text{Hardware}/5 = \$23,000 / 5 = \$4,600/\text{year}$$

Electricity:

- Server power: $500\text{W} \times 24\text{h} \times 365 \text{ days} = 4,380 \text{ kWh/year}$
- Cost at \$0.15/kWh = \$657/year

$$\text{Cooling: } \sim 50\% \text{ of electricity} = \$328/\text{year}$$

HR:

- 0.25 FTE DevOps engineer at \$120k/year = \$30,000/year

Facilities:

- Data center space: \$200/month = \$2,400/year

$$\begin{aligned} \text{TCO_On-Prem_yearly} &= 4,600 + 657 + 328 + 30,000 + 2,400 \\ &= \$37,985 \text{ per year} \end{aligned}$$

3.4.9 - נקודת האיזון tnioP nevE-kaerB

$$\text{Break-Even} = \frac{\text{TCO}_{\text{On-Prem}}}{\text{TCO}_{\text{Cloud}}}$$

נוסחה זו מראה כמה שנים ייקח עד שההשקעה ב-Cloud תחל לחשוף לעומת השקעה ב-Cloud. ערך.

דוגמה מהדוגמאות לעיל:

$$\text{Break-Even} = \$37,985 / \$13,479.6 = 2.82 \text{ years}$$

משמעותו: ערך מאפשר שינוי מהירים; Cloud דורש תכנון מוקדם, רתווי תזוחף מס. ה云端 יהיה(Cloud) מילא(Cloud) מילא(Cloud).

שיקולים נוספים מעבר לנושאות:

-- גמישות: ערך מאפשר שינויים מהירים; Cloud דורש תכנון מוקדם

-- סיכון: השקעה מוקדמת ב-Cloud לעומת תשלום שוטף בערך.

-- רגולציה: חלק מהתעשיות מחייבות Cloud

-- אבטחה: שליטה מלאה Cloud לעומת תלות בספק בערך.

-- מומחיות: האם יש לנו את הכספיים להניל Cloud?

rekcoD :sreniatnoC 5.9 כמנוע הփירה

אם נחזור רגע לאנלוגיה של בניין, דמיינו שבמוקום לבנות כל דירה מאפס, היינו יכולים להביא קופסאות סטנדרטיות שמכילות את כל מה שצורך - קירות, חשמל, אינסטלציה - ופשוט "להרכיב" אותן על המבנה. זה בדיקת מה ש-sreniatnoC עושים בעולם הփירה הטכנולוגי.

rekcoD, הכליל המוביל לניהול sreniatnoC, שינה לחלוטין את אופן הփירה של "ישומים" [2]. במקום להתקין את nohtyP, את כל הספריות, את המודול, ואת כל התלוויות על כל שרת בנפרד - תהליך מסובך ושגיאתי - אנחנו "אורזים" הכל לתוך reniatnoC: תמונה סטנדרטית ש מכילה את כל מה שצורך. reniatnoC זה יכול לזרוץ על המחשב הנייד של המפתח, על שרת בעירוב, או על מאות שרתים בענן - והוא יתנהג בבדיקה אותו דבר.

1.5.9 MADoC sreniatnoC קרייטיים ל-IA ?tnemyolpeD

1. עיקריות סביבה:

- מודל שעבוד במעבדה יעבד בייצור
- "enihcam ym no skrow ti tuB" - הבעה נעלמת
- אותה גרסה nohtyP, אותן ספריות, אותה התנהגות

2. ניידות:

- העברה קלה בין merP-n merP לענן
- תמייה בכל פלטפורמות העננים הגדולות
- אין rodneV-kcoL-iw ברמת התשתיות

3. בידוד:

- כל reniatnoC פועל במרחב מבודד
- אין התנגשויות בין גרסאות ספריות
- אבטחה משופרת - בעיה ב-C reniatnoC אחד לא משפיעה על אחרים

4. יעילות משאבים:

- קלילים יותר ממכוונות וירטואליות
- התחלה מהירה (שניות)
- ניצול טוב יותר של החומרה

2.5.9 GAR למערכת elifrekcoD דוגמה:

הבה נראה איך כתובים elifrekcoD - "מתכוון" לבניית reniatnoC - למערכת GAR פשוטה:

```
# Dockerfile for RAG System

# ימסר Python תונומתם לחתה
FROM python:3.11-slim

# הדובע תיינית רדגה
WORKDIR /app
```

```

# תווילתה צבוק תא קתעה
COPY requirements.txt .

# תווילת וקתה Python
RUN pip install --no-cache-dir -r requirements.txt

# היצקיים דוק תא קתעה
COPY . .

# תינוציה פרויקט 8000 טרוף פושט
EXPOSE 8000

# הביבס ינתשם רדga
ENV MODEL_NAME="sentence-transformers/all-MiniLM-L6-v2"
ENV VECTOR_DB_PATH="/app/data/vectordb"

# הזרה תדוקפ
CMD ["python", "app.py"]

```

קובץ stnemeriuqer.txt כולל:

```

langchain==0.1.0
chromadb==0.4.22
sentence-transformers==2.2.2
fastapi==0.109.0
uvicorn==0.27.0
python-dotenv==1.0.0

```

בנייה והרכה:

```

# הנומתה תיינב
docker build -t rag-system:v1 .

# הזרה Container
docker run -d \
--name rag-prod \
-p 8000:8000 \
-v /data/documents:/app/data/documents \
-e OPENAI_API_KEY=${OPENAI_API_KEY} \
rag-system:v1

# מיגול תקידב
docker logs rag-prod

# הרינו
docker stop rag-prod

```

תזמור מערכת מורכבת: esopmoC rekcoD 3.5.9

בש망רקט ה-IA שלנו מרכיבת מספר רכיבים - אפליקציה, בסיס נתונים וקיטורי, sideR לאזכור xnigN,ehcac לניהוב - ניהולים בנפרד הופך למסורבל. esopmoC rekcoD מאפשר להציג כל המערכת בקובץ LMAX אחד ולהפעיל הכל בפקודה אחת.

דוגמה: lmy.esopmoc-rekcod מלאה:

```

version: '3.8'

services:
  # RAG Application
  rag-app:
    build: .
    container_name: rag-application
    ports:
      - "8000:8000"
    environment:
      - OPENAI_API_KEY=${OPENAI_API_KEY}
      - VECTOR_DB_HOST=chromadb
      - REDIS_HOST=redis
    depends_on:
      - chromadb
      - redis
    volumes:
      - ./data:/app/data
    restart: unless-stopped

  # ChromaDB Vector Database
  chromadb:
    image: chromadb/chroma:latest
    container_name: vectordb
    ports:
      - "8001:8000"
    volumes:
      - chroma-data:/chroma/chroma
    environment:
      - IS_PERSISTENT=TRUE
    restart: unless-stopped

  # Redis for caching
  redis:
    image: redis:7-alpine
    container_name: cache
    ports:
      - "6379:6379"
    volumes:
      - redis-data:/data
    restart: unless-stopped

```

```
# Nginx reverse proxy
nginx:
  image: nginx:alpine
  container_name: gateway
  ports:
    - "80:80"
    - "443:443"
  volumes:
    - ./nginx.conf:/etc/nginx/nginx.conf:ro
    - ./ssl:/etc/nginx/ssl:ro
  depends_on:
    - rag-app
  restart: unless-stopped

volumes:
  chroma-data:
  redis-data:
```

הפעלה:

```
# תכרעהמה לכ תעלפה
docker-compose up -d

# סיתוריהשה לכ לש מיגולב הייפז
docker-compose logs -f

# הרסהו הריצע
docker-compose down

# מינוטן תקיחם מע הריצע
docker-compose down -v
```

6.9 ניהול סביבות: vne וקונפיגורציה

אחד הטעויות הנפוצות בפרישה היא "הרדקודינג" - השיבוץ של ערכים קבועים ישירות בקוד. מפתח IPA של AnepO כתוב ישירות בקוד הפיתוח? זהו סיכון אבטחה ענק. כתובות בסיס הנתונים שונות בין פיתוח ליצור? הקוד צריך להשתנות בכל סביבה. הפתרונו? ניהול סביבות נכון באמצעות קבועים vne ומשתני סביבה.

1.6.9 קבוע vne - ניהול קונפיגורציה

קבוץ vne הוא קבוע טקסט פשוט ששומר משתני סביבה - ערכים שימושיים בין סביבות שונות. הוא לעולם לא מתווסף ל-Git (נוסיף אותו ל-.erongitig) כדי למנוע דליפת סודות. דוגמה: vne לפיקוח:

```
# .env.development
```

פרק 9. חפריטה - מעבדה לייצור בעולם האמיתי

```
# API Keys
OPENAI_API_KEY=sk-proj-xxxxxxxxxxxxxxxxxxxx
ANTHROPIC_API_KEY=sk-ant-xxxxxxxxxxxxxxxxxxxx

# Database Configuration
VECTOR_DB_TYPE=chromadb
VECTOR_DB_HOST=localhost
VECTOR_DB_PORT=8001

# Redis Cache
REDIS_HOST=localhost
REDIS_PORT=6379
REDIS_TTL=3600

# Application Settings
APP_ENV=development
DEBUG=True
LOG_LEVEL=DEBUG
MAX_CONCURRENT_REQUESTS=10

# Model Settings
DEFAULT_MODEL=gpt-4o-mini
EMBEDDING_MODEL=text-embedding-3-small
MAX_TOKENS=4096
TEMPERATURE=0.7
```

דוגמה: לייצור: vne.

```
# .env.production

# API Keys (from secrets manager)
OPENAI_API_KEY=${SECRET_OPENAI_KEY}
ANTHROPIC_API_KEY=${SECRET_ANTHROPIC_KEY}

# Database Configuration
VECTOR_DB_TYPE=pinecone
VECTOR_DB_HOST=rag-prod-xyz123.pinecone.io
VECTOR_DB_PORT=443
VECTOR_DB_INDEX=production-embeddings

# Redis Cache
REDIS_HOST=redis-cluster.internal.company.com
REDIS_PORT=6379
REDIS_TTL=7200

# Application Settings
APP_ENV=production
DEBUG=False
LOG_LEVEL=WARNING
```

```
MAX_CONCURRENT_REQUESTS=100

# Model Settings
DEFAULT_MODEL=gpt-4o
EMBEDDING_MODEL=text-embedding-3-large
MAX_TOKENS=8192
TEMPERATURE=0.3
```

2.6.9 שימוש בـ vne בקוד P_nohtyp

הספרייה `vnetod-nohtyp` מאפשרת לטעון מעתני סביבה מהקובץ:

```
# config.py
import os
from dotenv import load_dotenv

# Environment variables
OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")
VECTOR_DB_HOST = os.getenv("VECTOR_DB_HOST", "localhost")
MAX_TOKENS = int(os.getenv("MAX_TOKENS", "4096"))
DEBUG = os.getenv("DEBUG", "False").lower() == "true"

# Configuration
if not OPENAI_API_KEY:
    raise ValueError(
        "OPENAI_API_KEY must be set in environment"
    )

# Usage
from openai import OpenAI
client = OpenAI(api_key=OPENAI_API_KEY)
```

3.6.9 tseB secitcarP סודות ניהול

.1. גייגר לא פעם אף :

-- הווסף .vne.erongitig

-- שמור .vne.elpmaxe עם דמה בעית

.2. יצור - sterceS reganaM :

-- SWA sterceS reganaM / yeK eruzA / PCG tluaV / terceS reganaM

-- proCihsaH tluaV לפתרון אגנוסטי

-- ביצור ישיר vne. קובץ לא פעם אף --

3. הפרדת סביבות:

-- noitcudorp.vne. ,gnigats.vne. ,tnempoleved.vne. --
-- טעינה דינמית לפי VNE_PPA --

4. רוטציה של מפתחות:

-- החלף IPA syek מעת לעת
-- אוטומציה עם ספקי sterceS --

gnilacS 7.9: כשההצלחה מביאה אתגרים

תארו לעצמכם: הסוכן שפיתחتم לשירות לקוחות זכה להצלחה מסחררת. במקומות 001 משתמשים ביום, יש עכשו 0,000,01. תשתיית השספקה בנוחות קורשת תחת העומס. زمنי התגובה מזנים מ-2 שניות ל-30 דקות. לקוחות מתלוננים. זה רגע האמת של gnilacS - היכולת של המערכת לצמוח עם הביקוש.

gnilacS lacitreV 1.7.9 - גדילה אנכית

MAR, gnilacS lacitreV, או "elacS", פירשו להוסף יותר משבבים לשרת קיימים: יותר UPC, יותר UPG חזק יותר. זה הפתרון הפשטוט יותר - לא צריך לשנות את הארכיטקטורה, רק להחליף החומרה או לשדרג את hnatsni בענו. **יתרונות:**

- פשוט לישום - אין צורך לשנות קוד
- אין מרכיבות של תקשורת בין שרתיים
- שמירה על עקביות נתונים

חסרונות:

- מוגבל פיזית - יש גבול לכמה שמכונה אחת יכולה
- צריך לעצב את השירות כדי לשדרג emitnwoD --
- אם השירות נופל, הכל נופל eruliaF fo tnioP elgniS --
- לא חסכוני מעבר לנקודה מסוימת

gnilacS latnoziroH 2.7.9 - גדילה אופקית

HgnilacS latnoziroH, או "elacS", פירשו להוסיף עוד שרתיים במקום לשדרג את הקיימים. במקומות אחד עם MAR BG46, נשתמש ב-8 שרתיים עם BG8 כל אחד. זו הגישה המעודפת במערכות ענן ומיקרו-שירותים מודרניות. **יתרונות:**

- גדילה כמעט אינסופית - פשוט מוסיפים עוד שרתיים
- עמידות גבוהה - שרת אחד נופל, האחרים ממשיכים
- גמישות - ניתן להוסיף/להסיר שרתיים דינמית

-- חסכוני - משתמשים רק במה שצרכי

חסרונות:

-- מרכיבות ארכיטקטונית - צרי^L, ניהול^S, recnalaB, daoL, etatS

-- עלויות פיתוח גבוהות יותר

-- אתגרי סyncronization נתונים

3.7.9 - חלוקת העומס gnicnalaB daoL

כשיש לנו מספר שרתים, צרי^L מנגן שיחלק את הבקשות ביניהם באופן חכם. recnalaB daoL הוא כמו פקיד קבלה במלון שמחולק אורחים בין החדרים הפנויים. אלגוריתמים נפוצים:

-- niboR dnuoR: כל בקשה לשרת הבא בתור (פשוט ויעיל)

-- snoitcennoC tsaeL: לשרת עם הכי פחות חיבורים פעילים

-- dethgieW: שרתים חזקים יותר מקבלים יותר בקשות

-- etats noisses hsaH: אותו משתמש תמיד לאותו שרת (חשוב ל-

4.7.9 - התאמה דינמית gnilacS-otuA

בעולם האידאלי, התשתיות תגדל ותקטן אוטומטית לפי הביקוש. בשעות השיא - יותר שרתים. בלילה - פחות. זה בדיק מה ש-gnilacS-otuA עשו. בענו, אפשר להגדיר כללים:

```
# Scaling Policy Example (AWS Auto Scaling)
```

```
Min instances: 2
```

```
Max instances: 20
```

```
Desired: 5
```

```
Scale-out rule:
```

```
IF CPU > 70% for 5 minutes
THEN add 2 instances
```

```
Scale-in rule:
```

```
IF CPU < 30% for 10 minutes
THEN remove 1 instance
```

```
Scale-out rule (predictive):
```

```
IF day == "Monday" AND hour == 8
THEN set desired = 10
```

5.7.9 דוגמה מעשית: gnilacS מערכת GAR

חברת טכנולוגיה פרסה מערכת GAR לשירות לקוחות. בהתחלה, שרת אחד הספיק. אך כמספר הלקוחות גדל, הם יישמו אסטרטגיית gnilacS מותחכמת:

. רמה 1 - :gnilacS lacitreV --

(MAR BG61 ,UPC 4) egralx.3t- ל (MAR BG4 ,UPC 2) muidem.3t- --
שדרוג מ-3t- --
הספק עד 1,000 בקשות ליום --

. רמה 2 - :gnilacS latnoziroH --

-- פיצול ל-3t- --
recnalaB daoL noitacilppA secnatsni --
כל יכול לטפל ב-1,000 בקשות --
סה"כ קיבולת: 3,000 בקשות ליום --

. רמה 3 - :gnihcaC + gnilacS-otuA --

-- הוספת ehcac sideR לשאלות חוזרות (%05) --
(tih ehcac secnatsni 01-2 :gnilacS-otuA --
לפי עומס --
קבולת: 52,000,000,5 בקשות ליום --

. רמה 4 - :NDC + noigeR-itluM --

-- פרישה בשני אזורי (tsaE-SU, tseW-UE) --
NDC tnorFduoIC להזאות סטטיות --
קבולת: מעל 1,000,000 בקשות ליום --

8.9 מעבר לייצור: tsilkcehC eviL-oG

הרגע הגדול מתקרב - העלאת המערכת לייצור. זה רגע שמעורר התרגשות ופחד בעת ובעוונה אחת. tsilkeehC מסודר יכול להפוך את התהליך ממלאץ ואקראי למסודר ובטוח. הנה רשות הבדיקות המלאה שכל מנהל צריך לעבורי לפני לחיצה על הcptor.

1.8.9 שלב א': טכני

.1. בדיקות ביצועים:

- האם המערכת עומדת בעומס הצפוי? gnitseT daoL --
- מה קורה כשחרוגים מהצפוי? gnitseT ssertS --
- זמני תגובה תחת תרחישים שונים gnitseT ycnetaL --
- ביצועים לאורץ זמן (84-42 דקות) gnitseT ecnarudnE --

.2. אבטחה:

- חיפוש פרצות אבטחה gnitseT noitarteneP --
- syeK IPA לא בקוד בסודות,

-- SLT/SPTTH בכל התקשורות
 -- esubA gnitimiL etaR --
 -- IIP gniggoL --
 -- לא חשיפת IP

3. גיבויים והתאוששות:

-- תהליך גיבוי אוטומטי פעיל
 -- בדיקת erotseR מגיבוי
 -- evitcejbO emiT yrevoceR) OTR --
 -- evitcejbO tnioP yrevoceR) OPR --
 -- evitcejbO emiT yrevoceR) OTR --
 -- evitcejbO tnioP yrevoceR) OPR --

4. ניטור ותצפית:

-- מערכת gniggoL מרכזית (hetaWduolC ,knulpS ,KLE)
 -- krowteN ,ksiD ,yromeM ,UPC - scirteM --
 -- etar rorre ,tuphguorht ,ycnetal - scirteM noitacilppA --
 -- seilamona - התראות על strelA --
 -- draobhsaD לצוות התפעול

2.8.9 שלב ב': תהליכי**1. תיעוד:**

-- ארQUITטורה - תרשימים מעודכנים
 -- noitatnemucoD IPA --
 -- skoobnuR - מה לעשות בתקלות --
 -- הדרכות למשתמשים

2. צוות:

-- noitatoR llaC-nO - מי זמין מתי --
 -- htaP noitalacsE - למי להעביר בעיות --
 -- הדרכה טכנית לצוות התמיכה --
 -- mooR raW - מרחב תקשורת לחירום --

3. תקשורת:

-- הודעה למשתמשים על ההשקה
 -- egaP sutatS - דף סטטוס זמינות --
 -- ערוצי תמיכה ברורים --
 -- QAF מוקן מראש --

3.8.9 שלב ג': עסקי

1. ALS והבטחות:

-- הגדרת emitpU מחייבת (%99.99%?)

-- זמן תגובה מקסימליים

-- תהליך פיצוי על אי-עמידה

2. עלויות:

-- תקציב ברור לחודשים הראשונים

-- strelA על חרינה מתקציב

-- תוכנית אופטימיזציה

3. מדידת הצלחה:

-- sIPK - מה מודדים?

-- enilesaB - מה המצב לפני ההשקה?

-- slaoG - מה הצלחה?

4.8.9 שלב ד': אסטרטגיה

1. nalP kcablloR :

-- איך חווירים למערכת הישנה אם משחזר?

-- בדיקת kcablloR בסביבת gnigatsS

-- זמן מקסימלי להחלטה על kcablloR

2. tuolloR desahP :

-- שבוע 1: 01% מהמשתמשים (ateB)

-- שבוע 2: 05% אם הכל עובד

-- שבוע 3: 001%

-- sgalF erutaeF לשיליטה דינמית

9.9 תכנון אסונות: yrevoceR retsasiD

"מה הדבר הגורע ביוטר שיכל לקרות?" - זו השאלה שמנהלים מעדייפים לא לשאול, אך חייבים. שריפה במרכז הנתונים. מתקפת סייבר. מחיקת בסיס נתונים בטעות. סופת שלג שמשביתה את ozamA. תכנון yrevoceR retsasiD (RD) הוא הביטוח שלנו - תוכנית מפורשת איך להחזיר את המערכת לפעילות במקרה הגורע ביוטר.

1.9.9 OPR ו-OTR - שני המבדדים הクリיטיים

א. כמה נתונים אנחנו מוכנים לאבד במקרה של אסון? (OPR) evitcejbO tnioP yrevoceR

(evisnepxe ,pukcab suounitnoc) = OPR --

(ecnalab) = OPR --

(:sruoh 42 :paehc ,אבל...) = OPR --

ב. תוך כמה זמן אנחנו חייבים לחזור לפעולות? (OTR) evitcejbO emiT yrevoceR

(ybdnatS toH :setunim) = OTR --

(ecnalab) ,תשתיית קיימת אך לא כל הנתונים (ybdnatS mraW :sruoh) = OTR --

(ybdnatS dloC :syad) = OTR --

2.9.9 אסטרטגיות RD**.1 (זול, איטי) erotseR dna pukcaB**

-- גיבויים תקופתיים לאחסן זול (reicalG 3S)

-- במקרה אסון: קנה תשתיית התקון, שחזור

-- OPR: שעות/ימים, OTR: ימים

-- עלות: נמוכה, סיכון: גבוהה

.2 (אייזון) thgiL toliP

-- מרכיבים קרייטיים תמיד פעילים (BD), שאר כבוי

-- במקרה אסון: הפעל את השאר, נתב טרפיך

-- OPR: דקotas-שעות, OTR: שעות

-- עלות: בינונית, סיכון: בינוני

.3 (מהיר יותר) ybdnatS mraW

-- מערכת מוקטנת תמיד פעולה באזור אחר

-- במקרה אסון: elacS pu ונתב

-- OPR: דקotas, OTR: דקotas-שעות

-- עלות: גבוהה, סיכון: נמוך

4. evitcA/evitcA etiS-itluM (זמיןנות מקסימלית)

- שני אתרים מלאים פעילים תמיד
- במקרה אסון: אתר אחד ממשיך
(revolaf tnerapsnart) 0 :OPR ,0 --
- עלות: מאוד גבוהה, סיכון: מינימלי

3.9.9 תרגיל RD - מה הולם את הארגון שלו?

טבלה 1.9 מסכמת את בחרת אסטרטגיית RD לפי סוג הארגון:

Scenario	RPO	RTO	Strategy
Startup Blog AI Tool	24h	2 days	Backup & Restore
E-commerce AI Chatbot	1h	4h	Pilot Light
Banking AI Fraud Detection	5m	30m	Warm Standby
Healthcare Critical AI Diagnosis	0	0	Multi-Site Active/Active

טבלה 1.9: בחרת אסטרטגיית RD לפי תרחיש

01.9 דוגמאות מעשיות

1.01.9 דוגמה 1: פריסת Oamall amall לרכיב מקומית של LLM

Oamall amall הוא כלי פופולרי להרצת מודלי LLM מקומי על המחשב שלו או על שירותים פרטיים [3]. זה אידיאלי לארגוני שרצו לרשום amall 3 או amall 3.0 ללא תלות בענן.

התקנה

```
# Linux/Mac
curl -fsSL https://ollama.com/install.sh | sh

# Windows - download from https://ollama.com/download

הקיידב #
ollama --version
```

הורדת מודל והרצה

```
# Llama 3.2 (3B parameters)
ollama pull llama3.2

# תיבית קארט ניא הזרה
ollama run llama3.2
```

?תיתוכאלם הניב הז המ >>>

```
[...אך גזות הבושותה]

>>> /bye

# תרשכ הזרה API
ollama serve

# שומיש n-Python

import requests
import json

def query_ollama(prompt, model="llama3.2"):
    url = "http://localhost:11434/api/generate"
    payload = {
        "model": model,
        "prompt": prompt,
        "stream": False
    }
    response = requests.post(url, json=payload)
    return response.json()["response"]

# שומיש
answer = query_ollama("מה זה תונורתה 3 ומה AI לש?")
print(answer)
```

פרק OmallO ב-Dockerfile

```
# Dockerfile
FROM ollama/ollama:latest

# מינכום מילודים קתעה
COPY models /root/.ollama/models

# פורט 11434
EXPOSE 11434

# הזרה
CMD ["ollama", "serve"]

# docker-compose.yml
version: '3.8'

services:
  ollama:
    image: ollama/ollama:latest
```

```
container_name: local-llm
ports:
  - "11434:11434"
volumes:
  - ollama-data:/root/.ollama
deploy:
  resources:
    reservations:
      devices:
        - driver: nvidia
          count: 1
          capabilities: [gpu]

volumes:
  ollama-data:
```

יתרונות Ollama מקומי:

- פרטיות מוחלטת - נתוניים לא עזבים את הארגון
- אין עלויות IPA חוותות
- אין ycnetal נמוך - אין תקשורת רשת
- עובד אופליין

חסרוןות:

- דרש חומרה חזקה (UPG מומלץ)
- מודלים קטנים יותר מ-TPG-4 (פחות מתוחכמים)
- צריך ניהול ולעדרן בעצמו

דוגמה 2: 2.01.9 rekcoD esopmoC למערכת GAR מלאה

הבה נבנה מערכת GAR שלמה עם כל הרכיבים:

```
# docker-compose-rag-full.yml
version: '3.8'

services:
  # Frontend - Streamlit UI
  frontend:
    build: ./frontend
    container_name: rag-ui
    ports:
      - "8501:8501"
  environment:
    - BACKEND_URL=http://backend:8000
  depends_on:
```

```

    - backend
  restart: unless-stopped

# Backend - FastAPI application
backend:
  build: ./backend
  container_name: rag-api
  ports:
    - "8000:8000"
  environment:
    - OPENAI_API_KEY=${OPENAI_API_KEY}
    - VECTOR_DB_HOST=chromadb
    - REDIS_HOST=redis
    - POSTGRES_HOST=postgres
  depends_on:
    - chromadb
    - redis
    - postgres
  volumes:
    - ./data:/app/data
  restart: unless-stopped

# ChromaDB - Vector database
chromadb:
  image: chromadb/chroma:latest
  container_name: rag-vectordb
  ports:
    - "8001:8000"
  volumes:
    - chroma-data:/chroma/chroma
  environment:
    - IS_PERSISTENT=TRUE
    - ANONYMIZED_TELEMETRY=FALSE
  restart: unless-stopped

# Redis - Caching layer
redis:
  image: redis:7-alpine
  container_name: rag-cache
  ports:
    - "6379:6379"
  volumes:
    - redis-data:/data
  command: redis-server --appendonly yes
  restart: unless-stopped

# PostgreSQL - Metadata storage
postgres:

```

```
image: postgres:15-alpine
container_name: rag-db
ports:
- "5432:5432"
environment:
- POSTGRES_DB=ragdb
- POSTGRES_USER=raguser
- POSTGRES_PASSWORD=${POSTGRES_PASSWORD}
volumes:
- postgres-data:/var/lib/postgresql/data
restart: unless-stopped

# Nginx - Reverse proxy & load balancer
nginx:
image: nginx:alpine
container_name: rag-gateway
ports:
- "80:80"
- "443:443"
volumes:
- ./nginx/nginx.conf:/etc/nginx/nginx.conf:ro
- ./nginx/ssl:/etc/nginx/ssl:ro
depends_on:
- frontend
- backend
restart: unless-stopped

# Monitoring - Prometheus
prometheus:
image: prom/prometheus:latest
container_name: rag-prometheus
ports:
- "9090:9090"
volumes:
- ./prometheus/prometheus.yml:\`/etc/prometheus/prometheus.yml:ro
- prometheus-data:/prometheus
command:
- '--config.file=/etc/prometheus/prometheus.yml'
restart: unless-stopped

# Visualization - Grafana
grafana:
image: grafana/grafana:latest
container_name: rag-grafana
ports:
- "3000:3000"
environment:
```

```

    - GF_SECURITY_ADMIN_PASSWORD=${GRAFANA_PASSWORD}
volumes:
  - grafana-data:/var/lib/grafana
depends_on:
  - prometheus
restart: unless-stopped

volumes:
  chroma-data:
  redis-data:
  postgres-data:
  prometheus-data:
  grafana-data:

networks:
  default:
    name: rag-network

```

הרצה:

```

# הביבס ינתשם תרדה #
export OPENAI_API_KEY="sk-..."
export POSTGRES_PASSWORD="secure_password"
export GRAFANA_PASSWORD="admin_password"

# הלופה #
docker-compose -f docker-compose-rag-full.yml up -d

# סוטט תקידב
docker-compose ps

# מיגול #
docker-compose logs -f backend

# הריצע #
docker-compose down

```

המערכת כוללת:

```

1058:tsohlacol//:ptth-ב dnetnorF --
0008:tsohlacol//:ptth-ב IPA dnekcaB --
0003:tsohlacol//:ptth-ב draobhsad anafarG --
                            לשאלות חוזרות gnihcac sideR --
                            למטאדרטה LQSergtsoP --
                            xnigN -- esrever yxorp --
                            מלא gnirotinoM --

```

3.01.9 דוגמה 3: מעבר מ-COP duolC ל-COP dirbyH

חברת בייטוח פיתחה COP של צ'טבוט IA ב-SWA. הוכחת הרעיון עבדה מצוין, אך כעת עולה שאלת: איך לעבור לייצור עם דרישות פרטיות מחמירות?

שלב 1: COP duolC (ylnO)

-- פלטפורמה: IPA IAnepO + rekaMegaS SWA

-- משתמשים: 05 עובדי פיתוח

-- נתוניים: דאטא סינטטי, ללא נתונים לקוחות אמתיים

-- עלות: \$005/חודש

שלב 2: החלטת dirbyH

nymokim:

-- רגולציה: נתונים בייטוח חיבבים להישאר בארץ

-- פרטיות: דרישות RPDG ממחמירות

-- עלות: שימוש אינטנסיבי ב-IPA יקר מדי (כפי: \$000,5/חודש)

-- שליטה: רצון לשלוט במודול וב-Pstpmor

התוכנית:

-- נתונים 1-nO :ecnerefnl

-- אימון ופיתוח: C duolC

-- B :pukcaB (МОЦП)

שלב 3: הטמעה

1. רכישת תשתיות:merP-nO

-- שרת UPG 001A AIDIVN עם egdErewoP lleD

-- אחסון BT02 SAN

-- רשת spbG01 פנימית

-- עלות: \$000,05 חד-פעמי

2. העברת המודול:

-- מעבר מ-OAnepI 4-TPG amalL-3 B07 3 Makomiy

-- על נתונים החברה gnut-eniF

-- merP-nO tnemyolped ,SWA

3. ארכיטקטורה היברידית:

-- merP-nO :noitcudorP (ישראל)

(lartneC-UE) SWA :gnigatS/veD --
 -- NPV מאובטח ביןיהם
 -- (detpyrcne יומי ל-3S (pukcaB --

4. תהליכי IC:DC

-- פיתוח בענן
 -- בדיקות אוטומטיות
 -- tnemyolpeD יدني ל-O-n-merP אחרא איישור

שלב 4: תוצאות אחרא שנה

טבלה 2.9 מציגה את תוצאות המעבר:

Metric	Cloud POC	Hybrid Prod
Users	50	2,000
Daily requests	500	20,000
Latency (avg)	800ms	200ms
Monthly cost	\$500	\$3,000*
Uptime	99.5%	99.9%
Data privacy	Medium	High

טבלה 2.9 השוואה: COP לעומת duolC dirbyH

*כולל פחות חומרה, חשמל, כוח אדם, ו-SWA. gnigats IOR מושג בשנה 5.2.

11.9. תרגילים

1.11.9 תרגיל 1: חישוב OCT לשלווה תרחישים (תיאורטי)

תרחיש: חברת משאבי אנוש מפתחת מערכת IA לסייע קורות חיים. عليك לחשב OCT ל-3 שנים עבור כל אפשרויות פרישה.
נתונים:

- שימוש צפוי: 0,01,000 קורות חיים לחודש
- זמןعيבוד ממוצע: 03 שניות לקורות חיים
- סה"כ זמן חישוב לחודש: $0,01 \times 000 = 0,003$ שעות

1. אופציה 1: (SWA) duolC

```
Compute: c6i.xlarge @ $0.17/hour
Hours/month: 83.3 (on-demand)
Storage: 100GB @ $0.10/GB = $10
Egress: 50GB @ $0.09/GB = $4.5
Support: $0 (community)
```

Monthly: $(0.17 \times 83.3) + 10 + 4.5 = \28.66
 Yearly: $\$28.66 \times 12 = \343.92
 3-Year TCO: $\$343.92 \times 3 = \$1,031.76$

אופציה 2 : sesimerP-nO

Hardware:

- Server: \$8,000
 - Storage: \$1,000
 - Network: \$500
- Total: \$9,500

Operating costs/year:

- Electricity: $200W \times 24h \times 365d \times \$0.15/kWh = \$262.8$
 - Cooling: 50% of electricity = \$131.4
 - HR: 0.1 FTE @ \$100k = \$10,000
 - Facilities: \$100/month = \$1,200
- Total/year: \$11,594.2

3-Year TCO:

- Hardware (depreciated): \$9,500
 - Operating: $\$11,594.2 \times 3 = \$34,782.6$
- Total: \$44,282.6

אופציה 3 : dirbyH

Production: On-Prem (80% of workload)

Development: Cloud (20% of workload)

On-Prem (same as option 2): \$44,282.6

Cloud (20% workload):

- Monthly: $\$28.66 \times 0.2 = \5.73
- 3-Year: $\$5.73 \times 12 \times 3 = \206.28

3-Year TCO: \$44,282.6 + \$206.28 = \$44,488.88

שאלות:

1. איזו אופציה היא הזולה ביותר ל-3 שנים?
2. באיזו נקודת merP-nO הופך משתלם יותר מ-Cloud?
3. אילו שיקולים נוספים (מעבר לעלות) צריכים לנקוט בחשבון?
4. מה יקרה אם השימוש יגדל פי 10? חשב מחדש.

2.11.9 תרגיל 2: תכנון מעבר Cloud ל-Hybrid (תיאורטי)

תרחיש: סטארטאפ פינטק מריצ' מערכת זיהוי הונאות ב-SWA. השימוש גדל, והעלויות מטפסות. הנהלה שוקלת מעבר להיברידית.

מצב הנוכחי:

(tennoS 5.3 edualC) kcordeB SWA --

- 000,000 בדיקות ליום
-- עלות: \$000,8\$/חודש
-- 2.1 שנים: ycnetaL

משימות:

1. תכנן ארכיטקטורה היברידית - מה נשאר בענן? מה עובר ל-Cloud? merP-nO-L?
2. תכנן את לוח הזמן - איזה שלבים? כמה זמן כל שלב?
3. זהה סיכונים - מה יכול להשתחש? איך למתן?
4. הגדר מודיע הצלחה - איך נדע שהמעבר הצליח?
5. חשב ROI - תוך כמה זמן ההשקעה תחזיר את עצמה?

3.11.9 תרגיל 3: כתיבת דרישות אבטחה (תיאורטי)

תרחישון: בית חולים מתכנן לפירוס מערכת AI לאבחון מדיקל. כתוב מסמך דרישות אבטחה.
עליך לבצע:

1. אבטחת נתוניים:

- הצפנה (tisnart ni, tser ta)
-- גיבויים (תדירות, מיקום)
-- גישה (מי רשאי? noitacitnehtua ?noitazirohtua ?)

2. רגולציה:

- ציות ל-HAPII
-- ציות ל-RPDG
-- תיעוד liarT tiduA

3. תשתיות:

- noitatnemges krowteN --
selur llaweriF --
noitceted noisurtnI --

4. תהליכי:

- nalp esnopser tnedicnI --
הדרכות אבטחה לצוות
gnitset noitarteneP --

4.11.9 תרגיל 4: תכנון tsilkcehC eviL-oG (תיאורטי)

תרחיוון: חברת SaaS עם 000,01 לקוחות משתמשים במערכת IA שלך 7/42. תכנן RD.

שלב 1: הגדרת דרישות

1. מה ה-OPR המקסימלי המותר? (כמה נתונים מותר לאבד?)
2. מה ה-OTR המקסימלי המותר? (תוקן כמה זמן חייבים לחזור?)
3. מהי עלות השבתה לשעה? (ללקוחות עזבים, אובדן מוניטין)

שלב 2: בחירת אסטרטגיה

erotseR & pukcaB --

thgiL toliP --

ybdnatS mraW --

evitcA/evitcA etiS-itluM --

הצדק את הבחירה על בסיס עלות/תועלת.

שלב 3: תכנון מפורט

1. איזה רכיבים יש לגבות? (sgol ,sgifnoc ,sledom ,sesabatad)
2. היכן לאחסן גיבויים? (אזור גאוגרפי אחר, ספק אחר?)
3. איך לבדוק את הגיבויים? (תרגיל RD תקופתיים)
4. מה תהליך ההטאושות? (koobnuR שלב-אחר-שלב)

שלב 4: תרגיל RD

כתבו koobnuR מפורט: "שרת הייצור נפל בשעה 00:41. מה עושים?"

5.11.9 תרגיל 5: בניית tsilkcehC eviL-oG (תיאורטי)

תרחיוון: אתה מנהל הפROYkt של מערכת IA חדשה שתשתחרר לייצור בעוד שבועיים. בנה tsilkcehC מקיף.

צור רשימת שימושות תחת הקטגוריות הבאות:

1. טכני:

(?gnitset ssertS ?gnitset daoL) --
(?tnemeganam sterceS ?gnitset neP) --
(?ydaer nalp RD ?detset pukcaB) --
(?strelA ?scirteM ?sgoL) --

2. תהליכי:

-- תיעוד (?skoobnuR ?scod IPA ?scod erutctihcrA)
-- צוות (?gniniarT ?noitator llac-nO) --

-- תקשורת (egap sutatS ?noitacnummoc resU)

.3. עסקים:

(?devorppA ?denifeD) ALS --
-- עלויות (strela tsoC ?tegduB)
-- מדידה (derusaem enilesaB ?denifed sIPK)

.4. אסטרטגי:

(airetirc noisiceD ?detseT) nalp kcablloR --
(?)%001 □ %05 □ %01) tuollor desahP --
(?ohW ?nehW) waiver hcnual-tsoP --

6.11.9 תרגיל 6:/esopmoC rekcoD - nohtyP/ למערכת IA בסיסית (קוז)

משימה: בנה מערכת IA בסיסית עם/esopmoC rekcoD/ הלקוחה:

IPA IAnepO dnekcab IPAtsaF .1

ל毛主席 sideR .2

xnigN .3

קובץ 1: yp.ppa/dnekcab

```
# backend/app.py
from fastapi import FastAPI, HTTPException
from pydantic import BaseModel
from openai import OpenAI
import os
import redis
import json
import hashlib

app = FastAPI(title="Simple AI API")

# רובי פ-OpenAI
openai_client = OpenAI(api_key=os.getenv("OPENAI_API_KEY"))

# רובי פ-Redis
redis_client = redis.Redis(
    host=os.getenv("REDIS_HOST", "localhost"),
    port=int(os.getenv("REDIS_PORT", "6379")),
    decode_responses=True
)

class ChatRequest(BaseModel):
    prompt: str
```

```
model: str = "gpt-4o-mini"
max_tokens: int = 500

class ChatResponse(BaseModel):
    response: str
    cached: bool

def cache_key(prompt: str, model: str) -> str:
    """בנין מפתחache key"""
    content = f"{model}:{prompt}"
    return hashlib.md5(content.encode()).hexdigest()

@app.get("/")
def root():
    return {
        "message": "AI API is running",
        "version": "1.0.0"
    }

@app.get("/health")
def health():
    """תומך בתקינות"""
    try:
        redis_client.ping()
        return {"status": "healthy",
                "redis": "connected"}
    except:
        return {"status": "unhealthy",
                "redis": "disconnected"}

@app.post("/chat", response_model=ChatResponse)
def chat(request: ChatRequest):
    """שאלה AI ומחזק קידוב"""
    # בקשת cache
    key = cache_key(request.prompt, request.model)
    cached_response = redis_client.get(key)

    if cached_response:
        return ChatResponse(
            response=cached_response,
            cached=True
        )

    # שאלת OpenAI
    try:
        completion = openai_client.chat.completions.create(
            model=request.model,
```

```

        messages=[
            {"role": "user", "content": request.prompt}
        ],
        max_tokens=request.max_tokens
    )

ai_response = completion.choices[0].message.content

# שמי חניה (TTL: 1 hour)
redis_client.setex(key, 3600, ai_response)

return ChatResponse(
    response=ai_response,
    cached=False
)

except Exception as e:
    raise HTTPException(status_code=500, detail=str(e))

@app.delete("/cache/clear")
def clear_cache():
    """לכין חניה"""
    redis_client.flushdb()
    return {"message": "Cache cleared"}

if __name__ == "__main__":
    import uvicorn
    uvicorn.run(app, host="0.0.0.0", port=8000)

```

קובץ 2 : txt.stnemeriuqer/dnekcab

```

fastapi==0.109.0
uvicorn==0.27.0
openai==1.10.0
redis==5.0.1
pydantic==2.5.3
python-dotenv==1.0.0

```

קובץ 3 : elifrekcoD/dnekcab

```

# backend/Dockerfile
FROM python:3.11-slim

WORKDIR /app

COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt

COPY . .

```

EXPOSE 8000

```
CMD ["uvicorn", "app:app", \
      "--host", "0.0.0.0", "--port", "8000"]
```

קובץ 4: fnoc.xnign/xnign

```
# nginx/nginx.conf
events {
    worker_connections 1024;
}

http {
    upstream backend {
        server backend:8000;
    }

    server {
        listen 80;

        location / {
            proxy_pass http://backend;
            proxy_set_header Host $host;
            proxy_set_header X-Real-IP $remote_addr;
            proxy_set_header X-Forwarded-For
                $proxy_add_x_forwarded_for;
            proxy_set_header X-Forwarded-Proto $scheme;
        }

        location /health {
            proxy_pass http://backend/health;
            access_log off;
        }
    }
}
```

קובץ 5: lmy.esopmoc-rekcod

```
# docker-compose.yml
version: '3.8'

services:
  backend:
    build: ./backend
    container_name: ai-backend
    environment:
      - OPENAI_API_KEY=${OPENAI_API_KEY}
      - REDIS_HOST=redis
```

```

      - REDIS_PORT=6379
depends_on:
  - redis
restart: unless-stopped

redis:
  image: redis:7-alpine
  container_name: ai-cache
  ports:
    - "6379:6379"
  command: redis-server --appendonly yes
  volumes:
    - redis-data:/data
  restart: unless-stopped

nginx:
  image: nginx:alpine
  container_name: ai-gateway
  ports:
    - "80:80"
  volumes:
    - ./nginx/nginx.conf:/etc/nginx/nginx.conf:ro
  depends_on:
    - backend
  restart: unless-stopped

volumes:
  redis-data:

```

קובץ 6 : vne.

```
# .env
OPENAI_API_KEY=sk-proj-your-key-here
```

קובץ 7 : yp.ipa_tset (בדיקות)

```

# test_api.py
import requests
import time

BASE_URL = "http://localhost"

def test_health():
    """תאימות תקינות"""
    response = requests.get(f"{BASE_URL}/health")
    print(f"Health check: {response.json() }")

def test_chat(prompt: str):
    """תאימות chat endpoint"""

```

```

payload = {
    "prompt": prompt,
    "model": "gpt-4o-mini",
    "max_tokens": 200
}

# First call - should NOT be cached
start = time.time()
url = f"{BASE_URL}/chat"
response1 = requests.post(url, json=payload)
duration1 = time.time() - start
result1 = response1.json()

print(f"\nFirst call (uncached):")
print(f"  Time: {duration1:.2f}s")
print(f"  Cached: {result1['cached']} ")
print(f"  Response: {result1['response'][:100]}...")

# Second call - SHOULD be cached
start = time.time()
response2 = requests.post(url, json=payload)
duration2 = time.time() - start
result2 = response2.json()

print(f"\nSecond call (cached):")
print(f"  Time: {duration2:.2f}s")
print(f"  Cached: {result2['cached']} ")
print(f"  Speedup: {duration1/duration2:.1f}x faster")

def test_cache_clear():
    """1. בדיקת תקינות"""
    response = requests.delete(f"{BASE_URL}/cache/clear")
    print(f"\nCache clear: {response.json()}")

if __name__ == "__main__":
    print("== Testing AI API ==")

    test_health()
    test_chat("Explain Docker in one sentence")
    test_cache_clear()

    print("\n== All tests completed ==")

```

הרצה:

```

# 1. הזרתו היינט
docker-compose up --build -d

# 2. בדיקת תקינות

```

```

docker-compose ps

# מיגול . 3
docker-compose logs -f backend

# (רזה לנימרטב) הקידב . 4
python test_api.py

# תינדי הקידב . 5
curl http://localhost/health
curl -X POST http://localhost/chat \
-H "Content-Type: application/json" \
-d '{"prompt": "What is AI deployment?"}'
'

# הריצע . 6
docker-compose down

```

שאלות:

1. הסבר איך ה-gnihcac משפר ביצועים
2. מה קורה אם sideR כתוב cigol kcabbllaf?
3. הוסף IAnep kcehc htlaeh גם את O IPA
4. הוסף etar gnitimil (מקסימום 01 בקשوت לדקה למשתמש)
5. הוסף gniggol מתקדם (לקובץ -elosnoc)

21.9 סיכום הפרק

הפרק הזהלקח אותנו למסע התיאוריה אל הפרקטייה האמיתית של הפעלת מערכות בינה מלאכותית [8]. למדנו שהפרישה אינה רק החלטה טכנית, אלא החלטה עסקית וסטרטגית המשפיעה על עליות, ביצועים, אבטחה וגישות לשנים קדימה.

ראינו שלוש הדריכים העיקריים לפרישה - dirbyH, duolC, sesimerP-nO; כל אחת עם יתרונות וחסרונות שלה, ושאי תשובה אחת נכונה. בנק שמחזיק נתונים לקוחות רגילים יבחר O-nmerP; סטרטטוף שורצוה לצמוח מהר יבחר C-duol; וארגון בוגר שמחפש איזון יבחר H-dirby.

למדנו שהעלות האמיתית אינה מה שכטוב על התג - OCT כולל גם חשמל, קירור, כוח אדם, ומאות פרטים קטנים שמצטברים. נוסחאות כמו OCT duolC-1-OCT merP-nO-OCT merP-nO מאפשרות למנהל להשותת תפוחים לתפוחים ולקלב החלטות מושכלות.

בפרט, sreniatnoC, rekcoD-1-esopmoC, rekcoD, lacitreV, gnilacS, gnilacS-latnoziroH המודרני - עקביות, נידות, ובידוד [1][4]. ניהול סביבות נכון בסכנות vne, sterceS-1-sreganaM מבטיח שלא נדיף סודות ושןוכל לנوع בין סביבות בקהלות.

gnilacS - המטפס של גדילה - דורש חשיבה מראש. gnilacS-otuA והחלטות חממות יכולות להפוך מערכת שקורסית תחת gnilacS מרכיב אך אינסופי. gnilacS-latnoziroH עומס המערכת שגדלה בבחן ובشكט.

ולבסוף, retsasiD-1-urevoceR - התכוון לדבר הגרוע ביותר שיכל לקרוות - הוא לא פסימיות, אלא אחריות מڪוציאית. OPR ו-OTR מגדירים כמה נתונים אנחנו מוכנים לאבד וכמה מהר אנחנו חיבים לחזר, ואלו החלטות שמספרידות מערכת רצינית מצטצוע.

הפרק הבא יעסוק בשיקולים אסטרטגיים - כיצד לבחור בין מודלים, איך לנחל זיכרון, איך להימנע מ-V-neL-rod-kcoL-ni. אבל לפני שנטקסם, קחו רגע להעריך: האם אתם מוכנים להעלות את המערכת שלכם לייצור? אם התשובה היא כן - מזל טוב, עברתם את המבחן הקשה ביותר. אם לא - חזרו על ה-C-ehC-tilkts, תקנו מה שחרר, וnipgash בצד השני.

הפריטה אינה סוף המסע, אלא התחלה. בייצור, הכל אמייתי - המשתמשים, הנטוונים, הלחצים. זהו הרגע שבו הבינה המלאכותית שבניתם מפסיקה להיות פרויקט ונהיית מוצר. והמוצר הזה, אם תנהלו אותו נכון, יכול לשנות את העסק שלכם - ואולי את העולם.