

# به نام خدا

## داتاکیومنت پروژه سیستم عامل

نام اعضا: هومن هراتی، عمامد یدالله پور حقیقی

### 1. معرفی کلی پروژه

در این پروژه یک سرویس ذخیره‌سازی محتوا محور (Content-Addressable Storage) پیاده‌سازی شده است که فایل‌ها را به قطعات (Chunk) تقسیم می‌کند، هر قطعه را بر اساس هش ذخیره می‌کند و سپس برای هر فایل یک Manifest تولید می‌کند. شناسه نهایی فایل (CID) نیز از روی مانیفست ساخته می‌شود.

معماری پروژه دو بخش دارد:

- **GATEWAY SERVER(main.py)**: یک سرور HTTP که درخواست‌های upload و download را از کاربر می‌گیرد و آن‌ها را به موتور C منتقل می‌کند که توسط خودتان پیاده‌سازی شده بود.
- **ENGINE(c\_engine.c)**: main.py با Unix Domain Socket صحبت می‌کند و عملیات assemble، chunking، هش، ذخیره و c\_engine.c پیاده‌سازی شده است.

### 2. ساختار پوشش‌ها و فایل‌های خروجی

داده‌ها را در پوشه‌ی store نگه می‌دارد. Engine store/blocks شامل فایل‌های بازیابی chunk می‌باشد. هر chunk با هش خودش ذخیره می‌شود. store/manifest/<CID>.json شامل مانیفست فایل‌ها می‌باشد. مانیفست شامل نسخه، الگوریتم هش، chunks و لیست filename, total\_size, chunk\_size است.

### 3. نحوه اجرا

برای کامپایل اینجن C دستور زیر را وارد می‌کنیم:

```
gcc -O2 -pthread -Wall -Wextra -o c_engine c_engine.c
```

برای اجرا دستور زیر را وارد میکنیم:

```
./c_engine /tmp/cengine.sock
```

برای اجرای Gateway دستور زیر را وارد میکنیم:

```
python3 main.py
```

برای تست دانلود/آپلود دستورات زیر را وارد میکنیم:

```
CID=$(curl -s -X POST "http://127.0.0.1:9000/upload" \
-H "X-Filename: hosts" \
--data-binary "@/etc/hosts" | python3 -c 'import sys,json; print(json.load(sys.stdin)[\'cid\'])')
```

```
curl -s "http://127.0.0.1:9000/download?cid=$CID" -o out.bin
```

```
diff -q /etc/hosts out.bin && echo "OK"
```

#### 4. ارتباط بین main.py و c\_engine.c

ارتباط بین main.py و c\_engine.c از طریق Unix Domain Socket یک پروتکل فریم‌بندی باینری پیاده‌سازی شده است. روی مسیری که هنگام اجرا داده می‌شود (مثلاً /tmp/cengine.sock) listen و bind می‌کند. main.py به همین سوکت connect می‌شود و پیام‌ها را با فرمت مشخص ارسال می‌کند. هر پیام از Engine به Gateway و بالعکس شامل این ساختار است:

- opcode(1 Byte)
- length(4 Byte, big-endian)
- payload(As much as length)

پس ساختار کلی به شکل زیر است:

```
[ opcode:1 ][ length:4 big-endian ][ payload:length ]
```

این فریمینگ باعث می‌شود پیام‌ها در یک stream (سوکت) قابل تفکیک باشند و Engine دقیقاً بداند چه مقدار داده باید برای هر پیام بخواند.

در main.py، Upload مراحل زیر را انجام می‌دهد:

1. ارسال UPLOAD\_START با payload شامل متادیتا (نام فایل و اندازه).
2. ارسال چند پیام UPLOAD\_CHUNK که بدن فایل را به تکه‌های باینری منتقل می‌کند.
3. ارسال UPLOAD\_END برای اعلام پایان آپلود.

در handle\_upload Engine دقیقاً همین ترتیب را می‌خواند:

- بعد از دریافت START، نام فایل و total\_size را ذخیره می‌کند.
- با دریافت هر CHUNK، همان bytes را روی یک بافر streaming می‌گیرد و به chunk های 256KB تقسیم می‌کند.
- برای هر chunk هش محاسبه می‌شود و اگر قبل از ذخیره شده باشد دوباره نوشته نمی‌شود.
- در پایان، مانیفست ساخته می‌شود و CID از هش مانیفست تولید می‌شود.
- پاسخ نهایی Engine به Gateway شامل JSON یا داده‌ی لازم برای برگرداندن cid است.

در main.py، Download مراحل زیر را انجام می‌دهد:

1. پیام DOWNLOAD\_START را با payload شامل cid ارسال می‌کند.
- :handle\_download در Engine
- مانیفست مربوط به cid را از store/manifests/<cid>.json می‌خواند.
  - سپس chunk ها را طبق index به ترتیب assemble می‌کند.
  - داده خروجی به صورت چند فریم ارسال می‌شود تا Gateway بتواند پاسخ HTTP را بدون Content-Length برگرداند.

اگر Engine به هر دلیل (mismatch، block missing، manifest missing) نتواند دانلود را کامل کند، یک پاسخ خطابه ارسال می‌کند تا Gateway HTTP error برگرداند. برای جلوگیری از باقی‌ماندن سوکت روی سیستم، هنگام توقف Engine مسیر سوکت unlink می‌شود. برای اینکه توقف با Ctrl+C درست کار کند، handling سیگنال به شکل اصولی پیاده‌سازی شده تا accept() آزاد شود و برنامه به مسیر cleanup برسد.

## 5. منطق Chunking و ذخیره‌سازی

فایل به صورت streaming خوانده می‌شود؛ به محض پر شدن chunk 256KB یک chunk ساخته می‌شود. آنچه ممکن است کوچک‌تر باشد. هر chunk هش می‌شود و با همان هش روی دیسک ذخیره می‌شود. اگر فایل‌های مختلف chunk‌های یکسان داشته باشند، آن chunk فقط یک بار روی دیسک وجود خواهد داشت. پس از پایان upload، لیست chunks و متادادا در قالب JSON نوشته می‌شود. CID از روی محتوای مانیفست ساخته می‌شود (هش متن JSON مانیفست) تا پایدار و قابل بازتولید باشد.

## 6. همزمانی (Concurrency)

برای پشتیبانی از چند درخواست همزمان Engine:

- اتصال‌های همزمان را می‌پذیرد (multi-client).
- عملیات‌های سنگین مثل هش/نوشتن chunk‌ها در thread pool انجام می‌شود.
- هنگام خواندن مانیفست و assemble دانلود، ترتیب خروجی حفظ می‌شود.

## 7. رابط گرافیکی (GUI)

رابط گرافیکی خارج از ترمینال نیز فراهم شده است:

GUI کلاینت (a)

فایل gui\_client.py امکان انتخاب فایل و upload و download CID را از طریق یک UI ساده فراهم می‌کند.

GUI اینجین (b)

فایل engine\_viewer.py امکان اجرای c\_engine از داخل پنجره، نمایش زنده‌ی لگهای engine start/stop و store را فراهم می‌کند. این GUI باعث می‌شود عملکرد Engine بدون ترمینال قابل مشاهده باشد.

## 8. استفاده از هوش مصنوعی

برای بخش ارتباطی بین main.py و c\_engine.c از یک ابزار هوش مصنوعی به عنوان راهنمای استفاده کردیم. روند کار به این شکل بود:

- ابتدا ساختار پیام‌ها و فریمینگ (Gateway (main.py) را از کد (opcode + length + payload) استخراج و با صورت پروژه تطبیق دادیم.

- سپس توابع خواندن/نوشتن دقیق فریم‌ها (read\_exact / write\_exact) را در Engine طراحی کردیم تا با stream سوکت سازگار باشد.
- در نهایت جریان Upload و Download را مرحله به مرحله مطابق ترتیبی که Engine Gateway ارسال می‌کند در پیاده‌سازی کردیم و با تست صحت آن را بررسی کردیم.
- سایر بخش‌ها مانند store، chunking، manifest و saxthar همین رویکرد (راهنمایی گرفتن برای طراحی، و سپس پیاده‌سازی و تست دستی) تکمیل شده‌اند.

## 9. تست‌های انجام شده

برای اطمینان از درستی پروژه این تست‌ها را اجرا کردیم:

- diff فایل کوچک و مقایسه با Upload/Download
- assemble فایل بزرگ‌تر از 256 KB و بررسی تعداد chunk‌ها و صحت Upload
- دانلود همزمان چندباره از یک CID و بررسی یکسان بودن خروجی‌ها
- بررسی وجود store و manifest block‌ها در مسیر
- تست توقف Engine با Ctrl+C و حذف شدن socket file