

## پروژه دوم درس



دانشکده مهندسی برق و کامپیوتر

سیستم‌های عامل - پاییز ۱۴۰۰

استاد:

مهلت تحویل:

طراحان:

دکتر مهدی کارگهی

یک‌شنبه ۳۰ آبان

امید بدایی، علیرضا سلامت

### مقدمه

هدف از انجام این پروژه آشنایی با نحوه مدیریت کردن پردازش‌ها و راه‌های ارتباطی میان آنها است. در این پروژه به پیاده‌سازی یک شمارنده کلمات با استفاده از چارچوب Map-Reduce پرداخته می‌شود.

### Map-Reduce

در دنیای امروز، به دلیل گسترش اینترنت و دستگاه‌های هوشمند، روزانه حجم زیادی از داده تولید می‌شود. در گذشته، داده‌های تولیدی قابلیت ذخیره و اجرا بر روی یک دستگاه سخت‌افزاری را داشتند اما امروزه برای بسیاری از موارد این امر غیر ممکن است. Map-Reduce یک چارچوب و مدل برنامه‌نویسی برای پردازش داده‌های حجیم است و بسیاری از مفاهیم آن از زبان‌های تابع‌گرا<sup>2</sup> مانند Lisp گرفته شده‌است. در صورت علاقه می‌توانید از این [لینک](#) برای آشنایی بیشتر با Map-Reduce استفاده کنید. در ادامه به توضیح این مدل خواهیم پرداخت:

<sup>1</sup> Process

<sup>2</sup> Functional

Map-Reduce از دو بخش نگاشت<sup>3</sup> و کاهش<sup>4</sup> تشکیل می‌شود.

- در قسمت نگاشت، تعدادی پردازش برای عمل نگاشت وجود دارد که کاملاً مستقل از یکدیگر عمل می‌کنند و هیچ ارتباطی با یکدیگر ندارند. خروجی این مرحله تعدادی key-value خواهد بود که برای استفاده به قسمت کاهش ارسال می‌شود. تعداد پردازش‌های قسمت نگاشت محدودیت خاصی ندارد و می‌تواند بر اساس منابع در دسترس و نوع داده‌ها انتخاب شود. هر کدام از پردازش‌های قسمت نگاشت می‌توانند به صورت موازی اجرا شوند.

- در قسمت کاهش، خروجی‌های مرحله قبل به عنوان ورودی دریافت می‌شود و سپس بر اساس کلید، داده‌ها تقسیم می‌شوند. داده‌هایی که کلید یکسان دارند، حتماً باید به یک پردازش داده شوند. هر پردازش کاهش، بر روی مجموعه داده‌های با کلید یکسان، عملیات موردنظر را انجام می‌دهد و خروجی را ایجاد می‌کند.

در ادامه با انجام یک مثال بسیار ساده‌شده که مربوط به پروژه است، با مفاهیم به صورت کامل آشنا خواهید شد. فرض کنید یک پوشه شامل چند فایل متنی بسیار حجیم در اختیار شما گذاشته شده‌است و شما باید تعداد کلمات درون این فایل‌ها را بدست بیاورید. برای مثال، فرض کنید دو فایل A.csv و B.csv وجود دارد. محتوای آنها به صورت زیر است:



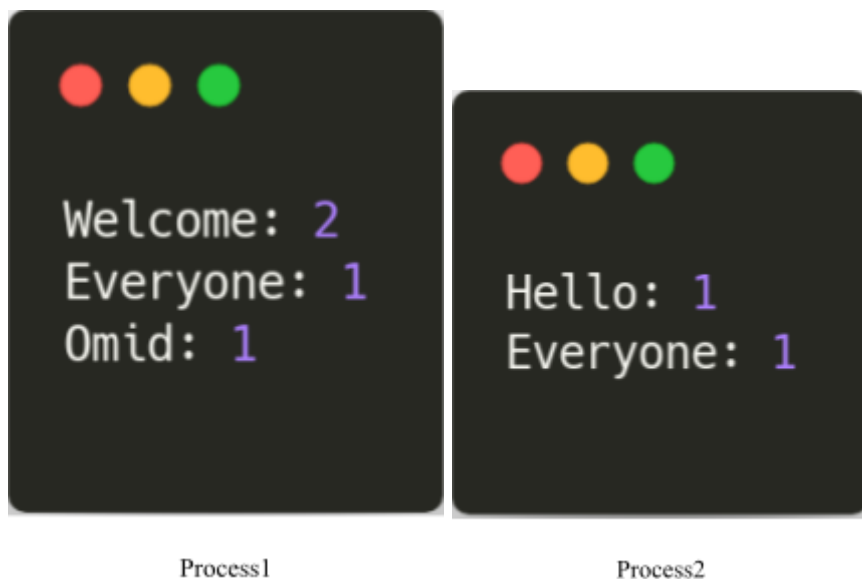
---

<sup>3</sup> Map

<sup>4</sup> Reduce

در مرحله نگاشت می‌توان برای هر فایل یا حتی برای هر خط از یک فایل، یک پردازش مجزا در نظر گرفت. حالت اول را در نظر بگیرید.

یک پردازش روی فایل اول و یک پردازش روی فایل دوم به صورت مجزا اجرا می‌شوند و خواهیم داشت:



حال در قسمت کاهش، کلیدهای یکسان به یک پردازش داده می‌شوند. البته ممکن است دو کلید متفاوت به یک پردازش داده شود اما امکان ندارد دو داده با کلید یکسان در دو پردازش مجزا قرار گرفته باشند. فرض می‌کنیم دو پردازش کاهش داریم.

کلیدهای Welcome و Everyone به پردازش اول و کلیدهای Omid و Hello به پردازش دوم داده می‌شوند. بر اساس نوع مساله که شمارش تعداد کلمات است، هر پردازش value هایی که کلید یکسان دارند را با یکدیگر جمع می‌کند و در خروجی می‌نویسد: پس خروجی هر پردازش به صورت زیر خواهد بود:



در این مرحله کار تمام شده است و می توان خروجی را ثبت کرد.

همانطور که گفته شد، این مدل بسیار ساده شده از مدل Map-Reduce است و بسیاری از مراحل برای سادگی در نظر گرفته نشده اند.

## شرح تمرین

در این پروژه، به پیاده سازی شمارنده کلمات با استفاده از روندی که در بالا توضیح داده شد خواهید پرداخت. دقت شود یک پوشه به نام testcases به شما داده می شود که حاوی تعدادی فایل csv می باشد. اسم فایل ها به صورت <num>.csv می باشد. هر کدام از این فایل ها حاوی تعدادی کلمه است که با , از یکدیگر جدا شده اند. برای مثال، فایل ها می توانند به شکل زیر باشند:

فایل 1.csv :

there,is,an,old,time,toast,which,is,golden,for,its,beauty

فایل 2.csv :

Hi,there,is,a,golden,goal,which,can,eat,toast

## معماری سامانه

برای این سامانه سه نوع پردازش با نام های پردازش اصلی، پردازش نگاشت و پردازش کاهش در نظر گرفته شده است که در ادامه توضیحات آنها به صورت کامل بیان می شود.

1. پردازش اصلی

این پردازش، پردازش والد سامانه محسوب می شود و وظیفه آن، بوجود آوردن پردازش های کاهش و نگاشت است. دقت شود فقط یک پردازش اصلی در سامانه وجود دارد. این پردازش، باید به تعداد فایل های csv موجود در دایرکتوری testcases، پردازش فرزند از نوع نگاشت ایجاد کند و نام فایل آنها را مشخص کند. دقت شود برای انتقال نام فایل به پردازش نگاشت حتما از **unnamed pipe**<sup>5</sup> استفاده شود. همچنین **یک** پردازش از نوع کاهش نیز در این پردازش ساخته می شود که خروجی نهایی را به پردازش اصلی ارسال می کند. در این مرحله پردازش اصلی باید نتایج را در فایل output.csv ثبت کند.

2. پردازش نگاشت

---

<sup>5</sup> <https://www.geeksforgeeks.org/pipe-system-call/>

هر کدام از پردازه‌های نگاشت نام فایلی را که باید بر روی آن عملیات انجام دهند، دریافت کرده‌اند. در این مرحله، هر پردازه تعداد رخداد کلمات در فایل مربوطه را محاسبه می‌کنند. خروجی این قسمت تعدادی کلید و مقدار متناظر آنها (یعنی کلمه‌ها و تعداد تکرار آن‌ها) است. هر پردازه نگاشت، خروجی خود را با استفاده از `named pipe`<sup>6</sup> به پردازه کاهش ارسال می‌کند.

### 3. پردازه کاهش

این پردازه که توسط پردازه اصلی ساخته شده‌است، منتظر می‌ماند تا خروجی تمام پردازه‌های نگاشت را دریافت کند. در این مرحله، خروجی‌ها با یکدیگر ترکیب می‌شوند و تعداد رخداد هر کلمه در میان تمام فایل‌ها محاسبه می‌شود. در نهایت، این خروجی با استفاده از `unnamed pipe` به پردازه اصلی فرستاده می‌شود.

مثال خروجی:

there: 2

is: 3

an: 1

old: 1

time: 1

toast: 2

which: 2

golden: 2

for: 1

its: 1

beauty: 1

Hi: 1

a: 1

---

<sup>6</sup> <https://www.geeksforgeeks.org/named-pipe-fifo-example-c-program/>

goal: 1

can: 1

- برای ساخت پردازنده‌ها توسط پردازنده اصلی، حتماً از فراخوانی‌های سیستمی fork و exec برای ساخت و اجرای آنها باید استفاده کنید.
- ترتیب کلمه‌ها در خروجی اهمیتی ندارد.
- فرمت انتقال داده‌ها میان پردازنده‌ها بر عهده خودتان است.
- بعد از استفاده از پایپ‌ها، آنها را ببندید.<sup>7</sup>
- به دلیل استفاده از نوع pipe ها در هر مرحله فکر کنید. در زمان تحویل سؤالاتی در این باره پرسیده خواهد شد.
- دقت شود تنها راه ارتباطی میان پردازنده‌ها استفاده از pipe است و هیچ راه دیگری قابل قبول نیست.
- هیچ نوع دیگری از پیاده‌سازی بجز مدلی که در بالا توضیح داده شد قابل قبول نیست.

نکات پایانی:

- برنامه حتماً باید با استفاده از makefile و کامپایلر g++ اجرا شود.
- برنامه باید در سیستم عامل لینوکس و در زمان معقول اجرا شود.
- تمامی نتایج را در یک فایل فشرده با اسم zip <#SID>-OS-CA2 در محل بارگذاری درس آپلود کنید.
- انجام این پروژه به صورت انفرادی است.
- نکاتی که در جلسه توجیهی یا فروم درس مطرح می‌شوند بخشی از صورت پروژه هستند لذا توصیه می‌شود که شرکت کنید.
- در صورت داشتن هرگونه سوال با تیمی‌های این پروژه در ارتباط باشید.

- [alireza.salamat@ut.ac.ir](mailto:alireza.salamat@ut.ac.ir)
- [omid.bodaghi79@gmail.com](mailto:omid.bodaghi79@gmail.com)

---

<sup>7</sup> close()