



دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر و فناوری اطلاعات

گزارش سمینار کارشناسی ارشد رشته مهندسی کامپیوتر (M.S.c)

گرایش نرم افزار

عنوان سمینار

بهبود نسبت فشرده سازی متن برای متن اسکی (بررسی و مرور)

استاد راهنما

دکتر سید علی رضوی ابراهیمی

نگارنده

مریم سادات موردگر

شهریور ۱۴۰۰



## فهرست مطالب

۱	چکیده
۱	کلمات کلیدی
۲	فصل اول مقدمه
۲	۱-۱ تاریخچه
۳	۱-۲ تعریف مسئله و بیان سؤالات اصلی تحقیق
۴	۱-۳ انگیزه و ضرورت تحقیق
۴	۱-۴ محدودیت‌ها
۵	۱-۵ پیشنهادها و هدف
۶	۱-۶ روش و مراحل تحقیق
۶	۱-۷ طرح کلی پایان‌نامه موردبررسی
۷	۱-۸ ساختار گزارش تحقیق
۸	فصل دوم مفهوم فشرده‌سازی داده‌ها و برخی از انواع آن
۸	۲-۱ مقدمه
۹	۲-۲ کدگذاری دیکشنری
۱۰	۲-۲-۱ الگوریتم‌های رایج کدگذاری دیکشنری
۱۱	۲-۳ فشرده‌سازی ASCII
۱۲	۲-۳-۱ الگوریتم شوکو (Shoco)
۱۳	۲-۴ کدگذاری هافمن
۱۳	۲-۴-۱ فرایند فشرده‌سازی
۱۵	۲-۵ کارهای مرتبط و تاریخچه تحقیق
۱۵	۲-۵-۱ تقویت فشرده‌سازی متن با رمزگذاری آماری مبتنی بر کلمه
۱۷	۲-۵-۲ فشرده‌سازی آنلاین فایل اسکی
۱۹	فصل سوم مروری بر کارهای انجام‌شده
۱۹	۳-۱ مقدمه
۱۹	۳-۲ روش‌ها و استراتژی‌های تحقیق
۲۰	۳-۳ جمع‌آوری داده‌ها
۲۰	۳-۴ تجزیه و تحلیل داده‌ها
۲۰	۳-۵ تضمین کیفیت

۲۱	۳-۶ روش‌های توسعه نرم‌افزار
۲۲	۳-۷ الزامات سیستم
۲۳	۳-۸ ویژگی مجموعه داده‌ها
۲۴	۳-۸-۱ تعریف مجموعه داده‌ها
۲۵	۳-۹ الزامات پیاده‌سازی
۲۵	۳-۱۰ طراحی TCS
۲۶	۳-۱۰-۱ طراحی ACM
۳۰	فصل چهارم ارزیابی و پیاده‌سازی الگوریتم‌ها و بررسی نتایج
۳۰	۴-۱ ارزیابی ماژول‌های کدگذاری دیکشنری و کدگذاری هافمن
۳۰	۴-۲ انتخاب پیاده‌سازی کدگذاری دیکشنری
۳۱	۴-۲-۱ مقایسه پیاده‌سازی‌ها
۳۲	۴-۳ انتخاب پیاده‌سازی کدگذاری هافمن
۳۳	۴-۳-۱ مقایسه پیاده‌سازی‌ها
۳۴	۴-۴ ارزیابی ACM
۳۵	۴-۴-۱ نتایج حاصل از ACM و Shoco
۳۶	۴-۵ ارزیابی TCS
۳۶	۴-۵-۱ نتایج حاصل از ارزیابی TCS
۳۷	۴-۵-۲ نتایج حاصل از مقایسه TCS
۳۹	۴-۶ بررسی نتایج پایان‌نامه
۳۹	۴-۶-۱ کار آیی TCS
۳۹	۴-۶-۲ مقایسه ACM و Shoco
۴۱	فصل پنجم جمع‌بندی و پیشنهادها
۴۱	۵-۱ مقدمه
۴۱	۵-۲ نتایج حاصل از تحقیق
۴۲	۵-۲-۱ پاسخ به سؤالات تحقیق
۴۳	۵-۳ پیشنهادها و کارهای آینده
۴۴	۵-۴ جمع‌بندی
۴۵	مراجع
۴۷	واژه‌نامه
۴۹	Abstrct

## فهرست جداول

۱۴	جدول ۲-۱ مثال فراوانی کاراکتر
۳۲	جدول ۴-۱ نتایج حاصل از پیاده‌سازی کدگذاری دیکشنری
۳۳	جدول ۴-۲ نتایج حاصل از پیاده‌سازی کدگذاری هافمن
۳۵	جدول ۴-۳ نتایج حاصل از مقایسه ACM و Shoco
۳۶	جدول ۴-۴ نتایج حاصل از ارزیابی TCS
۳۷	جدول ۴-۵ نتایج حاصل از ارزیابی LZW و هافمن
۳۸	جدول ۴-۶ نتایج حاصل از مقایسه TCS
۳۸	جدول ۴-۷. نتایج حاصل از مقایسه ACM + DEFLATE

## فهرست شکل‌ها

۱۴	شکل ۲-۱ درخت دودویی هافمن
۲۸	شکل ۳-۱ فهرست کاراکترهای متداول
۲۹	شکل ۳-۲ نمایش باینری از کاراکترهای ادغام‌شده

## فهرست علائم اختصاری

ASCII	American Standard Code for Information	استاندارد آمریکایی برای تبادل اطلاعات
TCS	Tool-less Compression System	یک روش فشرده‌سازی ترکیبی
ACM	ASCII Compression Module	ماژول فشرده‌سازی اسکی
PNG	Portable network graphics	گرافیک‌های قابل حمل در شبکه
Re-pair	Recursive Pairing	زوج بازگشتی
UTF-8	eight bit Unicode Transformation Format	فرمت تبدیل یونیکد ۸ بیتی
XML	Extensible Markup Language	زبان نشانه‌گذاری گسترش‌پذیر
JSON	JavaScript Object Notation	نشانه‌گذاری شی جاوا اسکریپت

## چکیده

فشرده‌سازی داده‌ها زمینه‌ای است که به‌طور گسترده مورد تحقیق قرار گرفته است. بیشتر الگوریتم‌های فشرده‌سازی که امروزه مورد استفاده قرار می‌گیرند، از چندین دهه پیش است، مانند کدگذاری هافمن<sup>۱</sup> و کدگذاری دیکشنری<sup>۲</sup> که همه‌منظوره هستند. الگوریتم‌های فشرده‌سازی می‌تواند در هر چیزی از داده‌های متنی گرفته تا تصاویر و ویدئو استفاده شود. با این حال، الگوریتم‌های بسیار کمی برای فشرده‌سازی انواع خاصی از داده‌ها، مانند متن اسکی<sup>۳</sup> وجود دارد که بدون اتلاف داده باشند. این پروژه در مورد ایجاد یک‌راه‌حل فشرده‌سازی متن با استفاده از ترکیبی از سه مورد است الگوریتم فشرده‌سازی کدگذاری دیکشنری، الگوریتم فشرده‌سازی اسکی و الگوریتم فشرده‌سازی کدگذاری هافمن.

این راه‌حل مخصوص فشرده‌سازی متن اسکی است، اما می‌تواند در هر شکلی از متن استفاده گردد. الگوریتم‌ها برای ایجاد یک نمونه اولیه ترکیب می‌شوند که بر اساس برنامه‌های فشرده‌سازی عمومی مورد ارزیابی قرار می‌گیرد. یک ارزیابی در برابر فشرده‌سازی اسکی از برنامه شوکو<sup>۴</sup> انجام خواهد شد. نتایج حاصل از ارزیابی نشان می‌دهد که ترکیب کدگذاری دیکشنری، فشرده‌سازی اسکی و کدینگ هافمن از نسبت فشرده‌سازی به‌دست‌آمده از برنامه‌های فشرده‌سازی عمومی بهتر نمی‌باشد.

(Haldar,2020)

## واژگان کلیدی

فشرده‌سازی متن، کدگذاری دیکشنری، اسکی، کدینگ هافمن

---

<sup>۱</sup>Huffman coding

<sup>۲</sup>Dictionary coding

<sup>۳</sup>Ascii

<sup>۴</sup>Program Shoco



# فصل اول

## مقدمه

فشرده‌سازی داده‌ها فرایند کاهش تعداد بیت‌های موردنیاز برای نمایش رسانه است (Sayood, 2006). به عبارت دیگر، کاهش حجم داده‌ها در حفظ تمامیت آن. فشرده‌سازی داده‌ها می‌تواند در هنگام پخش ویدئو از فضای ذخیره‌سازی، ذخیره کرده یا پهنای باند را کاهش دهد (Salomon, 2004). زیرمجموعه فشرده‌سازی داده‌ها فشرده‌سازی متن است. فشرده‌سازی متن فشرده‌سازی فایل‌های متنی است که فایل‌هایی هستند که شامل کاراکترها و نمادها می‌باشد. یک تفاوت مهم بین فشرده‌سازی متن و داده‌های باینری، مانند تصاویر و ویدئو، این است که فشرده‌سازی متن باید بدون اتلاف<sup>۱</sup> باشد. فشرده‌سازی بدون اتلاف به این معنی است که یک فایل را می‌تون فشرده کرد به طوری که فایل اصلی بدون از دست دادن اطلاعات بازیابی شود (Sayood, 2006). الگوریتم‌های فشرده‌سازی که قادر به فشرده‌سازی داده‌ها نیستند از فشرده‌سازی اتلافی استفاده می‌کنند. فشرده‌سازی با اتلاف<sup>۲</sup> بیشتر در داده‌های باینری استفاده می‌شود که در آن برخی از داده‌ها از بین می‌روند (Shukla and Prasad, 2011).

## ۱-۱ تاریخچه

فشرده‌سازی داده‌ها توسط الگوریتم‌های فشرده‌سازی انجام می‌شود. چندین الگوریتم متداول در برنامه‌های فشرده‌سازی استفاده می‌شود، و بسیاری از آن‌ها چندین دهه است که وجود دارند. الگوریتم متداولی که هنوز مورد استفاده قرار می‌گیرد، کدگذاری هافمن است که در سال ۱۹۵۲ اختراع شد (Huffman, 1952). کدگذاری هافمن کاراکترهای پرکاربردتر را با رشته‌های بیتی کوتاه‌تری نسبت به آن‌هایی که کاربردشان کمتر است، نشان داده می‌شوند. بنابراین یک متن انگلیسی ممکن است فشرده‌سازی خوبی به دست آورد، به عنوان مثال، حرف "e" بیشتر از حرف "z" استفاده می‌شود (Trost, 2020). یکی دیگر از الگوریتم‌ها، کدگذاری دیکشنری است که تکنیک جایگزینی زیررشته‌ها و کلمات با ارجاع به دیکشنری است. دیکشنری و منابع، فایل فشرده را تشکیل می‌دهند که می‌توان برای بازیابی فایل اصلی آن‌ها را از حالت فشرده خارج کرد. LZ77 و LZW از پرکاربردترین الگوریتم‌های

---

<sup>۱</sup>Lossless data compression

<sup>۲</sup>Lossy compression

کدگذاری دیکشنری هستند. هر دو الگوریتم مبتنی بر تکنیک جایگزینی زیررشته‌های تکراری با منابع هستند، اما از رویکردهای متفاوتی استفاده می‌کنند.

الگوریتم DEFLATE ترکیبی از کدگذاری دیکشنری (LZ77) (Ziv and Lempel, 1977) و کدگذاری هافمن است. DEFLATE در چندین برنامه فشرده‌سازی مانند Gzip, zip و WinRar استفاده می‌شود. باوجوداینکه در زمینه فشرده‌سازی داده‌ها تحقیقات گسترده‌ای انجام شده است، فناوری‌های جدید همچنان در حال ظهور هستند (Farina, et.al., 2012)، مانند شوکو، که در سال ۲۰۱۴ منتشر شد. شوکو یک برنامه فشرده‌سازی متن است که به‌ویژه بر روی متن ASCII مؤثر است. ASCII یک کدگذاری برای متن ساده است؛ یعنی متنی که از تنوع زیادی از کاراکترها و نمادها استفاده نمی‌کند. رمزگذاری کاراکتر یک متن به کامپیوتر می‌گوید چگونه داده‌ها باید خوانده شوند تا انسان بتواند آن را درک کند (Indurkha, et.al., 2010).

## ۲-۱ تعریف مسئله و بیان سؤالات اصلی تحقیق

انتقال حجم زیادی از داده‌ها از طریق اینترنت کار زمان بری است. حتی در شرایطی که یک کاربر دارای پهنای باند فوق‌العاده باشد، انتقال مگابایت یا گیگابایت داده بازهم چند دقیقه طول می‌کشد. این مشکل اغلب مربوط به فایل‌های باینری مانند فیلم‌ها است، اما ممکن است در مورد داده‌های مبتنی بر متن نیز صادق باشد. یک مثال می‌تواند شبیه‌سازی مخزن گیت باشد، هسته لینوکس مخزن گیت هاب نزدیک به سه گیگابایت حجم دارد. داده‌های متنی را می‌توان با استفاده از الگوریتم‌هایی مانند DEFLATE فشرده کرد، که توسط چندین برنامه فشرده‌سازی محبوب استفاده می‌شود. فشرده‌سازی داده‌ها حجم فایل را کاهش می‌دهد، پس انتقال آن از طریق اینترنت سریع‌تر خواهد بود الگوریتم DEFLATE برای کاربر روی انواع متن‌ها و داده‌های باینری طراحی شده است (Sayood, 2006). بااین‌حال یک الگوریتم فشرده‌سازی عمومی است که کار با آن می‌تواند بهینه شود. به‌عنوان مثال متن ASCII اکثر برنامه‌های فشرده‌سازی عمومی هستند و از الگوریتم‌های بدون اتلاف استفاده می‌کنند که به فشرده‌سازی هر نوع داده مانند متن، تصویر فایل‌های ویدئویی منجر می‌شود. چندین الگوریتم فشرده‌سازی اتلافی وجود دارد که به‌طور خاص برای هر نوع داده استفاده می‌شوند، اما تعداد آن‌ها کمتر از الگوریتم‌های بدون اتلاف است. این امکان وجود دارد که از الگوریتم‌های بدون اتلاف سفارشی برای فشرده‌سازی انواع داده‌های خاص، مانند متن ASCII، نسبت به الگوریتم‌های فشرده‌سازی عمومی، بیشتر استفاده شود.

پایان‌نامه به سؤالات تحقیق زیر پاسخ می‌دهد:

- چه ترکیبی از تکنیک‌ها می‌تواند نسبت فشرده‌سازی را بهبود بخشد؟
- آیا استفاده از ترکیب کدگذاری دیکشنری و کدگذاری هافمن برای فشرده‌سازی ASCII، نسبت به فشرده‌سازی DEFLATE فراتر می‌رود؟
- آیا این ترکیب فشرده‌سازی متن، فشرده‌سازی بالاتری نسبت به برنامه‌های فشرده‌سازی عمومی برای متن ASCII به دست می‌آورد؟
- آیا این ترکیب فشرده‌سازی متن، فشرده‌سازی بالاتری نسبت به برنامه‌های فشرده‌سازی که مخصوص متن ASCII هستند، به دست می‌آورد؟

### ۳-۱ انگیزه و ضرورت تحقیق

راه‌حل فشرده‌سازی متن پیشنهادی (که در اینجا به آن اشاره می‌شود TCS یا "ترکیبی") می‌تواند برای هرکسی که با مقدار زیادی (عمدتاً) متن اسکی سروکار دارد مفید باشد. این راه‌حل به گونه‌ای طراحی شده است که روی هر فرم متن کار کند، اما برای متن اسکی سفارشی شده است. این راه‌حل حجم فایل را کاهش می‌دهد تا فضای زیادی را اشغال نکند و همچنین هنگام بارگیری یا انتقال یک فایل فشرده از طریق اینترنت، پهنای باند را کاهش می‌دهد. کدگذاری دیکشنری (Langiu, 2013) و کدگذاری هافمن از قبل به طور گسترده‌ای در الگوریتم‌های فشرده‌سازی استفاده می‌شود (هر دو در DEFLATE)، اما الگوریتم‌های فشرده‌سازی کمی از روش فشرده‌سازی اسکی استفاده می‌کنند. بنابراین این ترکیب ممکن است فشرده‌سازی بهتری برای متن اسکی نسبت به الگوریتم‌های فشرده‌سازی عمومی باشد.

### ۴-۱ محدودیت‌ها

این پروژه به طور خاص برای فشرده‌سازی متن، و نه داده‌های باینری، مانند فایل‌های تصویری یا ویدیویی در نظر گرفته شده است. بنابراین آزمایش‌های انجام شده در این پروژه تنها روی داده‌های متنی می‌باشد. انتظار می‌رود که از ترکیب برای فشرده‌سازی مقادیر زیادی از متن استفاده شود، زیرا داده‌های متنی معمولاً فضای ذخیره‌سازی زیادی را در مقایسه با داده‌های ویدیویی یا یک فایل متنی کوچک که کاربرد چندانی ندارد؛ اشغال نمی‌کند. بنابراین فقط روی فایل‌های متنی با اندازه قابل توجه آزمایش می‌شود. اثربخشی راه‌حل (تأخیر) عاملی در ارزیابی نیست، زیرا

ترکیب فقط بر اساس نسبت فشرده‌سازی آن ارزیابی می‌شود، نه زمان اجرا. پایتون به‌عنوان زبان برنامه‌نویسی برای مازول فشرده‌سازی (ACM) ASCII انتخاب شده است تا کد منبع را ساده‌تر کند. پیاده‌سازی برنامه در C می‌توانست راه‌حل را سریع‌تر کند (Li Jun and Li Ling, 2010)، اما قابلیت فشرده‌سازی برای ACM اولویت بالاتری نسبت به تأخیر دارد (Salib, 2004). این پروژه پیاده‌سازی‌های مختلفی را برای استفاده در TCS ارزیابی کرد. پیاده‌سازی‌ها بر اساس میانگین نتایج حاصل از آزمایشات ارزیابی می‌شوند. این بدان معناست که موارد استفاده در آزمایش اهمیت برابر دارند. این پروژه با هیچ معضل اخلاقی مواجه نمی‌شود. این پروژه همکاری با یک شرکت نیست، بنابراین هیچ مشکلی با کپی‌رایت (حق چاپ) آن وجود نخواهد داشت. این پروژه تنها از کد منبع با مجوز باز استفاده می‌کند. آزمایشات روی فایل‌های ویکی‌پدیا فقط از یک زیرمجموعه از متن برای آزمایش فشرده‌سازی آن‌ها استفاده می‌کند، زیرا استفاده از چنین مجموعه داده بزرگی برای آزمایش غیرممکن است. زیرمجموعه به‌اندازه کافی بزرگ است که نتایج واضحی را ارائه می‌دهد.

## ۵-۱ پیشنهادها و هدف

اکثر برنامه‌های فشرده‌سازی از الگوریتم DEFLATE استفاده می‌کنند، که شامل کدگذاری دیکشنری و کدگذاری هافمن است. در این پروژه هم از این موارد استفاده خواهد شد زیرا مؤثر بودن این الگوریتم‌های استاندارد عملاً اثبات شده است (Fowers, et.al., 2015). الگوریتم‌های زیر در این پروژه استفاده می‌شود:

- کدگذاری دیکشنری؛ به‌ویژه در متن نوشتاری<sup>۱</sup> مؤثر است، مانند دایرةالمعارف (Shi and Sun, 1999).
- فشرده‌سازی ASCII؛ بر روی متن نوشتاری، به‌ویژه جایی که از کاراکترهای خاص بیشتر از بقیه استفاده می‌شود، مؤثر است (Trost, 2020).
- کدگذاری هافمن؛ تکنیکی که به‌طور گسترده در فشرده‌سازی داده‌ها استفاده شده است (Aggarwal, Narayan, 2000 و Klein, 2003).

این پروژه با استفاده از این سه روش، راه‌حلی را ایجاد می‌کند که برای فشرده‌سازی متن اسکی مطلوب است و سعی می‌کند نسبت فشرده‌سازی بیشتری برای متن اسکی نسبت به یک الگوریتم فشرده‌سازی عمومی داشته باشد.

---

<sup>۱</sup> "متن نوشتاری" در پایان نامه مورد بررسی به عنوان متنی که به برخی از زبان‌ها مانند یک زبان گفتاری یا یک زبان برنامه‌نویسی نوشته شده است تعریف شده است.

## ۶-۱ روش و مراحل انجام تحقیق

تحقیقات آماری که با داده‌های عددی سروکار دارند، کم هستند، درحالی‌که تحقیقات کیفی، تحقیقاتی و معمولاً غیر آماری است. آزمایشات برای سیستم هنگام مقایسه TCS با داده‌های عددی استفاده می‌شود، بنابراین تحقیقات کمی اندازه‌گیری‌های بهتری را در مقایسه با تحقیقات کیفی ارائه می‌دهد (Merriam, 2020). در مقابل آن برای نتیجه‌گیری از آزمایشات در تحقیقات کیفی به استدلال منطقی نیاز است که این استدلال می‌تواند استدلال قیاسی یا استدلال استقرایی باشد. استدلال قیاسی تشکیل یک نتیجه‌گیری بر اساس اظهارات یا حقایق عمومی پذیرفته‌شده است. استدلال استقرایی از منطق زیر پیروی می‌کند: اگر  $A = B$  و  $B = C$  سپس  $A = C$ . استدلال استقرایی شامل یک عنصر احتمال است و ممکن است به این نتیجه برسد که  $A$  احتمالاً برابر با  $B$  است (Merriam, 2020).

در پروژه موردبررسی تعدادی از تکنیک‌های فشرده‌سازی بر روی متون مشابه آزمایش می‌شوند، و نسبت فشرده‌سازی برای پیدا کردن تکنیک بیشترین فشرده‌سازی را به دست می‌آورد. آزمایشات داده‌های قطعی در مورد این‌که کدام راه‌حل بهترین نسبت فشرده‌سازی را به دست می‌آورد، ارائه می‌دهد و بنابراین نتیجه‌گیری با استدلال استقرایی یا قیاسی صورت می‌گیرد.

## ۷-۱ طرح کلی پایان‌نامه موردبررسی

فصل دوم این پایان‌نامه توضیح می‌دهد که فشرده‌سازی متن چیست و سابقه و کار مرتبط برای سه ماژول TCS را ارائه می‌دهد. کدگذاری دیکشنری، فشرده‌سازی اسکی و کدگذاری هافمن. فصل سوم روش‌های مختلفی که در تحقیقات دانشگاهی برای این پروژه استفاده می‌شود، و دلیل انتخاب این روش‌ها را، ارائه می‌دهد. فصل چهارم طراحی TCS و الزامات راه‌حل را ارائه می‌دهد. فصل پنجم فرایند ارزیابی TCS و ماژول‌های آن را شرح می‌دهد و نتایج حاصل از ارزیابی‌ها را ارائه می‌دهد. فصل ششم بحث نتایج حاصل از ارزیابی‌ها را پوشش می‌دهد و فصل هفتم نتیجه‌گیری و کار آینده پایان‌نامه را ارائه می‌دهد.

## ۸-۱ ساختار گزارش تحقیق

فصل اول به تعریف و مقدمه و دلایل نیاز به طرح ارائه‌شده پرداخته می‌شود.

فصل دوم به مفاهیم عمومی فشرده‌سازی پرداخته می‌شود.

فصل سوم مروری است بر کارهای انجام شده طرح پیشنهادی پایان نامه

فصل چهارم به ارزیابی و پیاده سازی الگوریتم ها و بررسی نتایج پرداخته می شود.

فصل پنجم به جمع بندی و نتیجه گیری پرداخته می شود.

## فصل دوم

### مفهوم فشردسازی داده‌ها و برخی از انواع آن

#### ۱-۲ مقدمه

فشردسازی داده‌ها فرایند کاهش تعداد بیت‌های موردنیاز برای نمایش رسانه است (Sayood, 2006). این بدان معناست که یک فایل فشرده کوچک‌تر از فایل اصلی و می‌تواند به‌عنوان مثال در مدت‌زمان کوتاه‌تری از طریق اینترنت منتقل شود. دو نوع تکنیک فشردسازی داده‌ها وجود دارد: اتلافی و بدون اتلاف. فشردسازی بدون اتلاف به این معنی است که می‌توان تمام داده‌های هنگام فشردسازی را، بازیابی کرد. این بدان معناست که پس از بازیابی فایل فشردشده با فایل اولیه کاملاً یکسان خواهد بود و بدون اتلاف می‌باشد. از طرف دیگر، فشردسازی ازدست‌رفته (اتلافی)، برخی از داده‌ها را در حین فرایند فشردسازی حذف می‌کند (Shukla and Prasad, 2011). هنگامی که فایل از حالت فشرده خارج می‌شود، داده‌های ازدست‌رفته قابل بازیابی نیستند. فشردسازی اتلافی بیشتر در مورد داده‌های دودویی استفاده می‌شود. فشردسازی اتلافی برای داده‌های دودویی در ترکیب با فشردسازی بدون اتلاف به کار می‌رود (Shukla and Prasad, 2011). به‌عنوان مثال فایل‌های صوتی، تصویری و ویدئویی (Sayood, 2006). مزیت فشردسازی اتلافی این است که معمولاً نسبت فشردسازی بسیار بالاتری نسبت به فشردسازی بدون اتلاف به دست می‌آورد، درحالی‌که کاربران هنوز می‌توانند اطلاعات داده‌های فشرده را درک کنند (Shukla and Prasad, 2011). اغلب، داده‌های ازدست‌رفته در فشردسازی برای کاربر قابل‌توجه نیست (Kleiman and Red, 2007). کدگذاری هافمن و کدگذاری دیکشنری نمونه‌هایی از فشردسازی بدون اتلاف هستند. این الگوریتم‌ها در بسیاری از برنامه‌های فشردسازی مانند Zip، Gzip یا Winrar استفاده می‌شوند، آنچه برای این الگوریتم‌ها متداول است این است که آن‌ها بالاترین نسبت فشردسازی را از فایل‌هایی که داده‌ها را تکرار می‌کنند یا توزیع ناهموار کاراکترها / نمادها را دارند، دریافت می‌کنند. این بدان معناست که داده‌های تصادفی یا یک متن "گیج‌کننده" ممکن است نسبت فشردسازی کمی داشته باشد. نسبت فشردسازی اندازه‌گیری میزان حجم فایل فشردشده در مقایسه با فایل اصلی است (Togneri and Desilva, 2002). به‌عنوان مثال، این الگوریتم‌ها هنگام فشردسازی متن نوشته‌شده به برخی از زبان‌ها مانند انگلیسی، آلمانی یا زبان برنامه‌نویسی مؤثر هستند؛ به این دلیل که برخی از کلمات، عبارات و کاراکترها بیشتر از بقیه استفاده می‌شوند (Togneri and Desilva, 2002). این

الگوریتم‌ها فقط در متن استفاده نمی‌شوند. فرمت PNG از این تکنیک‌ها برای یافتن مطابقت و الگوهای موجود در تصاویر استفاده می‌کند.

## ۲-۲ کدگذاری دیکشنری

کدگذاری یک تکنیک فشرده‌سازی است که رشته‌های تکراری متن را با منابع جایگزین می‌کند. درحالی‌که کدگذاری هافمن به کاراکترهای فردی و میزان استفاده آن‌ها در متن نگاه می‌کند، کدگذاری دیکشنری از رشته‌های تکراری استفاده می‌کند (Shi and Sun, 1999). به‌عنوان مثال، اگر متنی چندین بار حاوی کلمه "فشرده‌سازی" است، موارد تکراری کلمه را می‌توان با اشاره به کلمه در دیکشنری جایگزین کرد. هنگامی‌که متن از حالت فشرده خارج می‌شود، الگوریتم از مرجع برای جستجوی کلمه استفاده می‌کند و مرجع را با کلمه اصلی جایگزین می‌کند. دیکشنری لیستی از رشته‌های متنی است که اغلب مورد استفاده قرار می‌گیرد (Huffman, 1997). برخی از رمزگذارهای دیکشنری، مانند LZ77 از یک دیکشنری صریح استفاده نمی‌کنند (Ziv and Lempel, 1977)، اما در عوض از ارجاعات انجام‌شده قبلی در متن استفاده می‌کند (Hankerson, et. Al., 2003). رمزگذاری‌های دیکشنری نیز می‌توانند در داده‌های غیر متنی مانند تصاویر استفاده شوند. فرمت تصویر PNG از کدگذاری دیکشنری LZ77 استفاده می‌کند. در مورد تصاویر، الگوریتم باید دنباله‌های پیکسلی را به‌جای رشته‌های متنی تکرار کند.

مزیت استفاده از کدگذار دیکشنری این است که می‌تواند نسبت فشرده‌سازی بسیار بالایی را برای فایل‌های خاص به دست آورد. کدگذاری‌های دیکشنری به‌ویژه بر روی داده‌های تکراری، مانند متن نوشتاری مؤثر هستند (Shi and Sun, 1999). متن، که به زبان نوشته شده است، معمولاً برخی از کلمات یا ترکیب حروف را بیشتر از بقیه تکرار می‌کند که برنامه نویسان دیکشنری از این مزیت استفاده می‌کنند. عیب استفاده از برخی کدگذارهای دیکشنری، مانند LZW، این است که فایل فشرده ممکن است بزرگ‌تر از فایل فشرده نشده باشد (Senthil and Robert, 2011) و (Taleb, et.al., 2010). این اتفاق در صورتی می‌افتد که داده‌هایی که باید فشرده شوند هیچ اطلاعات تکراری نداشته باشند. اگر فشرده‌سازی یک فایل باعث کاهش حجم آن نشود، پس هدف شکست خورده است.

### ۲-۲-۱ الگوریتم‌های رایج کدگذاری دیکشنری

این بخش ویژگی‌ها و تفاوت‌های اصلی بین الگوریتم‌های کدگذاری دیکشنری را توصیف می‌کند. ویژگی‌های یک الگوریتم می‌تواند راه‌حل فنی آن یا چیزی باشد که آن را از بقیه متمایز کند. یک الگوریتم همچنین می‌تواند به این معنا باشد که کمابیش برای کارهای فشرده‌سازی مناسب است.



### الگوریتم LZ77

LZ77 یک کدگذار دیکشنری محبوب است که در بین سایر موارد، در فرمت تصویر PNG و الگوریتم DEFLATE استفاده می‌شود. الگوریتم LZ77 هنگام فشرده‌سازی داده‌ها از پنجره لغزان<sup>۱</sup> (Ziv and Lempel, 1977) استفاده می‌کند (به عنوان مثال یک فایل تصویری یا یک فایل متنی). sliding window یک بافر است که میزان متن (یا داده) را که الگوریتم در هر نقطه تجزیه و تحلیل می‌کند تا زیررشته‌های منطبق را پیدا کند را تعیین می‌کند. پنجره لغزان شامل یک بافر پیش‌بینی و یک بافر جستجو/است. بافر پیش‌بینی، متنی را که هنوز فشرده یا کدگذاری نشده است، تجزیه و تحلیل می‌کند، به عنوان مثال ۲۰ کاراکتر بعدی، در حالی که بافر جستجو متنی است که اخیراً کدگذاری شده است و معمولاً بسیار بزرگ‌تر از بافر پیش‌بینی است. اگر یک زیررشته در بافر جستجو باشد، نیز در بافر پیش رو ظاهر شود، سپس زیررشته در بافر پیش رو با اشاره به بافر قبلی جایگزین می‌شود (Ziv and Lempel, 1977).

### الگوریتم LZ78

در حالی که الگوریتم LZ77 با اشاره به یک نقطه قبلی در متن، زیررشته‌های مکرر ایجاد می‌کند، LZ78 از یک دیکشنری صریح استفاده می‌کند. الگوریتم LZ78 یک دیکشنری از زیررشته‌ها ایجاد می‌کند و زیررشته‌ها را در یک متن با اشاره به دیکشنری جایگزین می‌کند (Ziv and Lempel, 1978). الگوریتم LZ78 تجدیدنظر در الگوریتم LZ77 است، و به جز دیکشنری صریح در بقیه موارد مشابه‌اند (Atwal and Bansal, 2015).

### الگوریتم LZW

الگوریتم LZW یک نوع LZ78 است و همچنین از یک دیکشنری صریح استفاده می‌کند (Hankerson, et.al., ۲۰۰۳). فرض اساسی LZW این است که فرایند کدگذاری را با ۲۵۶ ورودی‌ها/کاراکتر در دیکشنری شروع می‌کند (Taleb, 2010). این معمولاً مجموعه کاراکتر اسکی است، که از یک بایت در هر کاراکتر استفاده می‌کند. این الگوریتم دائماً زیررشته‌هایی را به دیکشنری با مرجع بالاتر از ۲۵۶ اضافه می‌کند. این ورودی‌ها ممکن است در متن تکرار شوند یا اینکه تکرار نشوند. اگر یک زیررشته تکرار شود، از همان اشاره به زیررشته استفاده می‌شود (Hankerson, et.al., 2003).

### الگوریتم Re-pair

الگوریتم Re-pair از یک دیکشنری صریح استفاده می‌کند و شبیه الگوریتم LZ78 است. تفاوت Re-pair با LZ78 این است که Re-pair فقط زیررشته‌های دو رشته‌ای را ذخیره می‌کند (Ziv and Lempel, 1978). Re-pair همه

---

<sup>۱</sup>sliding window

جفت‌های کاراکتر را که بیش از یک‌بار در یک متن وجود دارد، می‌یابد و با اشاره به دیکشنری آن‌ها را جایگزین می‌کند. این امر تا زمانی ادامه می‌یابد که هیچ جفت کاراکتری که بیش از یک‌بار در متن ظاهر شود، وجود نداشته باشد (Larsson and Moffat, 2020).

## ۲-۳ فشرده‌سازی ASCII

فشرده‌سازی ASCII در پایان‌نامه موردبررسی به‌عنوان الگوریتم‌های فشرده‌سازی تعریف‌شده است. به‌ویژه بر روی متن اسکی مؤثر است. کدگذاری اسکی دارای تغییرات زیادی برای زبان‌های مختلف است، اما پایان‌نامه موردبحث بر US-ASCII تمرکز می‌کند، مگر این‌که خلاف آن بیان شود. برای ارتباطات اینترنت US-ASCII به کدگذاری اسکی ارجحیت دارد. امروزه US-ASCII با ASCII کاربرد یکسانی دارد (Arregoces and Portolani, 2004) و (Carlsson and Miller, 2012). رمزگذاری ASCII یک کدگذاری متنی است که به‌طور گسترده برای متنی که کاراکترها و نمادهای خیلی متنوع ندارد، استفاده می‌شود (Arregoces and Portolani, 2004) و (Carlsson and Miller, 2012). این محدودیت در تعداد نمادها از نحوه ذخیره متن بر روی کامپیوتر نشأت می‌گیرد. یک نماد دقیقاً یک بایت (۸ بیت) برای ذخیره‌سازی استفاده می‌کند، که در آن ۷ بیت کد نماد و ۱ بیت "مهم‌ترین بیت" است. هدف از مهم‌ترین بیت می‌تواند به‌عنوان یک بررسی برابری، برای تشخیص خطاها در داده‌ها استفاده شود (Korpela, 2006). همچنین گروهی از رمزگذاری‌ها معروف به "اسکی توسعه‌یافته" از ۸ بیت برای نماد استفاده می‌کنند و در نتیجه مهم‌ترین بیت را کنار می‌گذارند. رمزگذاری‌هایی که از همه ۸ بیت استفاده می‌کنند، به‌جای ۱۲۸ نماد استاندارد از ۲۵۶ نماد استفاده می‌کنند (Ediger, 2008).

### ۱-۲-۳ الگوریتم شوکو (Shoco)

شوکو یک الگوریتم فشرده‌سازی متن است که توسط Christian Schramm توسعه‌یافته است، که به‌عنوان "کمپرسور سریع برای رشته‌های کوتاه" توصیف می‌شود و به‌ویژه بر روی متن اسکی مؤثر است. راه‌حل فشرده‌سازی شوکو از این واقعیت استفاده می‌کند که در هر زبان، برخی از کاراکترها بیشتر از بقیه استفاده می‌شود. همان‌طور که قبلاً ذکر شد، اولین بیت در یک کاراکتر اسکی اضافی است، مگر اینکه برای اهدافی مانند تشخیص خطا استفاده شود. بنابراین

الگوریتم از اولین بیت در یک بایت استفاده می‌کند تا نشان دهد که آیا بیت‌های زیر به یک کاراکتر مشترک اشاره می‌کنند یا خیر. اگر بیت اول روی ۰ تنظیم شود، بدین معنی که ۷ بیت زیر نشان‌دهنده یک کاراکتر غیرمعمول است؛ پس از فشرده شدن کاراکتر اسکی تغییر نخواهد کرد. اگر بیت اول روی ۱ تنظیم شود، بدین معنی که بیت‌های زیر دو کاراکتر مشترک را نشان می‌دهند. سپس دو کاراکتر با ۳ یا ۴ بیت نشان داده می‌شوند. همه باهم از ۸ بیت استفاده می‌کنند که به‌عنوان یک بایت ذخیره می‌شوند.

شوکو همچنین می‌تواند کاراکترهای غیر اسکی را فشرده کند، اما هزینه دارد. "اگر رشته ورودی شما به‌طور کامل (یا بیشتر) اسکی نباشد، خروجی ممکن است افزایش یابد". وقتی الگوریتم با کاراکتری مواجه می‌شود که بیش از ۷ بیت استفاده می‌کند (مانند یک کاراکتر UTF-8) استفاده می‌کند، درست قبل از کاراکتر "نشانگر" را وارد می‌کند. نشانگر یک کاراکتر خاص است که یک بایت را اشغال می‌کند. هدف آن این است که نشان دهد که کاراکتر بعدی اسکی نیست و بنابراین بیش از ۷ بیت استفاده خواهد کرد. نشانگر همچنین می‌گوید که کاراکتر بعدی از چند بایت استفاده می‌کند (به‌عنوان مثال ۱ یا ۲ بایت).

تکنیکی که شوکو از آن استفاده می‌کند شمارش تکرار (فراوانی) بیگرم در یک متن است. بیگرام‌ها، درزمینه شوکو، جفت‌های منحصربه‌فرد از دو کاراکتر متوالی در یک متن است. یک بیگرام رایج در انگلیسی "qu" است، زیرا بعد از "q" تقریباً همیشه "u" می‌آید. بعد از اینکه شوکو لیستی از متداول‌ترین کاراکترهای یک متن تهیه کرد، لیستی دیگر از کاراکترهایی که به‌احتمال زیاد به دنبال کاراکترهای رایج می‌آیند، تهیه می‌کند. اگر شوکو متوجه شد که به‌عنوان مثال "he" یک بیگرام رایج است، آنگاه می‌تواند تمام کلمات حاوی این بیگرام، مانند "the"، "she" و "then" را فشرده کند.

مزیت استفاده از شوکو این است که تا زمانی که ورودی صد درصد اسکی باشد، فایل فشرده هرگز بزرگ‌تر از فایل فشرده نشده نخواهد بود. علاوه بر این، شوکو را می‌توان به‌راحتی در ترکیب با سایر الگوریتم‌های فشرده‌سازی استفاده کرد، مانند الگوریتم‌های کدگذاری دیکشنری، یا کدگذاری هافمن برای دستیابی به فشرده‌سازی بیشتر. همچنین یک الگوریتم شوکو بسیار سریع است. به گزارش وب‌سایت، شوکو هنگام فشرده‌سازی یک فایل ۴,۹ مگابایتی تقریباً ۷ برابر سریع‌تر از Gzip است. نقطه‌ضعف شوکو این است که استفاده از شوکو به‌تنهایی نسبت فشرده‌سازی بسیار بدتری از برنامه‌های فشرده‌سازی استاندارد دارد. علاوه بر این، کاربر باید بداند که چه نوع متنی را فشرده می‌کند، اگر متن شامل مقدار زیادی از کاراکترهای چندبایتی، مانند UTF, 8 باشد، ممکن است فایل فشرده بزرگ‌تر از فایل فشرده نشده باشد، که هنگام استفاده از ابزار فشرده‌سازی بسیار نامطلوب است.

## ۴-۲ کدگذاری هافمن

کدگذاری هافمن به گونه‌ای طراحی شده است که پرکاربردترین کاراکترها را در یک متن تا حد امکان بیت، نشان دهد. اول، الگوریتم ورودی را با پیمایش متن کاراکتر به کاراکتر، تجزیه و تحلیل می‌کند. سپس الگوریتم لیستی از هر کاراکتر که در متن استفاده می‌شود، ایجاد کرده و لیست را بر اساس تکرارها مرتب می‌کند. پس از ایجاد لیست، الگوریتم یک درخت دودویی بر اساس جدول می‌سازد. درخت دودویی یک ساختار داده است که از "گره" و "شاخه" تشکیل شده است (Pai, 2008). هر گره می‌تواند حداکثر دو "گره فرزند" داشته باشد. در شکل ۱، گره‌ها مستطیل‌های آبی‌رنگ هستند. "گره‌های برگ" آخرین گره‌های یک شاخه هستند. یک درخت دودویی را به عنوان یک درخت در نظر بگیرید که وارونه است. همه شاخه‌های یک درخت از یک ریشه شروع می‌شوند، و برگ‌ها در انتهای شاخه‌ها هستند. هر گره برگ ارزش یک کاراکتر مورد استفاده در متن را دارد.

هنگامی که یک درخت هافمن ساخته می‌شود، فشرده‌سازی متن شروع می‌شود. این الگوریتم درخت را برای هر کاراکتر در یک متن پیمایش می‌کند. برای رسیدن به یک گره برگ، الگوریتم به سمت گره فرزند چپ یا راست حرکت می‌کند تا کاراکتر مشخص شده، پیدا می‌شود. این پیمایش به صورت دوتایی با یک و صفر نشان داده شده است، جایی که ۱ است، به معنای "رفتن به راست" و ۰ به معنای "رفتن به سمت چپ" است. به منظور فشرده‌سازی یک متن برای بازیابی داده‌های اصلی، درخت هافمن نیز باید ذخیره شود. به همین دلیل است که کدگذاری هافمن در متون بزرگ‌تر مؤثرتر است. به عنوان مثال وقتی یک متن بزرگ وجود دارد که فقط شامل ۲۶ حرف الفبای انگلیسی و چند علامت نقطه گذاری است، درخت صرف نظر از اندازه متن، بیش از ۲۶ گره برگ ندارد.

کدگذاری هافمن یک رمزگذار آنتروپی است، به این معنی که داده‌ها را بر اساس تکرار نماد، فشرده می‌کند؛ که این برخلاف کدگذاری دیکشنری است که به تکرارهای زیررشته‌ها یا توالی نمادها نگاه می‌کند. برنامه‌های آنتروپی مانند کدگذاری هافمن هنگامی سودمند است که از داده‌های خاص بیشتر از بقیه نمادها استفاده شود. کدگذاری هافمن می‌تواند به آسانی هم در داده متنی و هم در داده‌های دودویی مورد استفاده قرار گیرد و می‌تواند در ترکیب با سایر الگوریتم‌های فشرده‌سازی مانند کدگذاری دیکشنری استفاده شود (Desoky and Gregory, 1988). کدگذاری هافمن اغلب با کدگذاری حسابی مقایسه می‌شود، که نوع دیگری از رمزگذار آنتروپی است. نتایج چنین مقایسه‌هایی نشان می‌دهد که کدگذاری هافمن سریع‌تر از کدگذاری حسابی است، اما کدگذاری حسابی به طور کلی نسبت فشرده‌سازی بهتری دارد (Bookstein and Klein, 1993) و (Duda, 2014).

### ۴-۲-۱ فرایند فشرده سازی

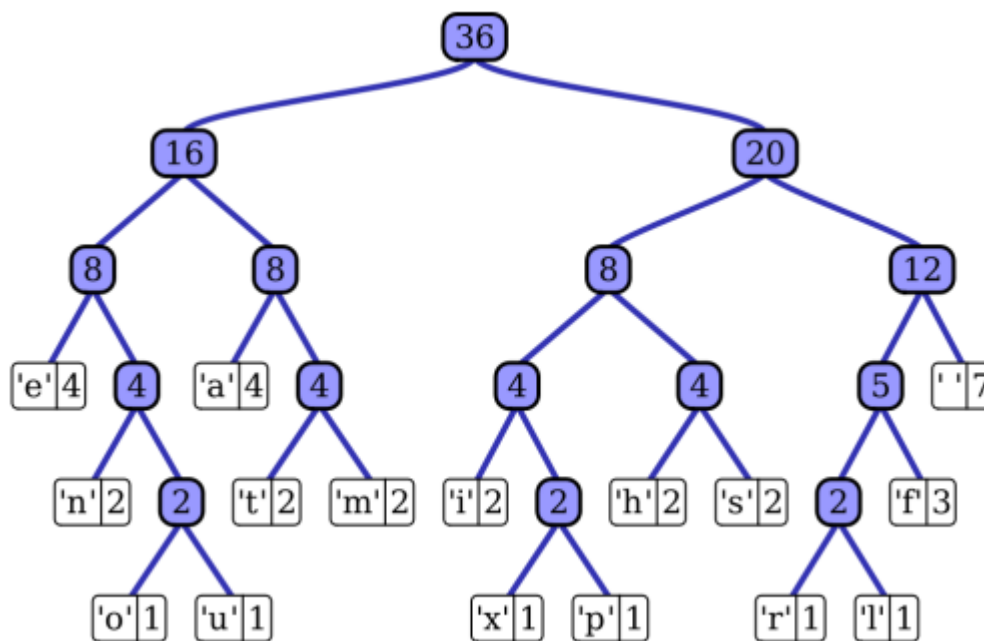
این بخش فرایند فشردگی برای کدگذاری هافمن را توضیح می‌دهد. یک مثال ورودی برای الگوریتم می‌تواند

"this is an example of a huffman tree" جدول ۱-۲ فراوانی از هر کاراکتر در جمله مثال است.

جدول ۱-۲. مثال فراوانی کاراکتر

کاراکتر	Space	A	E	F	H	I	M	N	S	T	L	O	P	R	U	X
فراوانی	۷	۴	۴	۳	۲	۲	۲	۲	۲	۲	۱	۱	۱	۱	۱	۱

از این جدول می‌بینیم که کاراکتر space بیشترین فراوانی را دارد، پس از آن "a" و "e" با چهار مورد. شکل ۱-۲ نشان می‌دهد که چگونه جمله مثال شبیه درخت دودویی هافمن است.



شکل ۱-۲. درخت دودویی هافمن

گره‌های این شکل همه دارای اعدادی هستند که در گره‌های فرزند نشان می‌دهد چند بار از هر کاراکتر استفاده شده است. این عدد هنگام ساختن درختان هافمن مهم است زیرا کاراکترهایی که بیشتر تکرار می‌شوند باید تا حد ممکن بالای درخت باشند (نزدیک به گره ریشه). جدول ۱-۲ نشان می‌دهد که فراوانی (تکرار) کاراکتر "a" و "e" بیشترین هستند و بنابراین به گره ریشه نزدیک هستند. این سه کاراکتر را می‌توان پس از سه عبور از گره ریشه به دست آورد.

جمله مثال "this is an example of a huffman tree"، که با کدگذاری اسکی ذخیره می‌شود، از ۳۶ بایت ذخیره‌سازی استفاده می‌کند. این اندازه فایل قبل از فشرده‌سازی است، جایی که هر کاراکتر از یک بایت (یا ۸ بیت) ذخیره‌سازی استفاده می‌کند. هنگامی که فشرده‌سازی شروع می‌شود، با اولین کاراکتر "t" شروع می‌شود، برای رسیدن به گره برگ که مقدار "t" را از گره ریشه نگه می‌دارد، الگوریتم ابتدا به گره فرزند چپ، سپس دو بار به گره فرزند راست و در نهایت به چپ برود. بنابراین نمایش "t" فشرده دودویی (باینری) "۰۱۱۰" خواهد بود، که فقط ۴ بیت است. الگوریتم این کار را برای هر کاراکتر در متن انجام می‌دهد و پس از اتمام، با یک جریان دودویی ۱۳۵ بیت به پایان می‌رسد. این جریان فقط از ۱۷ بایت حافظه در کامپیوتر استفاده می‌کند که کمتر از نیمی از متن فشرده نشده است. البته، درخت هافمن نیز باید برای فشرده‌سازی متن در آن ذخیره شود.

## ۵-۲ کارهای مرتبط و تاریخچه تحقیق

این بخش مقاله‌ها و راه‌حل‌های مربوط به فشرده‌سازی متن یا به‌طور خاص، فشرده‌سازی اسکی، را مورد بحث قرار می‌دهد، زیرا راه‌حل ارائه‌شده در پایان‌نامه مورد بررسی یک فشرده‌ساز متن است که برای فشرده‌سازی متن اسکی سفارشی شده است.

### ۱-۵-۲ تقویت فشرده‌سازی متن با رمزگذاری آماری مبتنی بر کلمه

این مقاله از نسخه ۲۰۱۲ مجله کامپیوتر ("The Computer Journal", 2020) است که توسط آنتونیو فاریا، گونزالو ناوارو و خوزه آر پاراما نوشته شده است. این مقاله پیشرفت احتمالی در برنامه‌های فشرده‌سازی را ارائه می‌دهد. هدف از راه‌حل پیشنهادی این است که نسبت فشرده‌سازی متن را افزایش داده و زمان فشرده‌سازی را برای فشرده‌سازهای متن کاهش دهد. این راه‌حل از تکنیک‌های فشرده‌سازی مبتنی بر کلمه و بایت جهت پیش‌پردازش برنامه‌های فشرده‌سازی عمومی استفاده می‌کند (Farina, 2012).

فشرده‌سازی مبتنی بر کلمه نوعی کدگذاری دیکشنری است که در آن الگوریتم کل کلمات را جستجو می‌کند (Moffat, 1989). اکثر برنامه نویسان دیکشنری، مانند الگوریتم‌های Lempel, Ziv، زیررشته‌های مکرر را جستجو می‌کنند (Ziv, Lempel, 1977). زیررشته‌ها می‌توانند هر طولی داشته باشند و می‌توانند بخشی از یک کلمه باشند. فشرده‌سازی مبتنی بر کلمه فقط کلمات کامل را به دیکشنری خود اضافه می‌کند. این مقاله بیان می‌کند که استفاده از فشرده‌سازی مبتنی بر کلمه، برخلاف کدگذاری استاندارد دیکشنری، زمان فشرده‌سازی و فشرده‌سازی بهتری

را ارائه می‌دهد. علاوه بر این، امکان جستجوی کلمات و عبارات روی فایل فشرده بدون نیاز به فشرده‌سازی آن وجود دارد (Farina, 2012). فشرده‌سازی بایت گرا دومین پیش‌پردازشی است که راه‌حل پیشنهادی از آن استفاده می‌کند. فشرده‌سازی بایت گرا تکنیکی است که هنگام نوشتن کدهای فشرده‌سازی فقط از بایت‌های کامل استفاده می‌کند و نه از بیت‌ها (De Moura, et. Al., 1998). برعکس، کدگذاری استاندارد هافمن یک تکنیک فشرده‌سازی است که از توالی بیت‌ها به‌عنوان کد برای حرکت در درخت هافمن استفاده می‌کند. بنابراین کدگذاری هافمن باید به‌عنوان دنباله‌ای از بیت‌ها خوانده شود تا بایت. که می‌تواند باعث افزایش تأخیر شود، زیرا داده‌ها به‌صورت بایت در کامپیوتر ذخیره و خوانده می‌شوند. راه‌حل پیشنهادی از یک الگوریتم بایت گرا به نام Tagged Huffman استفاده می‌کند (Silva de Moura, 2000). برچسب‌گذاری هافمن (Tagged Huffman) از تکنیک فشرده‌سازی مشابه کدگذاری استاندارد هافمن استفاده می‌کند، اما هر کد برای جستجوی درخت هافمن از تعداد مشخصی بایت استفاده می‌کند. بر اساس این مقاله، استفاده از Tagged Huffman به‌جای کدگذاری استاندارد هافمن، تأخیر را کاهش می‌دهد، اما نسبت فشرده‌سازی را در حدود ۰/۶ کاهش می‌دهد، که آن‌ها به‌عنوان یک معامله پذیرفتند (Farina, ۲۰۱۲).

نتایج آزمایشات نشان می‌دهد که یک برنامه فشرده‌سازی در فشرده‌سازی یک فایل از پیش‌پردازش شده، ۵ برابر سریع‌تر از فایل اصلی است. با این حال، زمان فشرده‌سازی پیش‌پردازنده را شامل نمی‌شود. مقایسه‌ها نشان می‌دهد استفاده از پیش‌پردازنده قبل از یک فشرده‌ساز عمومی (از جمله Gzip, Bzip2, 7-Zip, Re-pair)، فایل‌های فشرده نهایی بین ۵،۵ تا ۱۰ درصد کوچک‌تر از فایل فشرده‌شده با فشرده‌ساز عمومی بوده‌اند (Farina, 2012).

### مقایسه با TCS

فاریا (Fariña) و همکاران: راه‌حل (FS) و TCS هر دو فشرده‌ساز متنی هستند که نسبت فشرده‌سازی بالاتری را برای متن‌های نوشته‌شده به برخی از زبان‌ها به دست می‌آورد. روش فشرده‌سازی مبتنی بر کلمه به فشرده‌سازی FS می‌انجامد زیرا در زبان‌ها بیشتر از کلمات خاص استفاده می‌شود تا کلمات دیگر. همچنین TCS از ویژگی‌های زبان‌ها برای افزایش نسبت فشرده‌سازی استفاده می‌کند؛ ماژول فشرده‌سازی اسکی (ASCII compression module) (ACM)) از یک مدل فشرده‌سازی (توضیح داده‌شده در بخش ۱-۱۰-۳) بر اساس فراوانی حروف در زبان انگلیسی استفاده می‌کند، و ماژول کدگذاری دیکشنری زمانی فشرده‌سازی می‌کند که یک متن دارای ترکیبی از کاراکترها باشد. هدف FS افزایش نسبت فشرده‌سازی و کاهش تأخیر در برنامه‌های فشرده‌سازی است. هر دو روش مبتنی

بر کلمه و بایت جهت کاهش زمان فشردن سازی است (Farina, 2012). با این حال، TCS فقط نسبت فشردن سازی را در اولویت قرار می‌دهد، نه سرعت را.

## ۲-۵-۲ فشردن سازی آنلاین فایل‌های اسکی

این مقاله از کنفرانس بین‌المللی ۲۰۰۴ در زمینه فناوری اطلاعات: کدگذاری و محاسبه است که توسط جان ایستل، پاملا مندلباوم و اما رجنتوا (Istle, et. Al., 2004) نوشته شده است. در واقع راه‌حل پیشنهادی این است که در یک متن انگلیسی کدهای کوتاه‌تری به بیگرام‌های مکرر بدهید و در نتیجه به فشردن سازی برسید. این مقاله ادعا می‌کند که فشردن سازی باید رخ دهد زیرا ترکیبات خاصی از حروف بیشتر به زبان انگلیسی هستند نسبت به سایر زبان‌ها. این راه‌حل فقط بر روی متون دارای کد اسکی کار می‌کند و نه هیچ کدگذاری دیگری (Istle, et. Al., 2004). این الگوریتم با استفاده از ۲۸ دیکشنری استاتیک کار می‌کند. دیکشنری‌های استاتیک به این معنی است که با تجزیه و تحلیل یک متن خاص، برای یک مورد استفاده، ایجاد نشده‌اند، بلکه مجموعه‌ای از واژه‌نامه‌های عمومی هستند که باید فشردن شوند - مخصوصاً برای متون انگلیسی. یکی از دیکشنری‌ها برای اعداد، دیگری برای علائم نگارشی و کاراکترهای خاص و ۲۶ دیکشنری باقی مانده برای هر حرف الفبا استفاده می‌شود. اگر متنی که فشردن می‌شود حاوی حرف "a" باشد، الگوریتم، دیکشنری آن حرف خاص را جستجو می‌کند و پیدا می‌کند که حرف بعدی در چه شاخصی قرار دارد. به عنوان مثال، دیکشنری حرف "a" حرف "t" را به عنوان اولین فهرست دارد، زیرا "t" حرفی است که بیشتر از بقیه حروف به دنبال "a" می‌آید. اگر متن حاوی حرف "t" بود با شماره ۱ جایگزین می‌شد (Istle, et. Al., 2004).

برای هر کلمه جدید در یک متن، الگوریتم با جستجوی جایی که اولین حرف یا نماد در "دیکشنری پیش فرض" قرار دارد شروع می‌شود. دیکشنری پیش فرض بعد از علائم نگارشی یا کاراکترهای خاص استفاده می‌شود. اگر کلمه‌ای که قرار است فشردن شود "hat" باشد، "h" در دیکشنری پیش فرض فهرست ۱۳ خواهد بود، "a" در دیکشنری حرف "h" فهرست ۲ خواهد بود و حرف "t" در دیکشنری "a" فهرست ۱ خواهد بود. سپس اعداد ۱۳، ۲ و ۱ به دنباله‌های بیتی ترجمه می‌شوند (Istle, et. Al., 2004). دنباله‌های بیتی کوتاه‌تر برای هر کاراکتر کمتر از یک بایت استفاده می‌کنند، در مقابل کاراکترهای فشردن نشده اسکی که از بایت برای هر کاراکتر استفاده می‌کنند. جایگزینی کاراکترهای اسکی با دنباله‌های بیتی باعث فشردن سازی می‌شود.

آزمایش‌های انجام شده در مقاله نتایج حاصل از راه‌حل پیشنهادی را با برنامه فشردن سازی Winzip مبتنی بر DEFLATE مقایسه کرده است. یک فایل "کمتر از ۱۰۰ کیلوبایت" و "تا ۱ مگابایت" برای مقایسه استفاده شد



(Istle, et. Al., 2004). راه حل پیشنهادی نسبت فشرده سازی ۱,۶۷ را برای هر دو فایل به دست آورد، درحالی که WinZip نسبت فایل بزرگتر را ۲,۵ و نسبت فایل کوچکتر را کمتر از ۱ به دست آورد. از این نتایج، آن‌ها به این نتیجه رسیدند که "تکنیک دیکشنری چندگانه در فشرده سازی فایل‌ها در هر اندازه ای سازگار است" (Istle, et. Al., 2004).

### مقایسه با ماژول فشرده سازی ASCII

شباهت های زیادی بین روش دیکشنری چندگانه و ACM از TCS وجود دارد. هر دو فقط می تواند فایل های متنی را فشرده کنند، هر دو برای متن اسکی در نظر گرفته شده اند، و هر دو از دیکشنری های استاتیک که برای متون انگلیسی سفارشی شده اند استفاده می کنند. راه حل دیکشنری چندگانه از تکنیکی هم تراز با کدگذاری هافمن استفاده می کند. کدگذاری هافمن کدهای کوتاه تری به کاراکترهای مکرر می دهد. راه حل دیکشنری چندگانه، کدهای کوتاه تری را به بیگرم های تکراری می دهد، ACM کاراکترها را به کدها تبدیل نمی کند، بلکه از تکنیکی به نام ادغام استفاده می کند که در بخش ۱-۱۰-۳ توضیح داده می شود. تفاوت دیگر این است که ACM می تواند متن کدگذاری های مختلف را فشرده کند، درحالی که راه حل دیکشنری چندگانه نمی تواند. در مقاله آن‌ها این را به عنوان یک کار آینده بیان کردند (Istle, et. Al., 2004).

## فصل سوم

### مروری بر کارهای انجام شده

#### ۱-۳ مقدمه

در فصل ۳ پایان نامه مورد بررسی، روش ها و متدولوژی های مختلف را تشریح و مقایسه می کند و توضیح می دهد که کدام روش ها برای این پروژه انتخاب شده اند و چرا. همچنین روش های مختلف دانشگاهی را که برای پروژه های تحقیقاتی استفاده می شود و روش ها و مدل هایی را که بطور خاص برای توسعه نرم افزار استفاده می شود، ارائه می دهد و در فصل ۴ پایان نامه مورد بررسی، الزامات کلی سیستم و نیز الزامات مجموعه داده و پیاده سازی های مورد استفاده در TCS را شرح می دهد. سپس طراحی کلی TCS و راه حل فنی مازول فشرده سازی اسکی (ACM) را توضیح می دهد. در ادامه این مطالب را به تفکیک توضیح خواهم داد.

#### ۲-۳ روش ها و استراتژی های تحقیق

چندین روش تحقیق وجود دارد که می تواند برای پروژه های تحقیقاتی مورد استفاده قرار گیرد. روش های تحقیق روش هایی برای انجام وظایف تحقیق هستند و نحوه انجام تحقیق را توضیح می دهند. روش های رایج تحقیق شامل تجربی، توصیفی و بنیادی است. رویکرد تجربی علل، آثار و متغیرها را بررسی می کند. آزمایشات انجام شده و نتایج تجزیه و تحلیل می شوند. روش تحقیق توصیفی ویژگی های یک موقعیت را توصیف می کند، اما علل آن را توضیح نمی دهد. روش توصیفی اغلب از نظر سنجی یا مطالعات موردی استفاده می کند. روش تحقیق بنیادی ابتکاری است و ایده ها، اصول و نظریه های جدیدی را ایجاد می کند. تحقیقات بنیادی بر اساس کنجکاوی انجام می شود و در مورد مشاهده یک پدیده است (Hakansson, 2013). این پروژه از روش تحقیق تجربی استفاده می کند، زیرا آزمایش بخش مهمی از پروژه است. آزمایش روابط بین متغیرها و نحوه تغییر متغیرها بر نتایج را بررسی می کند.

استراتژی تحقیق، رویکرد مشخص تری برای نحوه انجام تحقیق است، در حالی که روش تحقیق چارچوبی برای تحقیق است. یک استراتژی تحقیق یک راهنما یا یک روش است. یک استراتژی تحقیق می تواند تجربی باشد، درست مانند روش تحقیق تجربی. رویکرد تجربی باهدف کنترل همه عواملی که ممکن است بر نتیجه یک آزمایش تأثیر بگذارد، کنترل می شود. این استراتژی بر اساس نتایج یک آزمایش، یک فرضیه را تأیید یا رد می کند. یکی دیگر از

استراتژی‌های رایج تحقیق استفاده از نظرسنجی یا پرسشنامه است. این استراتژی شامل جمع‌آوری اطلاعات از افراد است. این نظرسنجی‌ها می‌تواند عمیق باشد و ممکن است فقط شامل چند نفر باشد که روش کیفی است یا می‌توان نظرسنجی‌ها را برای تجزیه و تحلیل داده‌های تعداد بیشتری از افراد طراحی کرد که روش کمی است (Hakansson, 2013). این پروژه به همان دلیلی که از روش تحقیق تجربی استفاده می‌کند، از استراتژی تحقیق تجربی استفاده می‌کند.

### ۳-۳ جمع‌آوری داده‌ها

پایان‌نامه مورد بررسی از تحقیقات کمی استفاده می‌کند؛ مناسب‌ترین روش جمع‌آوری داده‌ها برای پروژه‌های تحقیقاتی کمی، آزمایش، پرسشنامه، مطالعات موردی و مشاهدات است. آزمایش‌ها مجموعه داده بزرگی را که برای متغیرها استفاده می‌شود جمع‌آوری می‌کند. پرسشنامه‌ها داده‌ها را از طریق سؤالات کمی یا کیفی جمع‌آوری می‌کند، مطالعات موردی تجزیه و تحلیل عمیق یک یا چند شرکت‌کننده است. مشاهدات رفتار را با تمرکز بر موقعیت‌ها و فرهنگ بررسی می‌کند (Hakansson, 2013). این پروژه از آزمایشات برای جمع‌آوری داده‌ها استفاده می‌کند، همان‌طور که از استراتژی تحقیق تجربی استفاده می‌کند.

### ۳-۴ تجزیه و تحلیل داده‌ها

برای تحقیقات کمی، متداول‌ترین روش‌های تجزیه و تحلیل داده‌ها آمار و ریاضیات محاسباتی است. آمار استنباطی است و شامل نتایج تجزیه و تحلیل برای نمونه‌های آماری و همچنین ارزیابی اهمیت نتایج است. ریاضیات محاسباتی شامل محاسبه روش‌های عددی، مدل‌سازی و استفاده از الگوریتم‌ها است (Hakansson, 2013). ریاضیات محاسباتی از کد کامپیوتری برای تجزیه و تحلیل داده‌ها استفاده می‌کند، برخلاف آمار، که توسط افراد قابل تجزیه و تحلیل است. نتایج حاصل از آزمایشات پایان‌نامه مورد نظر، داده‌های عددی را ارائه می‌دهد. تجزیه و تحلیل این داده‌ها شامل مقایسه داده‌ها است. ریاضیات محاسباتی برای تجزیه و تحلیل داده‌ها مورد نیاز نیست. بنابراین در پروژه برای تجزیه و تحلیل داده‌ها از روش آمار استفاده شده است.

### ۳-۵ تضمین کیفیت

تضمین کیفیت اعتبار و تأیید تحقیق است. داده‌ها باید قابل اعتماد، معتبر، قابل تکرار و اخلاقی باشند تا در تحقیق مورد استفاده قرار گیرند. پایی به ثبات و ثبات اندازه‌گیری‌ها اشاره دارد. این بدان معناست که هنگام در نظر گرفتن همه متغیرهای آزمایش، اندازه‌گیری‌های مختلف باید نتایج قابل قبولی را ارائه دهند. قابلیت اطمینان در هر آزمایشی که در این پروژه صورت گرفته ارزیابی شده است.

تصدیق (روایی) اطمینان می‌دهد که ابزارهای یک آزمایش آن‌گونه که انتظار می‌رود اندازه‌گیری شود، انجام شود. اعتبار این پروژه با ارزیابی مراحل و ابزارهای مورد نیاز برای انجام آزمایشات حفظ شده است.

تکرارپذیری اطمینان حاصل می‌کند که یک آزمایش با تکرار دقیق متغیرها نتایج یکسانی را به همراه خواهد داشت. در این پروژه، همه آزمایش‌ها تکرار شده است تا بررسی شود آیا نتایج یکسان هستند یا خیر.

تحقیق همچنین باید اخلاقی باشد. اخلاق اصول اخلاقی برنامه‌ریزی، انجام و گزارش نتایج حاصل از تحقیقات و همچنین حفظ حریم خصوصی و برخورد محرمانه با مطالب را پوشش می‌دهد (Hakansson, 2013). این جنبه‌های اخلاقی در تحقیق پروژه مورد بررسی نیز ارزیابی شده است.

## ۳-۶ روش‌های توسعه نرم‌افزار

یکی از روش‌های توسعه نرم‌افزار، یک برنامه (یا مجموعه‌ای از دستورالعمل‌ها) برای توسعه یک سیستم از ابتدا تا پیاده‌سازی است. مزایای استفاده از چنین روشی عبارتند از: این به درک چرخه حیات فرایند کمک می‌کند، یک رویکرد ساختاریافته برای توسعه را اجرا می‌کند و برنامه‌ریزی منابع را که در پروژه مورد استفاده قرار می‌گیرد، امکان‌پذیر می‌سازد (Agarwal, et. al., 2010).

مدل آبشار از فازهای متوالی و خطی استفاده می‌کند، که در آن هر مرحله ادامه فاز قبلی است. مراحل در مدل آبشار عبارتند از: امکان‌سنجی (درک مشکل)، الزامات و مشخصات، طراحی راه‌حل، کدگذاری و تست مازول، آزمایش سیستم، تحویل و نگهداری (Agarwal, et. al., 2010).

روش نمونه اولیه، نسخه‌های ناقص از یک محصول را که نمونه اولیه نامیده می‌شود، در حین توسعه ایجاد می‌کند. مزیت استفاده از نمونه‌های اولیه این است که توسعه‌دهندگان می‌توانند بازخورد کاربران را در حین توسعه محصول دریافت کنند و تغییرات احتمالی که باید انجام شود به راحتی انجام می‌شود. نمونه‌های اولیه همچنین می‌توانند توسعه‌دهندگان را در مورد قطعات ناتمام قسمت بهتری به توسعه‌دهندگان ارائه دهند، که می‌تواند برآورد بهتری در مورد زمان به پایان رساندن محصول ارائه دهد (Guide, et. al., 1999).

مدل چابک (سریع) یک رویکرد سازگار و انعطاف‌پذیر است که بر تحویل زودهنگام و مداوم نرم‌افزار به مشتری تمرکز دارد. درحالی‌که مدل آبشار خطی و سازگار است، مدل چابک انعطاف‌پذیر و قابل‌تغییر است. مدل چابک ۱۲ اصل را تعریف می‌کند که "بیانیه" توسعه چابک را تشکیل می‌دهند (Highsmith, 2002).

همه این روش‌ها معمولاً برای تیم توسعه‌دهندگان سازمان‌دهی می‌شود که در آن محصول برای مشتری اجرا می‌شود. پایان‌نامه موردبررسی یک پروژه انفرادی است و محصولی را برای مشتری ایجاد نمی‌کند، بنابراین پیروی از تمام اصول مرتبط با این روش‌های توسعه غیرضروری و یا غیرممکن است. در عوض، ایده کلی پشت روش‌ها موردتوجه قرار گرفته است. روش نمونه اولیه نمونه‌هایی را ایجاد می‌کند که برای آزمایش در حین توسعه محصول در نظر گرفته شده است. این پروژه شامل سه مازول است که به‌قدری ساده‌اند که ایجاد نمونه‌های اولیه ضروری نخواهد بود. در عوض نرم‌افزار در حین توسعه، آزمایش شده است. بنابراین این پروژه از روش نمونه اولیه استفاده نکرده است. انعطاف‌پذیری مدل چابک برای این پروژه مناسب‌تر از مدل آبشار است زیرا این پروژه فازهایی ندارد که به مراحل قبلی وابسته باشد تا تکمیل شود. بنابراین این پروژه از ایده کلی پشتیبان مدل چابک استفاده کرده است.

## ۳-۲ الزامات سیستم

الزامات سیستم الزامات عملکردی و غیر عملکردی هستند که باید مورد استفاده قرار گیرند تا TCS قابل استفاده و قوی باشد. الزامات بر اساس موارد استفاده برای TCS تعریف می‌شود. الزامات عملکردی خدماتی را که راه‌حل باید ارائه دهد و نحوه عملکرد راه‌حل در شرایط خاص را توصیف می‌کند. الزامات غیرکاربردی محدودیت‌هایی در خدمات ارائه‌شده توسط یک راه‌حل است. الزامات غیرکاربردی ممکن است بر عملکرد، قابلیت اطمینان و قابلیت استفاده متمرکز شوند (Sommerville, 2007).

### ۱- الزامات عملکردی TCS عبارتند از:

- در طول فشرده‌سازی هیچ داده‌ای از بین نمی‌رود.
- TCS از فشرده‌سازی بدون اتلاف استفاده می‌کند، بدین معنی که داده‌های اصلی باید پس از فشرده‌سازی قابل بازیابی باشد. اگر داده‌ها در طول فشرده‌سازی از بین بروند، راه‌حل فایده‌ای ندارد.
- TCS می‌تواند متن را با رمزگذاری‌های معمول فشرده کند.

رمزگذاری‌های رایج به کدگذاری UTF-8 و ASCII اشاره می‌کند. TCS باید بتواند رایج‌ترین رمزگذاری‌های متنی را مدیریت کند تا یک‌راه‌حل مفید باشد. چهار فایل با رمزگذاری UTF-8 و یک فایل با رمزگذاری ASCII برای آزمایش این کار استفاده می‌شود.

- **TCS می‌تواند متن را با رمزگذاری‌های غیرمعمول فشرده کند.**

در صورتی که TCS قادر به فشرده‌سازی متن با رمزگذاری‌هایی باشد که اغلب استفاده نمی‌شوند، علاوه بر رمزگذاری‌های رایج، این راه‌حل محکم است و می‌تواند انواع بیشتری از موارد استفاده را اداره کند. برای آزمایش این مورد از فایلی با کدگذاری cp037 استفاده می‌شود.

- **ماژول‌های TCS به راحتی متصل می‌شوند.**

سه واحد ماژول TCS برنامه‌های مستقل هستند که می‌توانند در ترکیب با یکدیگر استفاده شوند. ماژول‌ها باید به‌طور جداگانه به هم متصل شوند تا به‌صورت جداگانه ارزیابی شوند.

۲- الزامات غیرکاربردی برای TCS عبارتند از:

- **TCS می‌تواند فایل‌های بالای ۲۰ مگابایت را فشرده کند.**

کاهش حجم فایل بیشتر از فایل‌های کوچک‌تر روی فایل‌های بزرگ‌تر تأثیر می‌گذارد. برای ارزیابی این نیاز، یک فایل XML با حجم ۲۱,۶ مگابایت استفاده می‌شود.

- **TCS می‌تواند فایل‌ها را حداقل به نصف اندازه اصلی خود فشرده کند.**

چندین تست روی هر یک از متون انجام می‌شود و اگر حداقل یکی از تست‌ها دارای ضریب فشرده‌سازی حداقل ۲ باشد، شرایط موردنیاز برآورده می‌شود.

- **تمام کدهای شخص ثالثی که در TCS استفاده می‌شود باید منبع باز باشد.**

کد شخص ثالثی که در TCS استفاده می‌شود باید دارای مجوز باز باشد که اجازه اصلاحات را می‌دهد. به‌منظور ادغام ماژول‌ها در یک‌راه‌حل، ممکن است تغییرات در کد لازم باشد.

## **۸-۳ ویژگی مجموعه داده‌ها**

ویژگی مجموعه داده‌ها در پایان‌نامه موردبررسی به‌عنوان مجموعه‌ای از متون مختلف که TCS هنگام ارزیابی قابلیت فشرده‌سازی از آن‌ها استفاده می‌کند، تعریف شده است. این متون از رمزگذاری‌های مختلف استفاده می‌کنند، به‌عنوان مثال UTF-8 و ASCII، کدگذاری‌های معمولی هستند که ممکن است کاربر از آن‌ها استفاده کند. UTF-8 یک رمزگذاری کاراکتر بسیار رایج است که شامل انواع زیادی از نمادها است. ASCII برای متنی استفاده می‌شود که از تنوع زیادی از کاراکترها و نمادها استفاده نمی‌کند و معمولاً به زبان انگلیسی نوشته می‌شود (McEnergy and Xiao, 2005). از کدگذاری‌های مختلف برای ارزیابی متفاوت در فشرده‌سازی استفاده می‌شود و اینکه چه نوع کدگذاری برای راه‌حل بهینه است. متون از کدگذاری ASCII، cpo37 و UTF-8 استفاده خواهند کرد.

### ۱-۸-۳ تعریف مجموعه داده‌ها

مجموعه داده‌های مورد استفاده در ارزیابی پایان‌نامه موردبررسی، نمونه‌هایی از فایل‌های متنی بزرگ‌تر هستند که کاربر ممکن است TCS را روی آن‌ها اعمال کند. مجموعه داده‌ها شامل ۶ متن با ویژگی‌های مختلف است؛

۱- یک فایل XML ویکی‌پدیا. رمزگذاری UTF-8 ؛ ۲۱,۶ مگابایت

(<https://dumps.wikimedia.org/backupindex.html>)

۲- یک فایل XML ویکی‌پدیا. کدگذاری cpo37 ؛ ۱۸,۱ مگابایت

(<https://dumps.wikimedia.org/backupindex.html>)

۳- یک فایل کد نوشته‌شده با C. کدگذاری ASCII ؛ ۶۴ کیلوبایت

(<https://github.com/torvalds/linux/blob/master/kernel/sys.c>)

۴- کتابی که به زبان انگلیسی نوشته شده است. رمزگذاری UTF-8 ؛ ۶۸۰ کیلوبایت (Thus spake Zarathustra)

۵- کتابی که به زبان ایتالیایی نوشته شده است. رمزگذاری UTF-8 ؛ ۶۲۶ کیلوبایت (Magugliani, La divina commedia)

۶- کتابی که به زبان چینی نوشته شده است. رمزگذاری UTF-8 ؛ ۲۸۵ کیلوبایت (Salt Lake City)

فایل XML ویکی‌پدیا یک نسخه محلی و فقط متنی از مجموعه‌ای از مقالات (صفحات جداگانه) از ویکی‌پدیا است. فایل XML بیشتر شامل کاراکترهای ASCII است، بدین معنا که اکثر کاراکترهای متن را می‌توان در اسکی کدگذاری کرد. مجموعه داده‌ها نیز شامل همان فایل XML است که با cpo37 رمزگذاری شده است. Cpo37 رمزگذاری یک فایل پس از فشرده‌سازی آن با ماژول فشرده‌سازی ASCII است. این متن برای تست وجود دو ماژول دیگر راه‌حل (کدگذاری هافمن و کدگذاری دیکشنری) با ACM سازگار است. فایل کد C از مخزن

سیستم عامل لینوکس GitHub گرفته شده است و فقط شامل کاراکترهای ASCII است. این سه کتاب برای آزمایش وجود تفاوت بین زبان کد ساختاریافته و زبان طبیعی گنجانده شده است. کتاب انگلیسی بیشتر شامل کاراکترهای ASCII است، کتاب ایتالیایی نیز شامل بیشتر کاراکترهای ASCII است، اما از چندین کاراکتر استفاده می کند که در طرح کدگذاری ASCII گنجانده نشده است، و کتاب چینی شامل هیچ کاراکتر اسکی نمی شود. متون مختلف دارای اندازه های بسیار متنوعی هستند که برخی از آن ها چند کیلوبایت بزرگ و برخی دیگر چند مگابایت هستند. صرف نظر از این، متن ها اندازه کافی برای ارزیابی صحیح TCS دارند.

### ۹-۳ الزامات پیاده سازی

TCS شامل سه ماژول است. ACM، کدگذاری دیکشنری و کدگذاری هافمن. ACM به طور خاص برای پروژه مورد بررسی ایجاد شده است، اما کدگذاری دیکشنری و کدگذاری هافمن پیاده سازی هایی است که توسط افراد دیگر انجام شده و در راه حل یکپارچه شده است. پیاده سازی برنامه هایی است که دیگران برای پیاده سازی الگوریتم ها ساخته اند، پیاده سازی ها باید مجموعه ای از الزامات را برآورده کنند، اما همچنین باید بتوانند با سایر ماژول های TCS ادغام شوند تا راه حل به طور کلی بتواند الزامات سیستم را برآورده کند. چندین پیاده سازی برای الگوریتم های کدگذاری دیکشنری و کدگذاری هافمن یافت می شود. در پایان نامه مورد بررسی شرایط زیر را برای پیاده سازی ها در نظر گرفته است:

پیاده سازی ها باید دارای مجوز باز باشند.

پیاده سازی ها باید درست عمل کنند، حتی زمانی که داده های ورودی می توانند چند مگابایت باشند.

پیاده سازی ها باید با متن کدگذاری های مختلف سازگار باشند.

پیاده سازی ها همچنین باید بتوانند با متن کدگذاری های مختلف کار کنند. پیاده سازی هایی که الزامات را برآورده می کنند در برابر یکدیگر آزمایش می شوند و پیاده سازی هایی که بهترین نتیجه را می گیرند (بالاترین نسبت فشرده سازی) به عنوان کدگذاری دیکشنری و ماژول کدگذاری هافمن برای TCS انتخاب می شوند.

### ۱۰-۳ طراحی TCS



TCS یک راه حل فشرده سازی است که از سه ماژول تشکیل شده است: ماژول فشرده سازی (ACM) ASII، ماژول کدگذاری دیکشنری و ماژول کدگذاری هافمن. این ماژول ها به هم پیوسته هستند، به این معنی که آن ها برنامه های مستقل هستند که می توانند در یک راه حل ادغام شوند، یا می توانند به عنوان برنامه های مستقل استفاده شوند (Kaye, ۲۰۰۳). داشتن ماژول های به هم پیوسته، چابکی را در جستجوی بالاترین نسبت فشرده سازی ایجاد می کند و برای ارزیابی ماژول ها به صورت جداگانه ضروری است. ماژول ها در C یا Python نوشته شده اند و می توانند به صورت متوالی با استفاده از یک برنامه پایتون یا یک اسکریپت پوسته یونیکس اجرا شوند.

TCS ابتدا برای ACM در نظر گرفته شده است، سپس رمزگذار دیکشنری و در نهایت ماژول کدگذاری هافمن استفاده کند. ACM ابتدا استفاده می شود زیرا فقط قادر به فشرده سازی کاراکترهای اسکی است. خروجی رمزگذار دیکشنری یا ماژول های کدگذاری هافمن، داده های دودویی هستند که شامل کاراکترهای اسکی نیستند، بنابراین این ماژول ها نباید قبل از ACM استفاده شوند. خروجی ACM یک فایل متنی در کدگذاری cpo37 است. این خروجی به کدگذار دیکشنری ارسال می شود که زیررشته های تکراری را در متن جستجو می کند. سپس خروجی دودویی از رمزگذار دیکشنری به ماژول کدگذاری هافمن ارسال می شود که مقادیر بایت را که بیشتر از بقیه استفاده می شود، جستجو می کند. فایل خروجی از ماژول کدگذاری هافمن آخرین فایل فشرده شده از TCS است.

### ۱-۱۰-۳ طراحی ACM

ACM برنامه ای است که در Python 3.6.9 نوشته شده است و شامل دو فایل اسکریپت است: `ascii-compression.py` و `charprocessing.py`. `ascii-compression.py` مسئول خواندن و نوشتن فایل هاست، در حالی که `charprocessing.py` فشرده سازی متن واقعی را کنترل می کند.

`Ascii-compression.py` الگوریتم فشرده سازی است که ACM از آن استفاده می کند مشابه الگوریتم شوکو است و با ادغام کاراکترها با یکدیگر کار می کند. ادغام کاراکتر به این معنی است که دو کاراکتر از یک فایل متنی، به عنوان مثال "a" و "b"، می توانند با هم در فضای یک کاراکتر واحد ذخیره شوند. این الگوریتم تنها قادر به ادغام کاراکترهای اسکی است. همان طور که در بخش ۳-۲ ذکر شد، کاراکترهای اسکی فقط از ۷ بیت آخر برای کد نماد استفاده می کنند، در حالی که از اولین بیت می توان برای تشخیص خطاها در داده ها استفاده کرد. خروجی ACM (فایل فشرده) یک فایل متنی با کدگذاری cpo37 است. برخلاف ASCII، کدگذاری cpo37 از همه ۸ بیت برای کد نماد استفاده می کند. در مورد ACM، اولین بیت از یک کاراکتر cpo37 برای نشان دادن ادغام یا نبودن کاراکتر استفاده می شود و ۷ بیت آخر کد نماد یک یا دو کاراکتر است. رمزگذاری cpo37 به عنوان فرمت خروجی انتخاب

می‌شود زیرا از همه ۸ بیت به جای ۷ بیت ASCII استفاده می‌کند. cpo37 همچنین یک کدگذاری استاندارد در پایتون است و درست مانند US-ASCII برای زبان انگلیسی در نظر گرفته شده است.

#### Ascii-compression.py

اسکرپت ascii-compression.py یک فایل متنی با کدگذاری ASCII، UTF-8 یا cpo37 به عنوان ورودی می‌گیرد. فایل‌های باینری یا دیگر کدگذاری‌های متنی را نمی‌پذیرد. سپس هر خط از فایل متنی را به اسکرپت charprocessing.py ارسال می‌کند که فشرده‌سازی را کنترل می‌کند. پس از نوشتن فایل خروجی، یک آزمایش ساده را انجام می‌دهد تا بررسی کند که آیا فایل خروجی بزرگ‌تر از فایل اصلی است یا خیر. این می‌تواند در صورتی اتفاق بیفتد که فایل ورودی دارای تعداد زیادی کاراکتر باشد که نمی‌توان آن‌ها را در اسکی کدگذاری کرد، اگر فایل خروجی بزرگ‌تر از فایل اصلی باشد، داده‌های فشرده‌شده حذف می‌شوند، زیرا هدف برنامه فشرده‌سازی ایجاد فایل‌های کوچک‌تر است. در عوض، فایل خروجی با داده‌های اصلی بازنویسی می‌شود و فایل فشرده‌شده را به اندازه فایل اصلی می‌سازد. در نهایت، فایل فشرده یک خط جدید اضافه می‌شود که شامل 0 یا 1 است. این عدد نشان می‌دهد که آیا فایل فشرده تغییر کرده است (0) یا کپی از فایل اصلی است (1). این اطلاعات برای فشرده‌ساز مهم است، که الگوریتم فشرده‌سازی فایل را کاهش می‌دهد یا همان‌طور که هست کپی می‌کند. هنگامی که یک فایل فشرده از حالت فشرده خارج می‌شود، عدد حذف‌شده و فایل متنی به حالت اصلی خود بازگردانده می‌شود.

#### charprocessing-py

اسکرپت charprocessing.py وظیفه فشرده‌سازی و رفع فشار را بر عهده دارد. الگوریتم فشرده‌سازی هر کاراکتر یک متن را تکرار کرده و آن‌ها را در یک فایل فشرده با کدگذاری cpo37 ذخیره می‌کند. نحوه ذخیره هر کاراکتر به دو عامل بستگی دارد. آیا می‌توان کاراکتر را در ASCII رمزگذاری کرد و آیا bigram زیر فقط شامل کاراکترهای معمولی است. اگر کاراکتر در برنامه کدگذاری ASCII گنجانده نشده باشد، به این معنی است که کاراکتر ممکن است بیش از یک بایت ذخیره‌سازی را اشغال کند. مثال UTF-8 کاراکتر "e" است که از دو بایت استفاده می‌کند. وقتی این کاراکتر فشرده می‌شود، الگوریتم باید علامت دهد که از بیش از یک بایت استفاده می‌کند. در غیر این صورت، هنگام خروج از فشرده‌سازی الگوریتم تصور می‌کند که دو کاراکتر وجود دارد که هرکدام از یک بایت استفاده کرده‌اند. قبل از نوشتن کاراکتر، دو علامت بک اسلش "\" به کاراکتر اضافه می‌شود، که نشان می‌دهد بایت بعدی کاراکتر اسکی نیست. سپس کاراکتر ویژه به عددی که نشان‌دهنده کاراکتر است تبدیل می‌شود. این امر ضروری است زیرا اگر کاراکتر خاص بیشتر از یک بایت باشد، نمی‌تواند در cpo37 کدگذاری شود. سرانجام، یک بک اسلش "\" پس از کاراکتر اضافه می‌شود، به این معنی که کاراکتر ویژه به پایان رسیده است. داشتن علامت قبل

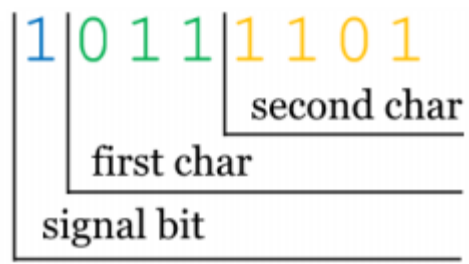
و بعد از کاراکتر ضروری است زیرا طول کاراکتر می‌تواند بین ۲ تا ۴ بایت باشد. نسخه فشرده‌شده کاراکتر "e" به‌صورت \332\ نشان داده می‌شود.

اگر فایل اصلی حاوی کاراکترهای بک اسلش "\" است، باید آن‌ها را علامت دهید تا هنگام خروج از فشرده‌سازی الگوریتم برای یک کاراکتر غیر ASCII اشتباه نگیرد. به‌عنوان مثال اگر دو بک اسلش "\" در یک سطر توسط کمپرسور به‌عنوان علامت نباشند، دی کمپرسور انتظار دارد کاراکتر بعدی یک کاراکتر ویژه باشد. روشی که سیگنال کمپرسور به بک اسلش نشان می‌دهد این است که قبل از آن دو بک اسلش دیگر اضافه می‌کند. این بدان معناست که هر بک اسلش در فایل اصلی با سه بک اسلش جایگزین می‌شود. هنگامی که دستگاه فشرده‌ساز در یک فایل فشرده یا بک اسلش مواجه می‌شود، می‌داند که کاراکتر بعدی نیز یک بک اسلش خواهد بود، اما یکی پس‌از آن می‌تواند آغاز یک کاراکتر ویژه یا سومین بک اسلش باشد. اگر کاراکتر خاصی باشد، به‌درستی از حالت فشرده خارج می‌شود، و اگر سومین بک اسلش باشد، دو بک اسلش برداشته می‌شود. اگر یک کاراکتر ASCII و کاراکتر بعدی آن هر دو در لیست کاراکترهای مشترک گنجانده شوند، آن‌ها باهم ادغام می‌شوند. ادغام تنها در صورتی کار می‌کند که دو کاراکتر متوالی مشترک باشند. اگر هیچ‌یک، یا فقط یک مورد، از کاراکترهای متوالی در لیست کاراکترهای معمولی نباشند، ادغام نمی‌شوند. در عوض، آن‌ها همانند فایل اصلی در فایل فشرده ذخیره می‌شوند، به‌عنوان مثال، کاراکتر "a" دارای یک بیت ارزش است ۰۱۱۰۰۰۰۱. این نشان‌دهنده باینری عدد ۹۷ است و به این صورت است که ACM کاراکترها را می‌خواند. اولین بیت ۰ است که به الگوریتم می‌گوید که ۷ بیت بعدی نشان‌دهنده یک کاراکتر ادغام نشده است.

```
2 firstCommonChars = ["e","t","a","o","n","i","h"," "]
3
4 secondCommonChars = ["e","t","a","o","n","i","h"," ","s","r","l","d","u","c","m","w"]
5
```

شکل ۳-۱. فهرست کاراکترهای متداول

اگر یک کاراکتر در لیست firstCommonChars گنجانده شده باشد، همان‌طور که در شکل ۳-۱ نشان داده شده است، و کاراکتر بعدی در لیست secondCommonChars گنجانده شده باشد، باهم ادغام می‌شوند. این لیستی از رایج‌ترین حروف (به‌علاوه نماد فاصله) در زبان انگلیسی است (Troost, 2020) و مدل فشرده‌سازی ACM است. یک مدل فشرده‌سازی می‌تواند لیستی از تکرار کاراکترها، نمادها یا بیگرم‌ها باشد. ACM برای همه انواع متون از یک مدل استفاده می‌کند، به این معنی که متون انگلیسی باید نسبت فشرده‌سازی بهتری نسبت به متون سایر زبان‌ها داشته باشند.



شکل ۲-۳ نمایش باینری از کاراکترهای ادغام شده

یک کاراکتر ادغام شده از ۳ قسمت تشکیل شده است: بیت سیگنال، کاراکتر اول و کاراکتر دوم. برای یک کاراکتر ادغام شده، اولین بیت روی ۱ تنظیم می شود. سه بیت بعدی نشان دهنده موقعیت اولین کاراکتر در لیست firstCommonChars است. در مثال شکل ۲-۴، اولین بیت های کاراکتر ۰۱۱ هستند، که مربوط به عددی است که اولین کاراکتر دارای کاراکتر با شاخص ۳ در firstCommonchars است که "o" است. کاراکتر دوم از ۴ بیت استفاده می کند، به این معنی که می تواند دو برابر کاراکتر اول را نشان دهد. در مثال، بیت های کاراکتر دوم ۱۱۰۱ هستند که با عدد ۱۳ مطابقت دارد. کاراکتر با شاخص ۱۳ در secondCommonChars کاراکتر "c" است. نتایج حاصل از فشرده سازی نشان می دهد که کاراکتر ادغام شده در شکل ۲-۴ نشان دهنده بیگرم "oc" است.

## فصل چهارم

### ارزیابی و پیاده‌سازی الگوریتم‌ها و بررسی نتایج

در فصل پنجم پایان‌نامه موردبررسی روند ارزیابی واحدهای مختلف واحد TCS و ارزیابی سیستم را به‌طورکلی شرح می‌دهد. ارزیابی ماژول‌های کدگذاری دیکشنری و کدگذاری هافمن شامل انتخاب پیاده‌سازی‌هایی با بیشترین نسبت فشرده‌سازی برای استفاده در TCS است. ارزیابی ACM شامل اندازه‌گیری نسبت فشرده‌سازی آن و اندازه‌گیری نسبت فشرده‌سازی Shoco و مقایسه نتایج است. ارزیابی کل سیستم (TCS) شامل اندازه‌گیری نسبت فشرده‌سازی آن و مقایسه عملکرد آن با سایر برنامه‌های فشرده‌سازی عمومی است.

#### ۴-۱ ارزیابی ماژول‌های کدگذاری دیکشنری و کدگذاری هافمن

مجموعه‌ای از پیاده‌سازی‌های احتمالی برای کدگذاری دیکشنری و کدگذاری هافمن باهم مقایسه شده است. مجموعه داده‌ها روی هر یک از پیاده‌سازی‌ها آزمایش شده است. مجموعه داده‌ها به‌صورت محلی در کامپیوتر ذخیره شده و نسبت فشرده‌سازی اندازه‌گیری شده است. نسبت فشرده‌سازی اندازه‌گیری میزان کوچک‌تر بودن فایل فشرده در مقایسه با فایل فشرده نشده است و با تقسیم اندازه فشرده نشده بر اندازه فشرده شده یک فایل محاسبه شده است. به‌عنوان مثال، اگر یک فایل فشرده نشده ۱۰ مگابایت و فایل فشرده آن ۸ مگابایت است، نسبت پیاده‌سازی برای آن فایل خاص ۱,۲۵ است. اگر نسبت فشرده‌سازی ۱ باشد، فایل فشرده شده دارای اندازه یکسانی با فایل فشرده نشده است، بدین معنی که الگوریتم قادر به فشرده‌سازی فایل نیست. پیاده‌سازی‌هایی که بالاترین نسبت فشرده‌سازی را به‌طور متوسط برای هر متن در مجموعه داده دارند، به‌عنوان ماژول‌های کدگذاری دیکشنری و کدگذاری هافمن برای TCS انتخاب شده‌اند. میانگین نسبت فشرده‌سازی به‌عنوان معیار انتخاب شد.

#### ۴-۲ انتخاب پیاده‌سازی کدگذاری دیکشنری

در پایان نامه بررسی شده چندین کدگذاری دیکشنری برای TCS در نظر گرفته شده است. این پیاده‌سازی‌ها برخی از الگوریتم‌های کدگذاری دیکشنری را نشان می‌دهند. الگوریتم دیکشنری LZ78، LZ77، LZ و Re-Pair است.

این الگوریتم‌ها توسط برنامه‌های فشرده‌سازی محبوب استفاده می‌شوند بنابراین موثر بودن آن‌ها ثابت شده است. الگوریتم‌های کدگذاری دیکشنری بیشتری وجود دارد، اما توسط برنامه‌های فشرده‌سازی محبوب، استفاده نمی‌شود. بسیاری از پیاده‌سازی‌های موجود در وبسایت کنترل نسخه GitHub مشخص نمی‌کند که راه حل آن‌ها دارای چه مجوزی است و بنابراین نمی‌توان آن‌ها را در نظر گرفت.

ماژول‌های استاندارد در پایتون برای برنامه‌های فشرده‌سازی محبوب مانند Gzip، Bz2 یا Zipfile وجود دارد، اما این ماژول‌ها همه از چند الگوریتم در ترکیب با یکدیگر استفاده می‌کنند. به عنوان مثال، همه این ماژول‌ها از کدگذاری هافمن استفاده می‌کنند و برخی از کدگذاری دیکشنری استفاده می‌کنند، اما نمی‌توان فقط از الگوریتم کدگذاری دیکشنری و کدگذاری هافمن استفاده کرد. پیاده‌سازی‌هایی که مورد استفاده قرار می‌گیرند نمی‌توانند بسته‌های چند الگوریتم باشند.

## ۱-۲-۴ مقایسه پیاده‌سازی‌ها

در پایان‌نامه ذکر شده یک پیاده‌سازی پایتون و دو پیاده‌سازی C شرایط مورد نیاز را برآورده می‌کند. آن‌ها دارای مجوز باز هستند، حتی برای فایل‌های بزرگ نیز به درستی کار می‌کنند و با متن کدگذاری‌های مختلف سازگار هستند. سپس پیاده‌سازی‌ها با یکدیگر مقایسه شده و نسبت فشرده‌سازی برای هر متن در مجموعه داده‌ها اندازه‌گیری شده است. پیاده‌سازی پایتون با پایتون نسخه ۲,۷,۱۷ اجرا می‌شود و پیاده‌سازی C با GNU Compiler Collection نسخه ۷,۵,۰ کامپایل شده است.

یک پیاده‌سازی الگوریتم LZ77 است و در پایتون نوشته شده است. پیاده‌سازی ۱۵ ژوئن ۲۰۲۰ بارگیری شده است. این پیاده‌سازی به کاربر این امکان را می‌دهد که اندازه پنجره الگوریتم (پنجره لغزان) (Ziv, Lempel, 1977) را تغییر دهد. اندازه پنجره تعیین می‌کند که چقدر الگوریتم متن در هر نقطه برای پیدا کردن زیررشته‌های منطبق تجزیه و تحلیل می‌کند. اندازه پنجره بزرگ‌تر نسبت فشرده‌سازی بیشتری را ارائه می‌دهد، اما الگوریتم را کندتر می‌کند. اندازه پنجره ۴۰۰ کاراکتر برای همه آزمایش‌ها استفاده شده، زیرا به گفته نویسنده پایان‌نامه این حداکثر اندازه ممکن برای پنجره بوده است.

پیاده‌سازی دیگر، پیاده‌سازی C الگوریتم LZW است. پیاده‌سازی در ۹ ژوئن ۲۰۲۰ بارگیری شد. الگوریتم LZW مانند LZ77 پنجره‌های لغزان استفاده نمی‌کند، اما در عوض زیررشته‌ها را به یک دیکشنری اضافه می‌کند که ممکن

است در متن تکرار شود یا نشود. اگر یک زیررشته تکرار شود، آن را با اشاره به ورودی در دیکشنری جایگزین می‌کند.

آخرین پیاده‌سازی C الگوریتم LZ77 است. پیاده‌سازی ۱۵ ژوئن ۲۰۲۰ بارگیری شد. درست مانند اجرای پایتون، امکان تغییر اندازه پنجره نیز وجود دارد. حداکثر اندازه پنجره ممکن برای این پیاده‌سازی ۱,۰۴۸,۵۷۶ کاراکتر است و بنابراین از این برای همه آزمایش‌ها استفاده شده است.

جدول ۴-۱. نتایج حاصل از پیاده‌سازی کدگذاری دیکشنری

Data set	مجموعه داده‌ها	LZ77 (Python)	LZW (C)	Lz77 (C)
XML file (21.6 MB)	فایل XML	۱.۴۷	۱.۹۲	۱.۸۵
cpo37 encoded XML file (18.1 MB)	فایل XML کدگذاری cpo37	۱.۴۵	۱.۶۱	۱.۷۵
C code file (64 KB)	فایل کد C	۱.۷۵	۲.۱۳	۲.۵
English book (680 KB)	کتاب انگلیسی	۱.۳	۱.۹۶	۱.۶۱
Italian book (626 KB)	کتاب ایتالیایی	۱.۳۳	۲.۱۳	۱.۵۶
Chinese book (285 KB)	کتاب چینی	۱.۱۶	۱.۵۹	۱.۳۴
Average	میانگین	۱.۴۱	۱.۸۹	۱.۷۷

جدول ۴-۱ نشان‌دهنده نسبت فشرده‌سازی سه کدگذاری دیکشنری برای هر متن در مجموعه داده است. اجرای پایتون از الگوریتم LZ77 بدترین نسبت فشرده‌سازی را دارد. پیاده‌سازی C برای الگوریتم LZ77 دارای نسبت فشرده‌سازی بسیار بهتری است، حتی اگر از الگوریتم LZ77 نیز استفاده شود. پیاده‌سازی LZW نسبت به پیاده‌سازی C در LZ77 برای برخی از متون دارای نسبت فشرده‌سازی بالاتری است، اما برای متون دیگر نسبت کمتری دارد. پیاده‌سازی که بهترین نتایج متوسط را از آزمایش داشته باشد، به‌عنوان ماژول کدگذاری دیکشنری برای TCS انتخاب شده بنابراین پیاده‌سازی C الگوریتم LZW به‌عنوان ماژول کدگذاری دیکشنری استفاده شده است.

### ۴-۳ انتخاب پیاده‌سازی کدگذاری هافمن

تعدادی پیاده‌سازی منبع باز از الگوریتم کدگذاری هافمن پیدا شده و مورد ارزیابی قرار گرفته، اکثر پیاده‌سازی‌هایی که پیدا می‌شوند مشخص نمی‌کنند که راه‌حل آن‌ها دارای چه مجوزی است و بنابراین نمی‌توان آن‌ها را در نظر

گرفت. مازول‌های استاندارد پایتون برای برنامه‌های فشرده‌سازی محبوب که از الگوریتم هافمن استفاده می‌کنند وجود دارد، اما با الگوریتم‌های دیگر ترکیب شده است. دو پیاده‌سازی پایتون یافت می‌شود که قادر به فشرده‌سازی متن ASCII و UTF-8 هستند، اما هیچ‌یک از آن‌ها قادر به فشرده‌سازی متن کد شده cpo37 نیستند. پیاده‌سازی در C یافت شد، اما قادر به ساخت کد منبع نیست، بنابراین نمی‌توان آن را ارزیابی کرد.

### ۱-۳-۴ مقایسه پیاده‌سازی‌ها

در پایان‌نامه موردبررسی سه پیاده‌سازی الزامات را برآورده می‌کنند و با یکدیگر مقایسه شده‌اند. یک پیاده‌سازی پایتون و دو پیاده‌سازی C. پیاده‌سازی پایتون با پایتون نسخه ۳,۶,۹ و پیاده‌سازی C با GNU Compiler Collection نسخه ۷,۵,۰ کمپایل شده است.

پیاده‌سازی پایتون Dahuffman نام دارد و توسط کاربر GitHub، Stefaan Lippens ساخته شده است. پیاده‌سازی ۲۶ ژوئیه ۲۰۲۰ بارگیری شده است. پیاده‌سازی این قابلیت را دارد که بر اساس فهرستی از تکرارهای نماد، یک جدول کد بسازد. فراوانی نماد می‌تواند توسط کاربر ارائه شود یا پیاده‌سازی می‌تواند آن را از یک متن ورودی محاسبه کند. اگر فراوانی نماد از یک متن ورودی محاسبه شود، پیاده‌سازی ابتدا یک جدول کد ایجاد می‌کند، و سپس متن را با استفاده از جدول کد ذکر شده فشرده می‌کند.

دو پیاده‌سازی C، توسط کاربران GitHub، Gagarine Yaikhom و Doug Richardson، جداول کد را به‌طور خودکار در طول فشرده‌سازی ایجاد می‌کنند. پیاده‌سازی Yaikhom در ۲۶ ژوئیه ۲۰۲۰، و اجرای Richardson در ۲۸ ژوئیه ۲۰۲۰ بارگیری شده است.

جدول ۲-۴. نتایج حاصل از پیاده‌سازی کدگذاری هافمن

Data set	مجموعه داده‌ها	Dahuffman	Yaikhom C impl.	Richardson C impl.
XML file (21.6 MB)	فایل XML	۱.۵۱	۱.۵	۱.۵
cpo37 encoded XML file (18.1 MB)	فایل XML کدگذاری cpo37	۱.۲	۱.۲	۱.۲
C code file (64 KB)	فایل کد C	۱.۴۸	۱.۴۷	۱.۴۷
English book (680 KB)	کتاب انگلیسی	۱.۸	۱.۷	۱.۷
Italian book (626 KB)	کتاب ایتالیایی	۱.۹۵	۱.۳۷	۱.۳۷
Chinese book (285 KB)	کتاب چینی	۲.۳۲	۱.۳۱	۱.۳۱
Average	میانگین	۱.۷۱	۱.۴۹	۱.۴۹



جدول ۲-۴ نشان‌دهنده نسبت فشرده‌سازی سه برنامه‌نویسی هافمن برای هر متن در مجموعه داده است که نشان می‌دهد، برای اکثر متون تفاوت چندانی بین پیاده‌سازی‌ها وجود ندارد. درواقع، دو پیاده‌سازی C نتایج دقیقاً یکسانی را به دست آوردند، اگرچه دو پیاده‌سازی متفاوت توسط افراد مختلف است. پیاده‌سازی پایتون، Dahuffman، نسبت فشرده‌سازی کمی بهتر از دو پیاده‌سازی C برای کتاب انگلیسی و ایتالیایی، اما نسبت بهتری برای کتاب چینی دارد. در مورد سه متن دیگر، Dahuffman کم‌ویش نتایج مشابهی با سایر پیاده‌سازی‌ها دارد. این می‌تواند به این دلیل باشد که این کتاب‌ها بیشتر از سایر متون دارای کاراکتر UTF-8 هستند، و Dahuffman در فشرده‌سازی کاراکترهای غیر اسکی کار آبی بهتری دارد. کتاب چینی به‌طور کامل از کاراکترهای UTF-8 تشکیل شده است، بنابراین این منطقی‌ترین نظریه به نظر می‌رسد. Dahuffman، دارای بالاترین نسبت فشرده‌سازی متوسط پیاده‌سازی‌ها است و بنابراین به‌عنوان ماژول کدگذاری هافمن برای TCS استفاده شده است.

#### ۴-۴ ارزیابی ACM

برنامه فشرده‌سازی TCS شبیه به Shoco است. هر دو TCS و Shoco برنامه‌های فشرده‌سازی متن بدون اتلاف هستند که در فشرده‌سازی متن ASCII تخصص دارند. یک تفاوت اصلی این است که Shoco یک الگوریتم است درحالی‌که TCS یک راه‌حل فشرده‌سازی است که از ۳ الگوریتم کدگذاری LZW، ACM و Huffman تشکیل شده است. ACM عملکرد و راه‌حل بسیار مشابهی با Shoco دارد. پایان‌نامه موردبررسی به‌منظور مقایسه عادلانه، فقط ACM در این قسمت با Shoco مقایسه کرده است.

مقاله موردبحث در ۲-۵-۲ یکی دیگر از برنامه‌های فشرده‌سازی ASCII است. باین‌حال، این راه‌حل در مقایسه با ACM به دلیل محدودیت‌های موجود نیست. راه‌حل پیشنهادی پایان‌نامه موردنظر فقط قادر به فشرده‌سازی متن ASCII است و هیچ کدگذاری دیگری ندارد. از متون مورداستفاده در مجموعه داده‌ها، فقط فایل کد C با راه‌حل سازگار است، بنابراین نمی‌توان مقایسه عادلانه‌ای با ACM داشت.

ACM از یک مدل فشرده‌سازی پیش‌فرض استفاده می‌کند که برای داده‌هایی که قرار است فشرده شود آموزش ندیده است. یک مدل فشرده‌سازی لیستی از کاراکترها، نمادها یا بیگرام‌ها است که اغلب در یک متن استفاده می‌شود. Shoco این قابلیت را دارد که یک مدل فشرده‌سازی را برای یک فایل متنی خاص آموزش دهد، به این معنی که نمادها و بیگرام‌های موجود در مدل اغلب در یک متن استفاده می‌شود. ACM این قابلیت را ندارد و از

یک مدل فشرده‌سازی عمومی استفاده می‌کند که برای کلمات انگلیسی بهینه‌شده است. در مقایسه، شوکو هم با مدل پیش‌فرض خود و هم با یک مدل آموزش‌دیده آزمایش می‌شود.

#### ۴-۴-۱ نتایج حاصل از مقایسه Shoco و ACM

جدول ۴-۳. نتایج حاصل از مقایسه Shoco و ACM

Data set	مجموعه داده‌ها	ACM	Shoco پیش‌فرض	Shoco آموزش‌دیده
XML file (21.6 MB)	فایل XML	۱.۱۹	۱.۲۳	۱.۴
cpo37 encoded XML file (18.1 MB)	فایل XML کدگذاری cpo37	۱	۰.۵۹	۰.۷۵
C code file (64 KB)	فایل کد C	۱.۲۱	۱.۱۷	۱.۳۵
English book (680 KB)	کتاب انگلیسی	۱.۴۶	۱.۳۲	۱.۵۶
Italian book (626 KB)	کتاب ایتالیایی	۱.۲۹	۱.۱۵	۱.۵
Chinese book (285 KB)	کتاب چینی	۱	۰.۵۳	۰.۷۵
Average	میانگین	۱.۱۹	۱	۱.۲۲

فایل XML با کد cpo37 و کتاب چینی دو متنی هستند که نسبت فشرده‌سازی ۱ برای هر دو الگوریتم دارند، به این معنی که حجم فایل فشرده کوچک‌تر نمی‌شود. این متون به ترتیب از نویسه cpo37 و UTF-8 تشکیل شده‌اند. ACM و Shoco هر دو فقط قادر به فشرده‌سازی کاراکترهای ASCII هستند، و بنابراین رمزگذاری‌های دیگر می‌توانند خروجی را بزرگ‌تر کنند. برای هر دو Shoco آموزش‌دیده و پیش‌فرض، نسبت فشرده‌سازی کمتر از ۱ است، به این معنی که فایل فشرده بزرگ‌تر از فایل اصلی است.

ACM دارای نسبت ۱ است، به این معنی که فایل فشرده دارای اندازه یکسان با فایل اصلی است. دلیل این امر این است که ACM آزمایشی را برای بررسی بزرگ‌تر بودن فایل فشرده از نسخه اصلی انجام می‌دهد. اگر چنین باشد، فشرده‌سازی کنار گذاشته می‌شود و فایل فشرده بدون تغییر باقی می‌ماند. Shoco این آزمایش را ندارد و بنابراین یک فایل فشرده بزرگ‌تر دریافت می‌کند.

کتاب انگلیسی بالاترین نسبت فشرده‌سازی را برای هر دو الگوریتم دارد. مدل‌های پیش‌فرض این دو الگوریتم برای کلمات انگلیسی بهینه‌شده‌اند، بنابراین نسبت بالایی برای کتاب انگلیسی به دست می‌آورند. این نسبت هنگام استفاده از یک مدل آموزش‌دیده برای شوکو حتی بیشتر است زیرا مدل برای آن متن خاص سفارشی‌شده است.

## ۴-۵ ارزیابی TCS

ارزیابی TCS اندازه‌گیری نسبت فشرده‌سازی در مجموعه داده‌ها و سپس مقایسه نتایج با سایر برنامه‌های فشرده‌سازی عمومی است. وقتی TCS ارزیابی شده است، از هر سه ماژول استفاده کرده است یعنی خروجی از یک ماژول به ماژول بعدی ارسال می‌شود، نسبت فشرده‌سازی برای خروجی ماژول قبلی اندازه‌گیری شده است. درنهایت، نسبت تراکم کل اندازه‌گیری می‌شود. این نسبت اندازه فایل اصلی در مقایسه با اندازه نهایی فایل فشرده است. ابتدا ACM، سپس کدگذاری دیکشنری (LZW) و درنهایت کدگذاری هافمن استفاده شده است.

### ۴-۵-۱ نتایج حاصل از ارزیابی TCS

جدول ۴-۴. نتایج حاصل از ارزیابی TCS

Data set	مجموعه داده‌ها	ACM	LZW	Huffman coding	Total ratio
XML file (21.6 MB)	فایل XML	۱.۱۹	۱.۶	۱.۰۱	۱.۹۳
cpo37 encoded XML file (18.1 MB)	فایل XML کدگذاری cpo37	۱	۱.۶	۱.۰۱	۱.۶۲
C code file (64 KB)	فایل کد C	۱.۲۱	۱.۸۳	۱.۰۱	۲.۲۳
English book (680 KB)	کتاب انگلیسی	۱.۴۶	۱.۳۶	۱	۱.۹۹
Italian book (626 KB)	کتاب ایتالیایی	۱.۲۹	۱.۵۷	۱	۲.۰۲
Chinese book (285 KB)	کتاب چینی	۱	۱.۵۹	۱	۱.۵۹
Average	میانگین	۱.۱۹	۱.۵۹	۱.۰۱	۱.۹

همان‌طور که گفته شد، نسبت فشرده‌سازی اندازه‌گیری شده برای هر ماژول بر اساس خروجی ماژول قبلی است. "نسبت کل" نسبت فشرده‌سازی است که از TCS (هر سه ماژول) برای متون موجود در مجموعه داده به دست می‌آید. نسبت کل نتیجه‌ای است که در ارزیابی با سایر برنامه‌های فشرده‌سازی در بخش بعدی مورد استفاده قرار می‌گیرد.

همان‌طور که جدول ۴-۴ نشان می‌دهد، برخی از فایل‌ها بیشتر از ACM و برخی دیگر از LZW سود بیشتری می‌برند. به عنوان مثال فایل کد C کتاب انگلیسی از ACM نسبت کوچک‌تری از LZW دارد. به این دلیل است که کتاب انگلیسی از کاراکترهای مکرر زبان انگلیسی بیشتری استفاده می‌کند و بنابراین نسبت بهتری از ACM دریافت می‌کنند، درحالی‌که فایل کد C دارای کلمات و زیررشته‌های تکراری‌تر است و بنابراین نسبت بهتری از LZW دریافت می‌کند.

جدول ۴-۴ همچنین نشان می‌دهد که کدگذاری هافمن تقریباً هیچ فشرده‌سازی را برای هیچ‌یک از فایل‌ها به دست نمی‌آورد. این می‌تواند به این دلیل باشد که برنامه‌نویسی هافمن برای فشرده‌سازی داده‌های باینری طراحی نشده است. همچنین ممکن است خروجی دودویی از LZW بسیار تصادفی و نامنظم باشد تا کدگذاری هافمن نتواند به فشرده‌سازی برسد. خروجی LZW قبلاً توسط ACM فشرده‌شده است و ممکن است ترکیب این دو مازول نیاز به کدگذاری هافمن را برطرف کند. نتایج حاصل از این ارزیابی نیاز به ارزیابی جدیدی دارد که فقط از کدگذاری LZW و هافمن استفاده می‌شود، تا به پاسخ این سؤال برسیم که آیا ACM قابلیت‌های مازول کدگذاری هافمن را مختل می‌کند یا این مازول قادر به فشرده‌سازی داده‌های باینری نیست.

جدول ۴-۵. نتایج حاصل از ارزیابی مازول کدگذاری LZW و هافمن

Data set	مجموعه داده‌ها	LZW	Huffman coding	Total ratio
XML file (21.6 MB)	فایل XML	۱.۹۲	۱	۱.۹۲
C code file (64 KB)	فایل کد C	۲.۱۳	۱.۰۱	۲.۱۴
English book (680 KB)	کتاب انگلیسی	۱.۹۶	۱	۱.۹۶
Italian book (626 KB)	کتاب ایتالیایی	۲.۱۳	۱	۲.۱۳
Average	میانگین	۲.۰۴	۱	۲.۰۴

فایل XML کد شده cpo37 و کتاب چینی در این ارزیابی گنجانده نشده‌اند زیرا ACM هیچ تأثیری بر روی این فایل‌ها نخواهد داشت و هدف از این ارزیابی مقایسه با نتایجی است که ACM روی آن تأثیر گذاشته است. جدول ۴-۵ نشان می‌دهد که حتی وقتی از ACM استفاده نمی‌شود، مازول کدگذاری هافمن عملاً قادر به فشرده‌سازی داده‌های باینری نیست.

## ۴-۵-۲ نتایج حاصل از مقایسه TCS

جدول ۴-۶ نشان می‌دهد که TCS نسبت به سایر برنامه‌ها برای هر فایل در مجموعه داده نسبت فشرده‌سازی کمتری دارد. Bzip2 و Zip بالاترین میانگین نسبت فشرده‌سازی را دارند زیرا از الگوریتم‌های متفاوتی نسبت به Zip و Gzip استفاده می‌کنند که بر اساس الگوریتم DEFLATE است (Korpela, 2006). TCS از راه‌حلی مشابه DEFLATE، اما با ACM به‌عنوان پیش‌پردازنده استفاده می‌کند. الگوریتم DEFLATE ترکیبی از رمزگذار دیکشنری LZ77 و کدگذاری هافمن است، درحالی‌که TCS از رمزگذار دیکشنری LZW و کدگذاری هافمن استفاده می‌کند.

جدول ۴-۶. نتایج حاصل از مقایسه TCS

Data set	مجموعه داده‌ها	TCS	Bzip2	Zip	v-Zip	Gzip
XML file (21.6 MB)	فایل XML	۱.۹۳	۳.۷۹	۲.۹۶	۴.۲۴	۲.۹۶
cpo37 encoded XML file (18.1 MB)	فایل XML کدگذاری cpo37	۱.۶۲	۳.۱۲	۲.۴۵	۳.۴۸	۲.۴۵
C code file (64 KB)	فایل کد C	۲.۲۳	۴.۲۲	۳.۸۸	۴.۲۵	۳.۹۱
English book (680 KB)	کتاب انگلیسی	۱.۹۹	۳.۵۲	۲.۶۵	۳.۱۸	۲.۶۵
Italian book (626 KB)	کتاب ایتالیایی	۲.۰۲	۳.۵۴	۲.۶۶	۳.۱۴	۲.۶۶
Chinese book (285 KB)	کتاب چینی	۱.۵۹	۲.۵۹	۲.۰۱	۲.۳۵	۲.۰۱
Average	میانگین	۱.۹	۳.۴۶	۲.۷۷	۳.۴۴	۲.۷۷

جدول ۴-۵ نشان می‌دهد که ماژول کدگذاری هافمن مورد استفاده برای TCS با ماژول LZW سازگار نیست و بنابراین TCS جایگزینی برای الگوریتم DEFLATE ندارد. برای بررسی اینکه آیا TCS در صورت جایگزینی ماژول کدگذاری LZW و Huffman که در TCS با یک الگوریتم DEFLATE خالص جایگزین شده است، نتایج بهتری خواهد گرفت یا خیر، باید ارزیابی جدیدی انجام شود. در این ارزیابی، از برنامه‌های Zip و Gzip با ACM به عنوان پیش پردازنده استفاده شده است. این ارزیابی به این سؤال پاسخ می‌دهد که آیا استفاده از ACM به عنوان پیش پردازنده باعث افزایش نسبت فشرده‌سازی الگوریتم DEFLATE می‌شود یا خیر.

جدول ۴-۷. نتایج حاصل از مقایسه ACM + DEFLATE

Data set	مجموعه داده‌ها	Zip	ACM + Zip	Gzip	ACM + Gzip
XML file (21.6 MB)	فایل XML	۲.۹۶	۲.۹۲	۲.۹۶	۲.۹۲
C code file (64 KB)	فایل کد C	۳.۸۸	۳.۷۷	۳.۹۱	۳.۸۲
English book (680 KB)	کتاب انگلیسی	۲.۶۵	۲.۶۶	۲.۶۵	۲.۶۶
Italian book (626 KB)	کتاب ایتالیایی	۲.۶۶	۲.۵۸	۲.۶۶	۲.۵۸
Average	میانگین	۳.۰۴	۲.۹۸	۳.۰۵	۳

مانند جدول ۴-۵، فایل XML کد شده cpo37 و کتاب چینی در این ارزیابی گنجانده نشده‌اند زیرا ACM هیچ تأثیری بر روی این فایل‌ها نخواهد داشت. همان‌طور که جدول ۴-۷ نشان می‌دهد، هنگامی که ACM به عنوان پیش پردازنده Zip یا Gzip استفاده می‌شود، نسبت فشرده‌سازی کمی پایین‌تر است. ACM فقط نسبت به کتاب انگلیسی

نسبت فشرده‌سازی را تا حدی بیشتر نشان می‌دهد، اما افزایش ۰,۰۱ نسبت فشرده‌سازی ناچیز است. استفاده از ACM به‌علاوه الگوریتم DEFLATE که Zip و Gzip استفاده می‌کند نسبت فشرده‌سازی بسیار بالاتری نسبت به TCS دارد.

## ۴-۶ بررسی نتایج پایان‌نامه

فصل ۶ پایان‌نامه موردبحث به بررسی نتایج حاصل از ارزیابی‌ها را پوشش می‌دهد. همچنین به کاستی‌های راه‌حل، مشارکت این پایان‌نامه و این‌که در چه شرایطی می‌تواند مورد استفاده قرار گیرد، می‌پردازد.

### ۴-۶-۱ کار آبی TCS

TCS مدعی الگوریتم DEFLATE بود، اما با تمرکز بر فشرده‌سازی متن ASCII. از نظر مفهومی، TCS از تکنیک‌های مشابه DEFLATE (کدگذاری دیکشنری و کدگذاری هافمن) استفاده می‌کند اما ACM را به‌عنوان یک گام جدید در کانال فشرده‌سازی معرفی می‌کند. تفاوت بین DEFLATE و TCS در این است که DEFLATE از رمزگذار دیکشنری LZ77 استفاده می‌کند در حالی که TCS از رمزگذار دیکشنری LZW استفاده می‌کند. ارزیابی TCS نشان می‌دهد که جایگزین DEFLATE که TCS از آن استفاده می‌کند ناقص است. بنابراین، استفاده از TCS (در حالت فعلی آن) به‌عنوان جایگزینی برای الگوریتم DEFLATE گزینه مناسبی نیست. مقایسه برنامه‌های فشرده‌سازی در جدول ۴-۶ همچنین نشان داد که برنامه‌هایی که از راه‌حل‌های جایگزین برای الگوریتم DEFLATE استفاده می‌کردند نسبت فشرده‌سازی بالاتری داشتند. بنابراین، حتی اگر TCS دارای یک الگوریتم کارکرد DEFLATE باشد، بازهم نمی‌تواند با نسبت‌های فشرده‌سازی Bzip2 و Zip-۷ رقابت کند.

### ۴-۶-۲ مقایسه Shoco و ACM

شوکو مدعی نزدیک ACM است و کار مرتبطی است که بیشتر در پایان‌نامه بررسی شده، موردبحث قرار گرفته است. هر دو الگوریتم دارای پیش‌فرض یکسان و راه‌حل‌های فنی مشابه هستند، اما یکسان نیستند. یک تفاوت پیچیدگی مدل‌های فشرده‌سازی است که دو الگوریتم از آن استفاده می‌کنند. تفاوت دیگر کار آبی الگوریتم‌ها است. الگوریتم شوکو به‌عنوان "کمپرسور سریع برای رشته‌های کوتاه" توصیف می‌شود و به نظر می‌رسد نقطه اصلی فروش شوکو سرعت آن باشد. ACM برای تأخیر اندازه‌گیری نشده است، اما آزمایش نشان داده است که Shoco (با مدل

فشرده‌سازی پیش‌فرض) سریع‌تر از ACM است. دلیل اینکه ACM برای تأخیر اندازه‌گیری نشده این است که تمرکز ACM بر قابلیت فشرده‌سازی آن است و نه بر کارایی آن. ACM به‌عنوان اثبات این مفهوم بود که الگوریتم فشرده‌سازی ASCII از لحاظ نظری امکان‌پذیر است. این قبل از اینکه شوکو برای نویسنده شناخته شود، ثابت شد که فشرده‌سازی ASCII امکان‌پذیر است.

### کاستی‌های ACM

تکنیک فشرده‌سازی که ACM برای کاراکترهای غیر ASCII استفاده می‌کند بهینه نیست. برخی از کاراکترهای UTF-8 مانند ایموجی‌ها می‌توانند تا ۴ بایت در یک فایل فشرده نشده استفاده کنند. این کاراکترها در صورت فشرده‌سازی از ۹ بایت استفاده می‌کنند. فشرده‌سازی غیر بهینه با تبدیل کاراکتر به عددی که آن را نشان می‌دهد و نوشتن رقم عدد برای رقم ایجاد می‌شود. این می‌تواند بهینه‌سازی شود، اما محدودیت زمانی منجر به این مسئله می‌شود که هنوز باید برطرف شود. فایلی که با ACM فشرده می‌شود دارای یک خط اضافه‌شده یا ۱ یا ۰ می‌باشد. این به این معنی است که آیا فایل تغییر کرده است یا نه، و هنگامی که فایل از حالت فشرده خارج می‌شود، حذف می‌شود. مشکلی که ممکن است رخ دهد این است که اگر یک فایل دارای خط ۱ یا ۰ به‌عنوان آخرین خط باشد، و مشخص نیست که فایل فشرده است یا نه. اگر کاربر این فایل فشرده نشده را با یک فایل فشرده اشتباه کند، پس از فشرده‌سازی تغییر می‌کند. یک فایل فشرده به‌عنوان یک فایل متنی (txt) ذخیره می‌شود و بنابراین احتمالاً دارای همان نوع فایل با یک فایل فشرده نشده است، بنابراین نمی‌توان گفت آیا فایل فشرده است یا نه.

ACM به زبان برنامه‌نویسی پایتون نوشته شده است، که در زمان اجرا به‌طور قابل‌توجهی برای مثال کندتر از زبان C است. زبان برنامه‌نویسی C به‌طور کلی سریع‌تر از پایتون است، به همین دلیل برنامه‌های پایتون اغلب به دلایل کارایی ماژول‌های C را ادغام می‌کنند (Beazley, 2009). اگر ACM به زبان C نوشته می‌شد، به‌طور بالقوه می‌توانست سریع‌تر باشد.

## فصل پنجم

### جمع‌بندی و پیشنهادها

#### ۱-۵ مقدمه

پایان‌نامه موردبررسی راه‌حلی با TCS برای فشرده‌سازی متن با تمرکز ویژه بر فشرده‌سازی متن ASCII ارائه می‌دهد. TCS مشکلات ذخیره و انتقال زمان زیادی از داده‌های متنی را برطرف می‌کند. TCS حجم فایل‌ها را کاهش می‌دهد بنابراین به فضای ذخیره‌سازی کمتری نیاز دارند. علاوه بر این، TCS هنگام بارگیری یا انتقال فایل فشرده از طریق اینترنت، پهنای باند را کاهش می‌دهد. TCS برای متن ASCII سفارشی‌شده است با این حال برای کاربر روی هر شکل از متن طراحی شده است. جاه‌طلبی TCS فشرده‌سازی فایل‌هایی برای کاربرانی است که با مقادیر زیادی متون ASCII یا ASCII سنگین سروکار دارند. TCS از سه ماژول ACM، ماژول کدگذاری دیکشنری و ماژول کدگذاری هافمن تشکیل شده است. TCS مبتنی بر الگوریتم DEFLATE است که همچنین کدگذاری دیکشنری و کدگذاری هافمن را ترکیب می‌کند. TCS از ACM به عنوان پیش پردازنده جایگزین DEFLATE برای افزایش فشرده‌سازی برای آزمایش ASCII استفاده می‌کند. ACM توسط نویسنده به طور خاص برای این پروژه توسعه یافته است، درحالی‌که دیکشنری و ماژول‌های کدگذاری هافمن پیاده‌سازی شده توسط افراد دیگر است. ماژول‌های خاص مورد ارزیابی قرار گرفته‌اند، و TCS ارزیابی شده و با برنامه‌های فشرده‌سازی عمومی مورد مقایسه قرار گرفته است.

#### ۲-۵ نتایج حاصل از تحقیق

نتایج حاصل از ارزیابی TCS نشان می‌دهد که در حالت فعلی، مدعی سایر برنامه‌های فشرده‌سازی عمومی نیست. مقایسه فشرده‌سازی اسکی با الگوریتم Shoco نشان می‌دهد که هنگام استفاده از مدل‌های فشرده‌سازی عمومی، ACM نسبت فشرده‌سازی متوسط بهتری دارد. بنابراین ACM می‌تواند مدعی رشته تخصصی تر فشرده‌سازی ASCII باشد. ارزیابی ACM (جدول ۳-۴) نشان می‌دهد که می‌توان آن را در هر شکلی از متن بدون افزایش اندازه فشرده استفاده کرد، اما نسبت فشرده‌سازی قابل توجهی فقط در متون سنگین ASCII به دست می‌آید. با این وجود، نسبت فشرده‌سازی ACM در مقایسه با برنامه‌های فشرده‌سازی عمومی بسیار کوچک است (جدول ۶-۴). ترکیب ACM با برنامه‌های فشرده‌سازی مبتنی بر DEFLATE (جدول ۷-۴) نسبت فشرده‌سازی متوسط را فقط تحت DEFLATE



نشان می‌دهد، با یک مثال از مزیت جزئی ACM. ممکن است بهبود الگوریتم فشرده‌سازی ACM نسبت فشرده‌سازی آن را افزایش دهد و حتی ممکن است از DEFLATE فراتر رود (در صورتی که ACM به‌عنوان پیش پردازنده DEFLATE استفاده شود).

## ۱-۲-۵ پاسخ به سؤالات تحقیق

"چه ترکیبی از تکنیک‌ها می‌تواند نسبت فشرده‌سازی متن ASCII را بهبود بخشد؟"

بدیهی است که ACM برای فایل‌های ASCII به فشرده‌سازی دست می‌یابد، اگرچه برخی از فایل‌های متنی نسبت فشرده‌سازی بیشتری نسبت به سایرین دارند (جدول ۳-۴). علیرغم شایستگی ACM، استفاده از آن در ترکیب با کدگذاری دیکشنری و ماژول کدگذاری هافمن (و همچنین برنامه‌های فشرده‌سازی مبتنی بر DEFLATE) نسبت فشرده‌سازی متن ASCII را بهبود نداده است. دلیل این امر این است که خروجی ACM کمتر از متن ساده قابل فشرده‌سازی است. ACM فقط به‌عنوان پیش پردازنده برنامه‌های فشرده‌سازی مبتنی بر DEFLATE آزمایش شده است، زیرا TCS جایگزین DEFLATE بود. ممکن است ترکیب الگوریتم‌های دیگر (مانند الگوریتم‌های مورد استفاده در Bzip2 و v-Zip) با ACM به‌عنوان پیش پردازنده بتواند نسبت فشرده‌سازی متن ASCII را افزایش دهد. این کار آینده TCS است.

"آیا استفاده از فشرده‌سازی ASCII در ترکیب با کدگذاری دیکشنری و کدگذاری هافمن از نسبت فشرده‌سازی DEFLATE فراتر می‌رود؟"

در حالت فعلی، TCS (ACM، ماژول کدگذاری دیکشنری و ماژول کدگذاری هافمن) از نسبت فشرده‌سازی DEFLATE فراتر نمی‌رود. همان‌طور که گفته شد، ماژول برنامه‌نویسی هافمن TCS با ماژول برنامه‌نویسی دیکشنری سازگار نیست و بنابراین TCS جایگزین DEFLATE نیست. در عوض، ACM در ترکیب با برنامه‌نویسی دیکشنری و برنامه‌نویسی هافمن برنامه‌های فشرده‌سازی مبتنی بر DEFLATE در جدول ۷-۴ آزمایش شده است. بر اساس میانگین نسبت فشرده‌سازی از ۴ متن مجموعه داده، با استفاده از ACM در ترکیب با رمزگذاری دیکشنری و هافمن کدگذاری از نسبت فشرده‌سازی DEFLATE فراتر نمی‌رود.

"آیا این ترکیب فشرده‌سازی متن نسبت فشرده‌سازی بالاتری نسبت به برنامه‌های فشرده‌سازی عمومی برای متن ASCII به دست می‌آورد؟"

در مقایسه بین TCS و ۴ برنامه فشرده‌سازی عمومی در جدول ۴-۶، ۴ متون بیشتر شامل کاراکترهای ASCII هستند. همان‌طور که گفته شد، نه ترکیب TCS و نه ACM + DEFLATE نسبت فشرده‌سازی بالاتری نسبت به برنامه‌های فشرده‌سازی مبتنی بر DEFLATE برای متون ASCII ندارند. علاوه بر این، Bzip2 و Zip-۷ از هر دو برنامه فشرده‌سازی مبتنی بر DEFLATE و TCS برای هر متنی در مورد استفاده بهتر هستند. در نتیجه، TCS نسبت فشرده‌سازی بیشتر از برنامه‌های فشرده‌سازی عمومی برای متن ASCII به دست نمی‌آورد.

"آیا این ترکیب فشرده‌سازی متن نسبت فشرده‌سازی بالاتری نسبت به راه‌حل‌های فشرده‌سازی که روی متن ASCII تخصص دارند، به دست می‌آورد؟"

در پایان‌نامه مورد بررسی، ACM تنها با فشرده‌سازی اسکی توسط الگوریتم Shoco مقایسه شده است. تا آنجا که نویسنده مطلع است، Shoco تنها الگوریتم فشرده‌سازی اسکی دیگری است که برای فشرده‌سازی متن اسکی سفارشی شده است و همچنین قادر به فشرده‌سازی کدگذاری‌های دیگر است. راه‌حل‌های فشرده‌سازی ASCII که قادر به فشرده‌سازی سایر کدگذاری‌ها نیستند، مانند مقاله مورد بحث در قسمت ۲-۵-۲، برای مقایسه در نظر گرفته شده است، زیرا تعداد بسیار کمی فایل‌های متنی امروزه فقط از کاراکترهای ASCII استفاده می‌کنند در مقایسه ACM و Shoco، ACM نسبت فشرده‌سازی بالاتری نسبت به Shoco دارد - با استفاده از مدل فشرده‌سازی پیش‌فرض خود، اما ACM نسبت به Shoco نسبت کمی کمتر دارد - هنگام استفاده از یک مدل آموزش‌دیده. در بخش ۴-۴-۱ ذکر شده است، مدل‌های پیش‌فرض و آموزش‌دیده شوکو دارای مزایا و معایبی هستند. هر دو مدل قابل اجرا هستند و بسته به شرایط می‌توان از آن‌ها استفاده کرد. بنابراین، این سؤال تحقیق را نمی‌توان با بله یا خیر قطعی پاسخ داد.

### ۳-۵ پیشنهادها و کارهای آینده

کارهای آینده در TCS شامل یافتن یک ماژول کدگذاری هافمن است که می‌تواند خروجی دودویی را از ماژول کدگذاری دیکشنری فشرده کند. دستیابی به این امر می‌تواند قابلیت فشرده‌سازی TCS را با الگوریتم DEFLATE برابر کند. آزمایش الگوریتم‌هایی که توسط DEFLATE برای بهبود نسبت فشرده‌سازی استفاده نشده است، نیز کار آینده است.

برای افزایش نسبت فشرده‌سازی می‌توان پیشرفت‌هایی در ACM انجام داد. ACM می‌تواند فشرده‌سازی کاراکترهای غیر ASCII را بهبود بخشد، زیرا تکنیک فعلی بهینه نیست. افزودن قابلیت ایجاد مدل‌های آموزش‌دیده ممکن است

نسبت فشردسازی آن را نیز به میزان قابل توجهی افزایش دهد. با این پیشرفت‌ها، ممکن است ACM در صورت استفاده از پیش پردازنده، نسبت فشردسازی برنامه‌های فشردسازی عمومی را افزایش دهد.

همچنین می‌توان افزایش تأخیر ACM را بهبود داد. کد را می‌توان در C بازنویسی کرد و بهینه‌سازی را برای فرآیند فشردسازی و خروج از فشردسازی انجام داد. اگر ACM بیشتر بر زمان تأخیر تمرکز داشت تا نسبت فشردسازی، استفاده از ACM در شرایطی که سرعت در اولویت بیشتری نسبت به ضریب فشردسازی است، مفید خواهد بود. به‌عنوان مثال، می‌توان آن را در ترکیب با کمپرسور اسنپی مبتنی بر LZ77 در زمینه و پس‌زمینه سرور استفاده کرد. Snappy نسبت فشردسازی را از نظر سرعت به خطر می‌اندازد و در فایل‌های متنی نسبت به فایل‌های باینری به نسبت بهتری دست می‌یابد. نتایج جدول ۴-۴ نشان می‌دهد که رمزگذار دیکشنری (LZW) می‌تواند نسبت فشردسازی مناسبی را با ACM به‌عنوان پیش پردازنده به دست آورد، بنابراین استفاده از ACM به‌عنوان پیش پردازنده Snappy می‌تواند نسبت فشردسازی را افزایش دهد.

## ۴-۵ جمع‌بندی

در آغاز مطالعه این پایان‌نامه به‌عنوان تحقیق سمینار دانشجویی تصور متفاوتی نسبت به فشردسازی داده‌ها و روش‌های تحقیق در این زمینه داشتم. با انجام مطالعه پایان‌نامه با موضوع "بهبود نسبت فشردسازی متن برای متن ASCII" اطلاعات مفیدی در این زمینه به دست آوردم. با مطالعات مرتبط بیشتری سعی کردم در این گزارش تحقیق سمینار مطالب درک شده را گردآوری کنم.

این گزارش شامل ۵ فصل است. فصل اول به تعریف و مقدمه، فصل دوم مفاهیم فشردسازی داده‌ها، فصل سوم مروری است بر کارهای انجام‌شده طرح پیشنهادی پایان‌نامه، فصل چهارم ارزیابی و پیاده‌سازی الگوریتم‌ها و بررسی نتایج و فصل پنجم نیز جمع‌بندی و نتیجه‌گیری نهایی است.

پس از درک موضوع و تجزیه و تحلیل آن‌ها در راستای ادامه کار موضوعات پیشنهادی مطرح شد.

- K. Sayood, Introduction to data compression, 3rd ed. Amsterdam ; Boston: Elsevier, 2006.
- D. Salomon, Data compression the complete reference. New York: Springer, 2004.
- K. K. Shukla and M. V. Prasad, Lossy image compression: domain decomposition-based algorithms. London ; New York: Springer, 2011.
- «GNU Gzip». <https://www.gnu.org/software/gzip/manual/gzip.html#Overview> (opened Sep. 13th, 2020).
- «WinRAR - WinRAR Documentation». <https://documentation.help/WinRAR/> (opened Sep. 14th, 2020).
- D. Huffman, «A Method for the Construction of Minimum-Redundancy Codes», Proc. IRE, vol. 40, nr. 9, s. 1098–1101, sep. 1952, doi: 10.1109/JRPROC.1952.273898.
- J. Ziv and A. Lempel, «A universal algorithm for sequential data compression», IEEE Trans. Inform. Theory, vol. 23, nr. 3, pp. 337–343, May 1977, doi: 10.1109/TIT.1977.1055714.
- Farina, G. Navarro, and J. R. Parama, “Boosting Text Compression with Word-Based Statistical Encoding”, The Computer Journal, vol. 55, no. 1, pp. 111–131, Jan. 2012, doi: 10.1093/comjnl/bxr096.
- «Shoco - a fast compressor for short strings». <https://ed-vonschleck.github.io/shoco/> (opened Sep. 14th, 2020).
- «torvalds/linux: Linux kernel source tree». <https://github.com/torvalds/linux> (opened Sep. 14th, 2020).
- Y. Q. Shi and H. Sun, Image and Video Compression for Multimedia Engineering, Boca Raton, FL: CRC Press, 1999, pp. 140–141.
- M. Aggarwal and A. Narayan, «Efficient Huffman decoding», in Proceedings 2000 International Conference on Image Processing (Cat. No.00CH37101), Vancouver, BC, Canada, 2000, vol. 1, pp. 936–939, doi: 10.1109/ICIP.2000.901114.
- Langiu, «On parsing optimality for dictionary-based text compression—the Zip case», Journal of Discrete Algorithms, vol. 20, pp. 65– 70, May 2013, doi: 10.1016/j.jda.2013.04.001.
- «bzip2 and libbzip2, version 1.0.8». <https://www.sourceware.org/bzip2/manual/manual.html> (opened Sep. 14th, 2020).
- R. Togneri and C. J. S. DeSilva, Fundamentals of information theory and coding design. Boca Raton: Chapman & Hall/CRC, 2002.
- R. Hankerson, G. A. Harris, and P. D. Johnson, Introduction to information theory and data compression, 2nd ed. Boca Raton, Fla: Chapman & Hall/CRC Press, 2003.
- M. Burrows and D. J. Wheeler, «A block-sorting lossless data compression algorithm», 1994.
- J. K. Korpela, Unicode explained, 1st ed. Sebastopol, CA: O'Reilly, 2006.
- «The Computer Journal». <https://academic.oup.com/comjnl> (opened Oct. 29th, 2020).

- Silva de Moura, G. Navarro, N. Ziviani, and R. Baeza-Yates, “Fast and flexible word searching on compressed text,” *ACM Trans. Inf. Syst.*, vol. 18, no. 2, pp. 113–139, Apr. 2000, doi: 10.1145/348751.348754.
- «Information Technology Coding and Computing – ITCC». <http://www.itcc.info/> (opened Nov. 3rd, 2020).
- Håkansson, «Portal of Research Methods and Methodologies for Research Projects and Degree Projects», 2013, [Online].
- B. Agarwal, S. P. Tayal, and M. Gupta, *Software engineering & testing: an introduction*. Sudbury, Mass: Jones and Bartlett, 2010.
- M. McEnery and R. Z. Xiao, «Character encoding in corpus construction.», in *Developing Linguistic Corpora : A Guide to Good Practice*, M. Wynne, Red. Oxford, UK: AHDS, 2005.
- «JSON Syntax». [https://www.w3schools.com/js/js\\_json\\_syntax.asp](https://www.w3schools.com/js/js_json_syntax.asp) (opened Sep. 16th, 2020).
- «HTML Charset». [https://www.w3schools.com/html/html\\_charset.asp](https://www.w3schools.com/html/html_charset.asp) (opened Sep. 16th, 2020). [72] «Wikimedia Downloads». <https://dumps.wikimedia.org/backupindex.html> (opened Sep. 24th, 2020).
- «torvalds/linux». <https://github.com/torvalds/linux/blob/master/kernel/sys.c> (opened Sep. 24th, 2020).
- F. Nietzsche and T. Common, *Thus spake Zarathustra*. Ware: Wordsworth Editions, 1997.
- Dante Alighieri, B. Garavelli, and L. Magugliani, *La divina commedia*. Milano: BUR Rizzoli, 2012.
- Ainajushi, 豆瓣閒話. Salt Lake City: Project Gutenberg, 2008.
- «Open license - Creative Commons». [https://wiki.creativecommons.org/wiki/Open\\_license](https://wiki.creativecommons.org/wiki/Open_license) (opened Sep. 16th, 2020).
- Kaye, *Loosely coupled: the missing pieces of Web services*. Marin County, Calif: RDS Press, 2003.
- «LZW compression». [https://rosettacode.org/wiki/LZW\\_compression#Python](https://rosettacode.org/wiki/LZW_compression#Python) (opened Sep. 16th, 2020).
- «gzip — Support for gzip files — Python 3.8.6rc1 documentation». <https://docs.python.org/3/library/gzip.html> (opened Sep. 16th, 2020).
- «dahuffman». <https://github.com/soxofaan/dahuffman> (opened Sep. 16th, 2020).
- T. Summers, “Hardware based GZIP compression, benefits and applications,” *CORPUS*, vol. 3, pp. 2–68, 2008.
- Belanger, K. Church, and A. Hume, “Virtual Data Warehousing, Data Publishing and Call Detail,” in *Databases in Telecommunications*, vol. 1819, W. Jonker, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2000, pp. 106–117.
- «7z Format». <https://www.7-zip.org/7z.html> (opened Oct. 18th, 2020).
- M. Salib, “Faster than C: Static type inference with Starkiller” in *PyCon Proceedings*, Mar. 2004, pp. 2–26.
- M. Beazley, *Python essential reference*, 4th ed. Upper Saddle River, NJ: Addison-Wesley, 2009.

## واژه‌نامه

### واژه‌نامه فارسی

Sliding window	پنجره کشویی
Binary data	داده‌های دودویی
Loss compression	فشرده‌سازی با اتلاف
Lossless compression	فشرده‌سازی بدون اتلاف
Data compression	فشرده‌سازی داده
Coding dictionary	کدگذاری دیکشنری
Coding Huffman	کدگذاری هافمن
Leaf node	گره برگ
Child node	گره فرزند
Data set	مجموعه داده‌ها
Git hab repository	مخزن گیت هاب

## واژه‌نامه انگلیسی

Binary data	داده‌های دودویی
Child node	گره فرزند
Coding dictionary	کدگذاری دیکشنری
Coding Huffman	کدگذاری هافمن
Data compression	فشرده‌سازی داده
Data set	مجموعه داده‌ها
Git hab repository	مخزن گیت هاب
Leaf node	گره برگ
Loss compression	فشرده‌سازی با اتلاف
Lossless compression	فشرده‌سازی بدون اتلاف
Sliding window	پنجره کشویی

## Abstract

Data compression is a field that has been extensively researched. Many compression algorithms in use today have been around for several decades, like Huffman Coding and dictionary coding. These are general-purpose compression algorithms and can be used on anything from text data to images and video. There are, however, much fewer lossless algorithms that are customized for compressing certain types of data, like ASCII text. This project is about creating a text-compression solution using a combination of the three compression algorithms dictionary coding, ASCII compression, and Huffman coding. The solution is customized for compressing ASCII text, but it can be used on any form of text. The algorithms will be combined to create a prototype that will be evaluated against general-purpose compression programs. An evaluation will also be made against the ASCII compression program Shoco. The results from the evaluation show that combining dictionary coding, ASCII compression, and Huffman coding does not surpass the compression ratio achieved from general-purpose compression programs.

Keywords Text compression, Dictionary coding, ASCII, Huffman coding





**Payam Noor University**

**Department of Computer Engineering and Information Technology**

**Seminar Report (M.S.c)**

**Title:**

**Improving the text compression ratio for ASCII  
text (Review)**

**Supervisor:**

**Dr. Seyed Ali Razavi Ebrahimi**

**By:**

**Maryam Sadat Mourdgar**

**September 2021**