

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه جامع علمی کاربردی
مرکز آموزش عالی علمی-کاربردی جهاد دانشگاهی مشهد

معرفی خلاصه‌سازی ماشینی متن و معرفی روش‌های خلاصه‌سازی

متون در وب

پایان نامه کارشناسی رشته برنامه نویسی تحت وب

طیبه هنرور دوقلعه

سیده زکیه اتحادی

سعید راه داران

شهره راد رحیمی

استاد راهنما:

دکتر محمد عبدالهی



دانشگاه جامع علمی کاربردی
مرکز آموزش عالی علمی-کاربردی جهاد دانشگاهی مشهد

پایان نامه کارشناسی رشته برنامه نویسی تحت وب خانم ها طیبه هنرور دوقلعه،

سیده زکیه اتحادی، شهره راد رحیمی و آقای سعید راه داران

تحت عنوان معرفی خلاصه سازی ماشینی متن و معرفی روش های خلاصه سازی

متون در وب

در تاریخ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهایی قرار گرفت.

۱- استاد راهنمای پایان نامه دکتر محمد عبدالهی

۲- استاد مشاور پایان نامه مهندس علی نیک فرجام

۳- استاد داور

۴- استاد داور

معاون شعبه:

تشر و قدردانی

به مصداق «من لم یشکر المخلوق لم یشکر الخالق» بسی شایسته است از استاد فرهیخته و فرزانه جناب آقای دکتر محمد عبدالهی که با کرامتی چون خورشید ، سرزمین دل را روشنی بخشیدند و گلشن سرای علم و دانش را با راهنمایی های کار ساز و سازنده بارور ساختند ، تقدیر و تشر نماییم.

همچنین از پدر و مادران عزیز ، دلسوز و مهربانمان که آرامش روحی و آسایش فکری را فراهم نمودند تا با حمایت های همه جانبه در محیطی مطلوب ، مراتب تحصیلی و نیز پایان نامه درسی را به نحو احسن به اتمام برسانیم ، سپاسگزاری می نماییم.

کلیه حقوق معنوی اعم از چاپ، تکثیر، نسخه برداری، ترجمه، اقتباس و حقوق مادی مترقب بر نتایج مطالعات ابتکارات و نوآوریهای ناشی از تحقیق موضوع پایان نامه متعلق به مرکز آموزش عالی علمی- کاربردی جهاددانشگاهی مشهد است. نقل مطالب با ذکر مأخذ بلامانع است.

تقدیم به

ماحصل آموخته هایم را تقدیم می کنم به آنان که مهر آسمانی شان آرام بخش آلام زمینی ام است
به استوارترین تکیه گاهم ، دستان پرمهر پدرم
به سبزترین نگاه زندگیم ، چشمان سبز مادرم
که هرچه آموختم در مکتب عشق شما آموختم و هرچه بکوشم قطره ای از دریای بی کران مهربانیتان را سپاس
نتوانم بگویم.

امروز هستی ام به امید شماسست و فردا کلید باغ بهشتم رضای شما
ره آوردی گران سنگ تر از این ارزان نداشتم تا به خاک پایتان نثار کنم، باشد که حاصل تلاشم نسیم گونه غبار
خستگیان را بزداید.

بوسه بر دستان پرمهرتان

فهرست مطالب

چکیده	۲
معرفی موضوع تحقیق	۲
۱-۱ مقدمه	۳
۲-۱ اهمیت و ضرورت تحقیق	۴
۳-۱ سوالات تحقیق	۴
۴-۱ فناوری تعریف شده در تحقیق	۵
۵-۱ خلاصه فصل	۵
فناوری های پیشین	۶
۲-۱ مقدمه	۷
۲-۲ مختصر اشاره به فناوریهای قبلی	۸
۳-۲ خلاصه فصل	۹
بررسی ابزارهای خودکار خلاصه سازی زبان های دنیا برای استفاده در خلاصه سازی متون زبان فارسی	۱۰
۱-۳ مقدمه	۱۱
۲-۳ فرایند خلاصه سازی	۱۱
۳-۳ فرایند خلاصه سازی کامپیوتری	۱۱
۴-۳ مراحل اساسی خلاصه سازی	۱۱
۵-۳ علوم لازم برای تجزیه و تحلیل متن	۱۲
۶-۳ چه مشکلاتی در تجزیه و تحلیل متن داریم؟	۱۲
۷-۳ زبان فارسی چه مشکلاتی در طراحی خلاصه ساز خودکار به همراه دراد؟	۱۳
۸-۳ انواع خلاصه	۱۴
۱-۸-۳ خلاصه استخراجی	۱۴
۲-۸-۳ خلاصه چکیده	۱۴
۹-۳ انواع خلاصه ساز	۱۴
۱-۹-۳ خلاصه سازهای عمومی	۱۴
۲-۹-۳ خلاصه سازهای موضوعی	۱۵
۱۰-۳ مراحل خلاصه سازی	۱۵
۱-۱۰-۳ در مرحله پیش پردازش متن چه رخ می دهد؟	۱۵

- ۱۶-۱۰-۳ مرحله پردازش متن ۲-۱۰-۳
- ۱۶-۱۰-۳ روش آماری ۱-۲-۱۰-۳
- ۱۷-۱۰-۳ روش موجودیتی-معنایی ۲-۲-۱۰-۳
- ۱۷-۱۰-۳ سطح کلامی ۳-۲-۱۰-۳
- ۱۸-۱۰-۳ روش ترکیبی ۴-۲-۱۰-۳
- ۱۸-۱۰-۳ مرحله تولید خلاصه ۳-۱۰-۳
- ۱۸-۱۱-۳ کارهای انجام شده ۱۱-۱۱-۳
- ۱۸-۱۱-۳ خلاصه چکیده ۱-۱۱-۳
- ۱۹-۱۱-۳ خلاصه تجربی ۱-۱-۱۱-۳
- ۱۹-۱۱-۳ گرامر موردی ۲-۱-۱۱-۳
- ۱۹-۱۱-۳ خلاصه استخراجی ۲-۱۱-۳
- ۱۹-۱۱-۳ روش های برپایه پردازش سطحی ۱-۲-۱۱-۳
- ۱۹-۱۱-۳ روش های برپایه کلمات و مکان ۱-۱-۲-۱۱-۳
- ۲۰-۱۱-۳ خلاصه ساز FARSISUM ۲-۱-۲-۱۱-۳
- ۲۱-۱۱-۳ روش های موجودیتی برپایه معنی ۲-۲-۱۱-۳
- ۲۱-۱۱-۳ استفاده از زنجیره لغوی ۱-۲-۲-۱۱-۳
- ۲۱-۱۱-۳ روش های برپایه گراف ۲-۲-۲-۱۱-۳
- ۲۲-۱۱-۳ روش های برپایه ساختار کلامی ۳-۲-۱۱-۳
- ۲۳-۱۱-۳ روش های ترکیبی ۴-۲-۱۱-۳
- ۲۳-۱۱-۳ روش های برپایه آموزش یادگیری ۱-۴-۲-۱۱-۳
- ۲۴-۱۱-۳ روش های افزایش و کاهش جملات ۲-۴-۲-۱۱-۳
- ۲۴-۱۲-۳ خلاصه سازی چندسنده ۱۲-۳
- ۲۴-۱۳-۳ روش های ارزیابی خلاصه ساز ۱۳-۳
- ۲۴-۱۴-۳ خلاصه فصل ۱۴-۳
- تولید خودکار عنوان مستقل از زبان برای متون با استفاده از روش فشرده سازی آماری جملات ۲۶
- ۱-۴ مقدمه ۲۷
- ۲-۴ متدولوژی ۲۷
- ۱-۲-۴ ویژگیهای مستقل از زبان برای امتیازدهی به جملات ۲۸
- ۱-۱-۲-۴ طول جمله ۲۸
- ۲-۱-۲-۴ موقعیت جمله ۲۸
- ۳-۱-۲-۴ کلمات اشاره ۲۸

۲۹ ۴-۱-۲-۴ کلمات کلیدی
۲۹ ۵-۱-۲-۴ وزن کلمه
۲۹ ۶-۱-۲-۴ اسامی خاص و مقادیر عددی
۲۹ ۷-۱-۲-۴ شباهت جمله به جمله
۳۰ ۸-۱-۲-۴ کلمات کم اهمیت
۳۰ ۲-۲-۴ انتخاب جمله
۳۰ ۳-۲-۴ فشرده سازی آماری جمله منتخب
۳۱ ۳-۴ نتایج ارزیابی
۳۲ ۴-۴ خلاصه فصل
۳۳ بررسی کلی خلاصه سازی خودکار متون به روش استخراجی
۳۴ ۱-۵ مقدمه
۳۴ ۲-۵ استخراجی
۳۵ ۳-۵ فرایند خلاصه سازی متن به روش استخراجی
۳۵ ۴-۵ مدل هایی که از روش استخراجی استفاده کرده اند
۳۶ ۵-۵ ویژگی های خلاصه سازی متن در روش استخراجی
۳۶ ۱-۵-۵ ویژگی کلمات کلیدی
۳۷ ۲-۵-۵ ویژگی کلمات عنوانی
۳۷ ۳-۵-۵ ویژگی محل استقرار مکان جمله
۳۷ ۴-۵-۵ ویژگی طول جمله
۳۷ ۵-۵-۵ ویژگی اسم مناسب
۳۷ ۶-۵-۵ ویژگی کلماتی که باحروف بزرگ نوشته شده اند
۳۷ ۷-۵-۵ ویژگی عبارت نشانه دار
۳۸ ۸-۵-۵ ویژگی کلمات جانبدارانه
۳۸ ۹-۵-۵ ویژگی مبتنی بر فونت
۳۸ ۱۰-۵-۵ ضمیر
۳۸ ۱۱-۵-۵ جملاتی که دارای همبستگی هستند
۳۸ ۱۲-۵-۵ جملاتی که مرکز ثقل آنها دارای انسجام است
۳۸ ۱۳-۵-۵ بروز اطلاعات غیر ضروری
۳۹ ۱۴-۵-۵ تجزیه و تحلیل گفتمان
۳۹ ۶-۵-۵ روش های خلاصه سازی استخراجی
۳۹ ۱-۶-۵ روش TF-IDF

۳۹	۲-۶-۵ روش مبتنی بر خوشه‌بندی
۴۰	۳-۶-۵ روش یادگیری ماشین
۴۰	۴-۶-۵ روش LSA
۴۰	۵-۶-۵ روشی برای بدست آوردن خلاصه‌ی مفهومی
۴۱	۶-۶-۵ روش خلاصه‌سازی متن با استفاده از شبکه‌های عصبی
۴۲	۷-۶-۵ روش خلاصه‌سازی مبتنی بر منطق فازی
۴۳	۸-۶-۵ خلاصه‌سازی استخراجی چند سندی
۴۴	۹-۶-۵ خلاصه‌سازی مبتنی بر پرس‌وجو
۴۵	۱۰-۶-۵ خلاصه‌سازی متون چند زبانه
۴۵	۷-۵ خلاصه فصل
۴۶	ابزارهای پردازش زبان طبیعی
۴۷	۱-۶ مقدمه
۴۷	۲-۶ مهم‌ترین ابزارهای پردازش زبان طبیعی در متون
۴۷	۱-۲-۶ تشخیص دهنده‌ی جمله (SENTENCE SPLITTER)
۴۷	۲-۲-۶ TOKENIZER
۴۸	۳-۲-۶ ریشه‌یابی (STEMMING)
۴۹	۶-۲-۴ POS Tagger
۵۱	۵-۲-۶ NAMED ENTITY RECOGNITION
۵۱	۶-۲-۶ WORD-NET
۵۲	۷-۲-۶ SIMILARITY RECOGNITION
۵۲	۸-۲-۶ PARSER
۵۲	۹-۲-۶ CHUNKER
۵۲	۱۰-۲-۶ SEMANTIC ROLE LABELER
۵۳	۱۱-۲-۶ ANNOTATOR
۵۳	۱۲-۲-۶ COREFERENCE RESOLUTION
۵۳	۳-۶ خلاصه فصل
۵۴	تحلیل سیستم خودکار کلمات کلیدی متون زبان فارسی
۵۵	۱-۷ مقدمه
۵۵	۱-۱-۷ مسائل و چالش‌های پردازش متن فارسی
۵۵	۲-۱-۷ پردازش لغوی
۵۶	۳-۱-۷ پردازش ساخت واژی

- ۵۶..... ۴-۱-۷ تهیه منابع زبانی
- ۵۷..... ۲-۷ بازیابی اطلاعات و استخراج کلمات کلیدی
- ۵۷..... ۱-۲-۷ بازیابی اطلاعات
- ۵۷..... ۲-۲-۷ تئوری لان
- ۵۹..... ۳-۲-۷ قانون ZIPF
- ۵۹..... ۴-۲-۷ کلمات کلیدی
- ۶۰..... ۱-۴-۲-۷ تقسیم‌بندی روش‌ها
- ۶۰..... ۱-۱-۴-۲-۷ تقسیم‌بندی ابزاری
- ۶۰..... ۲-۱-۴-۲-۷ تقسیم‌بندی تکنیکی
- ۶۱..... ۲-۴-۲-۷ مراحل استخراج کلمات کلیدی
- ۶۱..... ۱-۲-۴-۲-۷ حذف کلمات عمومی فارسی
- ۶۳..... ۲-۴-۲-۷ ریشه‌یابی
- ۶۳..... ۳-۲-۴-۲-۷ وزن دهی به کلمات
- ۶۳..... ۱-۳-۲-۴-۲-۷ پارامتر TF*IDF
- ۶۴..... ۲-۳-۲-۴-۲-۷ پارامتر سیگنال و نویز
- ۶۵..... ۳-۳-۲-۴-۲-۷ پارامتر مقدار تمایز
- ۶۷..... ۳-۴-۲-۷ پارامترها
- ۶۷..... ۱-۳-۴-۲-۷ پارامتر دربرگیری
- ۶۷..... ۲-۳-۴-۲-۷ پارامتر تعیین کنندگی
- ۶۸..... ۴-۴-۲-۷ داوری مبتنی بر کارشناس انسانی
- ۶۸..... ۵-۴-۲-۷ داوری مبتنی بر سیستمهای بازیابی اطلاعات
- ۶۸..... ۱-۵-۴-۲-۷ مقدار بازخوانی
- ۶۸..... ۲-۵-۴-۲-۷ مقدار دقت
- ۶۹..... ۳-۵-۴-۲-۷ پارامتر FMEASURE
- ۶۹..... ۴-۵-۴-۲-۷ پارامتر FALLOUT
- ۶۹..... ۵-۲-۷ روشهای آماری
- ۷۱..... ۳-۷ دشواریهای ریشه‌یابی فارسی و معرفی روشهایی برای ریشه‌یابی فارسی
- ۷۱..... ۱-۳-۷ قالبهای گوناگون پرونده‌های رایانه‌ای
- ۷۲..... ۲-۳-۷ استاندارد خط در رایانه
- ۷۳..... ۳-۳-۷ دستور خط فارسی
- ۷۴..... ۱-۳-۳-۷ «ی» پس از «ه»

۷۴ «ها»ی نشانه جمع
۷۴ فاصله گذاری
۷۴ کلمات مرکب
۷۴ حرکت گذاری در نوشتار فارسی
۷۵ دگرگونی در کلمه ها هنگام پیوند
۷۵ کلمات زبان های دیگر فارسی
۷۵ شناسایی ریشه فعل ها
۷۵ روش های ریشه یابی
۷۵ ریشه یاب های جدولی
۷۵ ریشه یابی به کمک روش های آماری
۷۶ ریشه یابی به کمک روش PORTER یا شبیه به آن
۷۶ ریشه یاب های کار شده در زبان فارسی
۷۶ روش های جست و جو
۷۶ مقدمه
۷۶ جست و جوی دو ارزشی
۷۷ تعمیم جست و جوی دو ارزشی فازی
۷۷ مدل آستانه ای
۷۷ توابع تطبیق
۷۷ روش دو ارزشی تعمیم یافته
۷۸ مدل فضا برداری
۷۸ خلاصه فصل
۷۹ نتیجه گیری
۸۰ مراجع

فهرست جدول‌ها

۸	جدول ۱-۲
۳۲	جدول ۱-۴
۶۲	جدول ۱-۷
۶۲	جدول ۲-۷
۶۶	جدول ۳-۷
۶۸	جدول ۴-۷

فهرست شکل‌ها

۹.....	شکل ۱-۲.....
۲۲.....	شکل ۱-۳.....
۳۵.....	شکل ۱-۵.....
۴۱.....	شکل ۲-۵.....
۴۲.....	شکل ۳-۵.....
۴۲.....	شکل ۴-۵.....
۴۳.....	شکل ۵-۵.....
۵۸.....	شکل ۱-۷.....
۶۱.....	شکل ۲-۷.....
۷۱.....	شکل ۳-۷.....

چکیده پایان نامه

هر تکنولوژی جدیدی که عرضه می‌شود سنگ بنایش شناخت دقیق علوم حاکم بر آن می‌باشد. پژوهشگر با شناخت کامل به تسط در حوزه تحقیقاتی خود می‌رسد و پس از تخصص ایده افکار فرد را در برمی‌گیرد و در نهایت منجر به تولید و نوآوری می‌شود. تحقیقی که پیش رو دارید اصول تکنولوژی خلاصه‌سازی خودکار را به طور جامعی بیان می‌کند. تمام شاخه‌ها و زیر شاخه‌ها به دقت بررسی شده و تقریباً تمامی محاسبات لازم برای شروع طراحی یک سیستم در اختیار مخاطب قرار می‌گیرد.

شما با مطالعه هر فصل با بنیان خلاصه‌سازی خودکار آشنا می‌شوید و می‌آموزید این فناوری پیشرفته در مراحل اولیه چه طور خود را با علوم مختلف دیگر پیوند می‌زند و شما را وادار می‌کند برای آشنایی با اصول آن علوم مختلف دیگر را هم بیاموزید. ما تحقیق را بگونه‌ای پیش می‌بریم که مخاطب ابتدا اصول کلی را بیاموزد و سپس مشکلاتی که بر سر راه طراحی این فناوری می‌باشد را دریابد، چون این فناوری پیوندی ناگسستنی با زبان شناسی دارد در نتیجه مشکلات آن در زبانی همانند انگلیسی با زبانی همانند فارسی متفاوت می‌باشد و به ناچار باید با مشکلات هر کدام آشنا شود، کاری که در این تحقیق به صورت مفصل دنبال کردیم و آن را در اختیار مخاطب قرار دادیم.

یکی از چالش‌های بزرگ در این حوزه علمی چگونگی امتیازدهی به جملات برای تعیین قسمت‌های مهم یک متن می‌باشد که در نهایت متن منتخب به صورت خلاصه به کاربر تحویل داده می‌شود. البته هنوز هم این چالش وجود دارد و همه پژوهشگران در سرتاسر جهان در حال مطالعه و نقد و بررسی این چالش هستند که هرچه می‌گذرد سیستم‌هایی با عملکرد بهتر عرضه می‌شود. در این تحقیق مهم‌ترین و کاربردی‌ترین محاسبات مورد نقد و بررسی قرار گرفت و به طور مفصل با توضیح هر فرمول در تحقیق ثبت شد تا مخاطب به راحتی بتواند اصول یاد گرفته در فصل‌های اولیه را با این فصل تلفیق کند و بنیان مدل سازی سیستم خلاصه‌ساز خودکار را به صورت علمی و عملی بیاموزد. این تحقیق تنها به سیستم‌های خلاصه‌ساز خودکار فارسی نمی‌پردازد بلکه سیستم‌هایی که به زبان فارسی وجود دارد را بررسی می‌کند که نمونه بارز آن ققنوس و پارسینا می‌باشد سپس هر دو مورد نقد قرار می‌گیرند و مزایا و معایب هر کدام برای مخاطب به سادگی همراه با آمار بیان می‌شود و بعد از اینکه سیستم‌های موجود فارسی توضیح داده شد قدرتمندترین سیستم‌های خارجی نیز مورد بررسی قرار می‌گیرند و قسمت‌های مختلف هر کدام با زبانی ساده در اختیار مخاطب قرار می‌گیرد. امید می‌رود با این تحقیق بتوانیم ضعف‌های سیستم‌های فارسی را برطرف کرده هم چنین بتوانیم سیستم‌های قدرتمندتری را طراحی کنیم که این بزرگترین هدف و مورد انتظارترین نتیجه‌ای است که می‌توانیم برای این تحقیق در نظر داشته باشیم.

کلمات کلیدی:

خلاصه‌سازی متن ، خلاصه‌سازی استخراجی ، خلاصه‌ساز، روش آماری

فصل اول

معرفی موضوع تحقیق

۱-۱ مقدمه

علم و دانش هر روز به پیشرفت خود در راستای بی‌نهایت حرکت می‌کند و مفاهیم جدیدی به روی انسان قرن بیستم می‌گشاید به طوری که تخصص داشتن در همه حوزه های علمی از عهده هیچ کس برنمی‌آید و جمع‌آوری علوم در یک کتاب به هیچ عنوان میسر نیست، بلکه هر علمی میلیون‌ها کتاب مخصوص خود را می‌طلبد، که در نتیجه بیانگر دنیایی بسیار گسترده است. اکنون در عصر حاضر یافتن پاسخ سوالات گوناگون از علوم متفاوت چگونه است؟

با ظهور الکترونیکی کردن علوم این کار تا حد نسبتاً زیادی آسان شد. ظهور اینترنت نیز کمک شایانی به توسعه علوم کرد و اکنون با جست‌وجو در اینترنت می‌توان به کلیه دانش‌های سرتاسر این کره خاکی دست یافت.

اما آیا همین مقدار کافی است؟ پاسخ بسیار روشن است. یافتن مقالات متنوع صرفاً ما را در هنگام جستجوی یک مطلب خاص یاری نمی‌کند زیرا مطالعه مطالب متفاوت برای دستیابی به پاسخ مطلوب کاری بسیار زمان‌بر است و در مواردی هم مفید واقع نمی‌شود، زیرا ممکن است پاسخ خود را در آن مقاله بدست نیآوریم و تنها زمان زیادی صرف مطالعه آن کرده باشیم.

آیا راهکاری برای این مشکل می‌توان یافت؟

مشکل فوق را با طرح جدیدی که به تازگی مطرح شده می‌توان برطرف کرد. هرچند که هنوز این ایده به طور کامل به منتهای کمال خود نرسیده ولی پیشرفت‌های چشمگیری را به جهان ارائه داده است. این طرح خلاصه‌سازی اتوماتیک می‌باشد که مطلب مورد نظر را به صورت خودکار خلاصه‌سازی شده در دسترس مخاطبین قرار می‌دهد. مهمترین مزیت استفاده از خلاصه‌سازی، کاهش زمان خواندن متن است. یک خلاصه خوب، باید موضوعات گوناگون یک سند را بدون داشتن افزونگی بیان کند. ابزارهای خلاصه‌سازی می‌توانند برای تشخیص عناوین و موضوعات کلیدی یک متن مورد استفاده قرار گیرد. از جمله کاربردهای خلاصه‌سازی متن می‌توان به خلاصه‌سازی پرونده‌های پزشکی بیماران، سرویس‌های صوتی برای ناشنویان، بازیابی اطلاعات، خلاصه‌سازی فایل‌های صوتی و تصویری و کاربردهای دیگر اشاره کرد.

۲-۱ اهمیت و ضرورت تحقیق

امروزه خلاصه‌سازی اتوماتیک مورد بررسی و تجزیه و تحلیل‌های زیادی قرار می‌گیرد. نمونه‌های مطرحی در جهان وجود دارد که به عنوان نسخه‌های اولیه طرح مذکور خودنمایی می‌کند با این که این نمونه محصول مایکروسافت می‌باشد و Microsoft word 's Anfo Summarize است و قدرت نسبتاً زیادی دارد اما هنوز نتوانسته است رضایت کامل را به دنبال داشته و انتظارات را تحقق بخشد. در این تحقیق هدف ما آشنایی فارسی‌زبانان با این طرح جدید و بسیار کارآمد می‌باشد تا بلکه بتوانیم در آینده خلاصه‌ساز اتوماتیک فارسی را راه‌اندازی کنیم زیرا سیستم‌های موجود در نسخه زبان انگلیسی نیز مشکلات بسیار زیادی دارند و می‌دانیم پیاده‌سازی این سیستم به زبان فارسی مشکلاتی به مراتب بیشتر از نسخه زبان انگلیسی خواهد داشت. این بدین معناست که باید با دقت بیشتری با این تکنولوژی آشنا شده و نکات علمی آن را فراگیریم تا بتوانیم از این تکنولوژی به نحو احسن استفاده کنیم. لازم به ذکر است که کمپانی‌های اروپایی تمرکز خود را بر روی ایجاد خلاصه‌ساز اتوماتیک انگلیسی و فرانسوی قرار داده‌اند و این ما هستیم که باید با تحقیق و بررسی این سیستم را به زبان فارسی آماده کنیم.

آماده‌سازی این سیستم می‌تواند تاثیر بسیار زیادی در رشد علمی کشور به همراه داشته باشد. با داشتن بستری مناسب برای محققین و پژوهشگران راه رسیدن به پاسخ سوالات راهورتر شده و هر کسی می‌تواند با صرف زمانی بسیار اندک به مطلب موردنظر خود دست یابد و لازم نیست ساعتها وقت برای جستجوی مطلب موردنظر خود بگذارد. یکی از دلایل پیشرفت کشورهای غربی توجه بسیار زیاد آنها به مدیریت زمان و راه‌های کاهش صرف زمان‌های کم کاربرد است. اتفاقاً همین دلیل هم باعث توجه مهندسين نرم افزار و هوش مصنوعی به سیستم خلاصه‌ساز اتوماتیک شده است.

۳-۱ سوالات تحقیق

سوالاتی که در این تحقیق به آنها پاسخ داده می‌شود از این قرار است:

- ۱- متن در یک سیستم الکترونیکی چگونه مورد بررسی قرار می‌گیرد؟
- ۲- چگونه می‌توان قسمت‌هایی از متن را جداسازی کرد؟
- ۳- چگونه می‌توان قسمت‌های مهم متن را تشخیص داد؟
- ۴- چگونه می‌توان قسمت‌های مهم متن را از متن کامل جدا کرد؟
- ۵- چه ویژگی‌هایی را می‌توان استخراج کرد؟

۴-۱ فناوری تعریف شده در تحقیق

هدف ما در این تحقیق ارائه مطالبی سودمند برای مهندسين می‌باشد. به همین دلیل سعی کردیم بهترین و معروف‌ترین خلاصه‌سازهای جهان را مورد نقد و بررسی قرار دهیم. مزایا و معایب هر یک را به صورت علمی بررسی می‌کنیم. نقاط ضعف و قوت آنها را به مخاطبینمان معرفی می‌کنیم و بگونه‌ای مطالب تحقیق را توسعه می‌دهیم که پس از پایان مطالعه این تحقیق هر نوع خلاصه‌سازی که به شما معرفی شد به راحتی فرآیند اجرای آن را بشناسید و آن را از لحاظ علمی به چالش بکشید. هدف نهایی ما نیز این است که با آشنایی شما با سیستم‌های فوق گام‌های خود را جهت طراحی سیستمی مبتنی بر ساپورت زبان فارسی بردارید.

این سیستم‌ها شامل:

- Microsoft word ' s Anfo Summarize
- British telecom prosum
- FRUMP
- Auto Summerized
- GistaSumm
- FarsiSum

لازم به ذکر است در تحقیق به خلاصه‌ساز FarsiSum به عنوان خلاصه‌ساز فارسی نیز توجه داشتیم و آن را نیز مورد بررسی قرار دادیم.

۵-۱ خلاصه فصل

با گسترش روزافزون حجم اطلاعات موجود در وب و افزایش چشم گیر مقالات منتشر شده در زمینه‌های مختلف علمی، دسترسی درست و مطالعه اطلاعات مورد نیاز، همواره یکی از مشکلات محققان و پژوهشگران قرن ۲۱ می‌باشد. اینکه چه طور از یک طرف با این حجم انبوه از داده‌ها و از طرفی دیگر با زمان محدودی که در اختیار داریم، بتوانیم مطالب مورد نیاز خود را مطالعه کنیم و یا اینکه چه طور می‌توان در روز چندین کتاب را مطالعه نمود و یا اینکه آیا می‌توان سیستمی طراحی نمود که بتواند با داده‌های موجود به تمامی سوالات ما پاسخ دهد، اینها سوالاتی است که پاسخ آنها را می‌توان در یک سیستم خلاصه‌ساز متن جستجو کرد.