

الآية

قال تعالى :

(قَالُوا سُبْحَانَكَ لَا عِلْمَ لَنَا إِلَّا مَا عَلَّمْتَنَا إِنَّكَ أَنْتَ الْعَلِيمُ الْحَكِيمُ)

سورة البقرة الآية (32)

الحمد لله

الحمد لله رب العالمين، أعطى اللسان، وعلم البيان، وخلق الإنسان، فبأي الأء ربكما تكذبان .. لك الحمد يا من هو للحمد أهل، أهل الثناء والمجد، أحق ما قال العبد وكلنا لك عبد، لك الحمد مادعونك إلا حسن ظن بك ومارجوناك إلا ثقة فيك، وماخفناك إلا تصديقاً بوعدك ووعدك لك الحمد حمداً كثيراً طيباً مباركاً فيه، وصلى الله على سيدنا محمد خاتم الأنبياء والمرسلين أجمعين بشر وأنذر ووعد وأوعد، أنقذ الله به البشر من الضلالة وهدى الناس الى صراط المستقيم، صراط الله الذي له مافي السموات ومافي الأرض الا الى الله تصير الأمور .

الإهداء

بدأنا بأكثر من يد وقاسينا أكثر من هم وعانينا الكثير من الصعوبات وهانحن اليوم والحمد لله نطوي سهر الليالي
وتعب الأيام وخالصة مشوارنا بين دفتي هذا العمل المتواضع إلى منارة العلم والإمام المصطفى إلى الذي علم إلى
سيد الخلق إلى رسولنا الكريم خاتم الأنبياء والمرسلين

إلى منارة العلم والإمام المصطفى إلى الأمي الذي علم إلى سيد الخلق إلى رسولنا الكريم سيدنا خاتم الأنبياء
والمرسلين محمد صلى الله عليه وسلم

إلى النبيوع الذي لا يمل العطاء إلى من حاكت سعادتي بخيوط منسوجة من قلبها إلى والدتي العزيزة

إلى من سعى وشقى لأنعم بالراحة والهناء الذي لم يبخل بشئ من أجل دفعي في طريق النجاح الذي علمني أن أرتقي
سلم الحياة بحكمة وصبر إلى والدي العزيز

إلى من حبههم يجري في عروقي ويلهج بذكرهم فؤادي إلى أخواتي وإخواني إلى من سرنا سوياً ونحن نشق الطريق
معاً نحو النجاح والإبداع إلى من تكاتفنا يداً بيد ونحن نقطف زهرة تعلمنا إلى زميلاتي وزملائي إلى من علمونا
حروفاً من ذهب وكلمات من درر وعبارات من أسمى وأجلى عبارات في العلم إلى من صاغولنا علمهم حروفاً ومن
فكرهم منارة تنير لنا سيرة العلم والنجاح إلى أساتذتنا الكرام

الشكر والتقدير

اللهم لك الحمد كله، ولك الملك كله وبيدك الخير كله وإليك يرجع الأمر كله علانيته وسره ،اللهم صل على سيدنا محمد وعلى آله وصحبه أجمعين، وعملاً بقول الرسول عليه أفضل الصلاة وأتم التسليم:

((لا يشكر الله من لا يشكر الناس))

أخرجه البخاري

لابد لنا ونحن نخطو خطواتنا الأخيرة في الحياة الجامعية من وقفة نعود إلى أعوام قضيناها في رحاب الجامعة مع أساتذتنا الكرام الذين قدموا لنا الكثير باذلين بذلك جهودا كبيرة في بناء جيل الغد لتبعث الأمة من جديد وقبل أن نمضي نقدم أسمى آيات الشكر والامتنان والتقدير والمحبة إلى الذين حملوا أقدس رسالة في الحياة إلى الذين مهدوا لنا طريق العلم والمعرفة إلي جميع أساتذتنا الأفاضل ونخص بالشكر والتقدير المشرف الدكتور : محمد المصطفى

الذي نقول له بشراك قول رسول الله صلى الله عليه وسلم:

"إن الحوت في البحر والطير في السماء ليصلون على معلم الناس الخير"
لتفضله بالإشراف على هذه الدراسة ومنحها الوقت والجهد والنصح رغم مسئولياته المتعدده فلم يبخل علينا بالتوجيه والإرشاد طوال فترة البحث وكما نخص بالشكر كل الشكر من ساعدتنا وسانددتنا لإكمال هذه الدراسه الأخت الفاضلة نسبية الهادي

المستخلص

نظراً للتطور السريع في الويب وانتشار المعلومات عليه بصورة واسعة أصبحت هنالك الكثير من التحديات والصعوبات التي تواجه أنظمة إسترجاع المعلومات وخصوصاً عند إنتشار الوثائق علي الويب بلغات مختلفة أو وجود الوثيقة الواحدة مكتوبة بأكثر من لغة(وثيقة متعددة اللغات) فكانت من التحديات كيفية فهرسة هذه الوثائق بالأخص متعددة اللغات بطريقة تؤدي إلى زيادة كفاءة الإسترجاع في أنظمة إسترجاع المعلومات لذلك تم إستخدام الطريقة المركزية الموزعة التي لها فوائد عدة خاصة في فهرسة الوثائق متعددة اللغات، لبناء فهارس أحادية اللغة للوثائق أحادية اللغة وفهرس للوثائق متعدد اللغات ثم البحث في هذه الفهارس بإستخدام إستعلامات المستخدمين وبهذه الطريقة نكون قد أتينا للمستخدم كتابة الإستعلام باللغة التي يراها مناسبة أو التعبير عن مصطلحاته باللغة التي يعرفها, ولكن عند البحث في مجموعة الفهارس تكون هنالك أكثر من قائمة مسترجعة(قائمة لكل فهرس) فكان لابد من إيجاد طريقة لدمج هذه النتائج المسترجعة من مجموعة الفهارس في قائمة واحدة تعرض للمستخدم، بحيث تكون فيها الوثائق ذات الأهمية الأكبر للمستخدم والصلة الأعلى بالإستعلام في أعلى القائمة حتى تكون النتيجة مقنعة ومفيدة للمستخدم، ولهذا تم إستخدام ثلاث خوارزميات لدمج النتائج ومقارنة نتائجها بإستخدام معيار DCG لقياس القيمة الخام) كفاءة الإسترجاع، ووجد أن خوارزمية Raw Score تعطي أفضل النتائج, وكان لطول الإستعلام (ذو اللغتين) عربي-إنجليزي) تأثير على نتائج خوارزمية CORI. لذلك نوصي بمعالجة الطول لتحسن كفاءة الإسترجاع.

ABSTRACT

Because of the rapid development and the spread of information widely in the web , there are a lot of challenges and difficulties faced by systems of information retrieval, especially when the spread of the documents on the web in different languages and the presence of a single document written in more than one language (Mixed Documents) was challenging how the indexing of these documents Especially Mixed Documents and then search in these indexes using queries of users, therefore users are allowed to write their queries in different languages and with their own way , but there will be more than a retrieved list (List each index) therefore it was necessary to find a way to merge these results which are retrieved from the indexes in a single list which will be presented to the user so that the most relevant documents to the user and to the query should occupies top of the list to make the result compelling and useful to the user, and therefore the need to use merging algorithms to merge the results and to compare the results using a standard to measure the efficiency of retrieval Discounted Cumulative Gain(DCG)

The length of a query languages influence the outcome of CORI so we recommend addressing the length which will increase the recovery efficiency.

شرح المصطلحات

الاختصار	المصطلح	شرح المصطلح
WWW	World wide web	الويب أو الشبكة العنكبوتية العالمية
IR	Information retrieval	البحث عن الوثائق، وعن داخل الوثائق، وعن المعلومات المتعلقة المعلومات الوصفية بالوثائق، بالإضافة إلى البحث في شبكة الإنترنت وقواعد البيانات
	mixed documents	وثيقة واحدة موجودة بأكثر من لغة
HTML	Hyper Text Markup Language	لغة النص المتشعب
	mixed query	الإستعلام الواحد بأكثر من لغة
	Relevant	أهمية الوثائق الوثائق بالنسبة لإستعلام المستخدم
	monolingual arabic	كتابة الإستعلام باللغة العربية
	monolingual English	كتابة الإستعلام باللغة الإنجليزية
	Indexing	جمع، تخزين البيانات لتسهيل سرعة ودقة استرجاع المعلومات.
	Matching	مطابقة الوثائق المستجدة مع

		إستعلام المستخدم
	Tokenization	عملية تقطيع تيار من نص للوصول إلى كلمات، أو رموز، أو العناصر ذات معنى تسمى القطع (tokens)
	Token	كلمة، أو رمز، أو عنصر ذو معنى ناتج من عملية تقطيع نص معين.
	Stopwords	كلمات تُستبعد قبل، أو بعد، تجهيز البيانات باللغة الطبيعية
	Normalization	عملية جعل الكلمات موحدة و متسقة إملائياً
	Stemming	الحد من تصريف الكلمة و لإيجاد أصلها أو جذرها
	Indexing	جمع، تخزين البيانات لتسهيل سرعة ودقة استرجاع المعلومات.
	Keyword	يتم تمثيل الوثائق المفهرسة في شكل مجموعة من الكلمات المفتاحية
	Term	الكلمات المفتاحية التي يتم تمثيلها في الفهرس
	Recall	نسبة تعبر عن عدد الوثائق المسترجعة من مجموعة الوثائق الكلية
	Precision	نسبة تعبر عن عدد الوثائق المسترجعة ذات الصلة من مجموعة الوثائق المسترجعة

	Ranking	ترتيب الوثائق المسترجعة تنازلياً بناءً على مدى ارتباطها بطلب المستخدم
	Term Frequency	عدد مرات ظهور المصطلح في الوثيقة
	Weighted	وزن الوثيقة
	Ranked Retrieval Model	إسترجاع يقوم على تقدير درجة أهمية الوثائق وتحديد أي منها هو المطابق للإستعلام
	Vector Space Model	هو نموذج يمثل كل من الوثائق والإستعلامات في شكل ناقل
	Similarity	المشابهة بين الإستعلام والوثائق
IDF	Inverse Document Frequency	يستخدم لتحديد أهمية المصطلح في مجموعة الوثائق
	Probabilistic Retrieval Model	هو ترتيب الوثائق على حسب الإحتمالية المقدره لهذه الوثيقة بأن تكون ذات صلة بالإستعلام
	Probabilty	إحتمالية أن الوثائق التي تم إسترجاعها تنتمي للوثائق ذات الصلة
	Best Matches ²⁵	أشهر خوارزميات نموذج الإسترجاع الإحتمالي
DCG	Discounted Cumulative Gain	معياري لكفاءة إسترجاع المعلومات في

		محركات البحث
	Google translator API	تقنية تستخدم للترجمة

فهرس الأشكال

رقم الصفحة	موضوع الشكل	رقم الشكل
5	العمليات الأساسية في نظام إسترجاع المعلومات	1.2
12	المراحل الأساسية في تمثيل المعلومات	2.2
14	مفهوم الإسترجاع ثنائية اللغة	3.2
15	أنظمة إسترجاع المعلومات ذات اللغات المتعددة MLIR	4.2
16	يوضح طريقة الفهرسة المركزية	5.2
17	الفهرسة الموزعة	6.2
18	الطريقة المركزية الموزعة	7.2
22	إستخدام الفهرسة التي تدمج بين الموزعة والمركزية	1.3
30	العمليات التي تتم في محرك بحث لوسين	2.3
31	الصفحة قبل إستخدام محلل الجريكو	1.4
32	يوضح عملية تعريف الTag	2.4
32	الصفحة بعد إستخدام الجريكو	3.4
34	عملية الحث ودمج النتائج في الفهارس المركزية الموزعة	4.4

37	نتائج الخوارزميات في الثلاث خوارزميات	2.5
----	---------------------------------------	-----

فهرس الجداول

رقم الصفحة	موضوع الجدول	رقم الجدول
33	السوابق والواحق في Light 10 stemmer	1.4
33	مكونات وثيقة لوسين	2.4