

## 1.1 المقدمة

تأتي قضية تجهيز ومعالجة المحتوى على الشبكة العنكبوتية العالمية (World Wide Web) من أبرز القضايا التي تؤرق مجتمعات المعلومات في الوقت الراهن، وينطوي هذا الأمر على عدد من القضايا الفرعية الأخرى كقضية نشر المحتوى وفعاليته وطرق إكتشافه ونظم إدارته ولكن تبرز قضية هامة على صعيد هذا الأمر وهي قضية إسترجاع المحتوى فلا توجد جدوى من وجود المحتوى إن لم يتم إسترجاعه وإستثماره من قبل المستخدمين.

وعليه تعد أدوات البحث والإسترجاع وعلى رأسها محركات البحث في الويب (Web Search Engines) بمثابة حجر الأساس لهذا المحتوى المعلوماتي، وحلقة الوصل بين طرفي النشر والإسترجاع للمحتوى.

وعلما بأن المعلومات الموجودة على الشبكة العنكبوتية أصبحت متواجدة بلغات متعددة ومختلفة وأن الوثائق الإنجليزية هي المهيمنة على الشبكة العنكبوتية ، ظهرت أهمية انظمة إسترجاع المعلومات ثنائية اللغة (CLIR) وفيها يقوم المستخدم بكتابه الإستعلام بلغة ويتم إسترجاع النتيجة بلغة أخرى خصوصاً في البلدان التي لا تنطق بالإنجليزية . وبما أن المحتوى المعلوماتي على الويب أصبح متواجد بأكثر من لغة كان لابد من السماح للمستخدم بكتابه الإستعلام بأي لغة وإسترجاع النتائج بأكثر من لغة ومن هنا ظهرت أهمية اأنظمة إسترجاع المعلومات متعددة اللغات (MLIR) ، وبما أن النتائج المسترجعة تكون من أكثر من مجموعة كان لابد من دمج هذه النتائج في قائمة واحدة علي حسب أهميتها للإستعلام أو صلتها به .

ومع تطور المحتوى المعلوماتي على الشبكة العنكبوتية أصبحت الوثيقة الواحدة موجودة بأكثر من لغة (mixed documents) فكان لابد من إيجاد طريقة لفهرسة هذه الوثائق حتى نتمكن من البحث فيها بطريقة فعالة ومن هنا كان لابد من تمكين المستخدم من كتابة الإستعلام الواحد بأكثر من لغة(mixed query) وأن تكون النتائج المسترجعة من البحث بهذا الإستعلام تضم الوثائق الأعلى صلة بالإستعلام بغض النظر عن لغة الوثيقة.

## 2.1 مشكلة البحث

من الصعوبات التي تواجه أنظمة إسترجاع المعلومات هي كيفية فهرسة الوثائق خصوصاً متعددة اللغات(التي تكون مكتوبة بأكثر من لغة) بطريقة مثل لتحسين كفاءة الإسترجاع ، وأيضاً عندما يقوم المستخدم بكتابه إستعلام بلغتين فإن ترتيب الوثائق المسترجعة من عملية البحث يتم على أساس المطابقة بين الإستعلام والوثيقة وليس على أساس أهميتها(relevant) بالنسبة للمستخدم حيث أن الوثائق ذات اللغات المتعددة تكون ذات وزن أعلى من الوثائق ذات اللغة الواحدة، مع معرفة أنه من المحتمل أن تكون هنالك وثائق ذات وزن أقل ولكنها أعلى صلة بالإستعلام من الوثائق متعددة اللغات، وبالتالي تكون النتيجة المسترجعة غير مقنعة للمستخدم، مثل لذلك إذا كان المستخدم يبحث عن (مفهوم ال deadlock) يقوم محرك البحث بتحويل الإستعلام إلى اللغة العربية (monolingual arabic) مرة بحيث يصبح (مفهوم الاقفال) وأخرى إلى الإنجليزية (monolingual English) (deadlock concept) وتصبح (deadlock concept) وإستعلام ثالث بدمج الإستعلام باللغتين العربية والإنجليزية (مفهوم الاقفال deadlock concept)، ويتم البحث بهذا

الإستعلام في مجموعة الوثائق العربية ويتم تجاهل وزن الجزء الإنجليزي من الإستعلام (أي وزنه يساوي صفر)، والبحث مرة أخرى في مجموعة الوثائق الإنجليزية وفي هذه الحالة يتم تجاهل وزن الجزء العربي من الإستعلام (أي وزنه يساوي صفر) ، والبحث أخيراً بالإستعلام في مجموعة الوثائق ذات اللغتين وهنا تكمن المشكلة في كيفية دمج النتائج المسترجعة من مجموعة الوثائق العربية والوثائق الإنجليزية والوثائق متعددة اللغات بناء على أهميتها بالنسبة للمستخدم.

## 3.1 أهمية البحث

رغم أن استخدام محركات البحث إنתרنر بصورة أكبر من السابق، وأن معظم المستخدمين لا يتجاوزوا استخدام أول صفحتين من نتائج محركات البحث التي تعرض محتوى العنكبوبية ، مردود هذا الأمر يعود إلى عدم تحقيق التطابق بين المحتوى المطلوب وبين المحتوى المسترجع من قبل محركات البحث ، و من الصعوبات والتحديات التي كانت سبباً ودافعاً لدراسة التحديات التي تواجه محركات البحث في إسترجاع المحتوى هو أن الوثائق أصبحت متواجدة بأكثر من لغة، وأن الوثيقة نفسها يمكن أن تكون مكتوبة بأكثر من لغة (خصوصاً الوثائق غير الإنجليزية) و كان لابد لمحركات البحث المستقبلية أن تضع في اعتبارها أن المستخدم قد يكون غير قادر على التعبير عن مصطلح ما يريد البحث عنه أو لا يعرف معنى التعبير بلغته(خصوصاً في البلدان التي لا تنطق بالإنجليزية) فيقوم بكتابة الإستعلام بأكثر من لغة مثلاً إذا كان يريد البحث عن "مفهوم ال network " ولا يعرف معنى كلمة network بالعربية فيقوم بكتابة الإستعلام كالتالي " مفهوم ال network " .

لذا كان لا بد من السعي نحو إيجاد طريقة تُمكِّن المستخدم من كتابة الإستعلام بأكثر من لغة و إسترجاع النتائج ودمجها على أساس صلتها بالمستخدم وليس على أساس المطابقة أو الوزن حيث تكون النتائج النهائية ذات صلة بطلب المستخدم .

## 4.1 أهداف البحث

1. تطبيق خوارزميات لدمج النتائج وإنتاج قائمة واحدة من الوثائق المرتبة على حسب أهميتها للمستخدم .
2. زيادة كفاءة النظام بزيادة كفاءة عملية الإسترجاع من الفهارس .
3. السماح للمستخدم الذي لا يستطيع التعبير عن المصطلحات التي يريد البحث عنها بأن يكتب الاستعلام بأي لغة واسترجاع نتائج تكون مقتنة وذات أهمية بالنسبة للمستخدم.
4. بناء وتطوير معمارية جديدة لفهرسة الوثائق(indexing).

## 5.1 منهجية البحث

في المشروع قيد الدراسة ستتم فهرسة الوثائق بإستخدام طريقة الفهرسة المركزية الموزعة وبناء عدة فهارس والبحث فيها بإستخدام الإستعلام ومن ثم تطبيق خوارزميات عدة لدمج النتائج المسترجعة من الفهارس المختلفة في قائمة واحدة .

## **6.1 حدود البحث**

سنتناول في هذه الدراسة أنظمة إسترجاع المعلومات وبالأخص أنظمة إسترجاع المعلومات ذات اللغات المتعددة بهيكلية فهرسة تتمجج بين طريقة الفهرسة المركزية والفهرسة الموزعة وقبل دمج الوثائق الناتجة من عملية البحث يتم تطبيق وزنها على حسب الخوارزمية التي أستخدمت في عملية الدمج .

## **7.1 هيكلية البحث**

- يتضمن البحث بالإضافة إلى هذا الباب الأبواب :
- الباب الثاني والذي يتضمن نبذة عامة عن إسترجاع المعلومات حيث يحتوي على فصلين هما المقدمة بالإضافة إلى الدراسات السابقة.
  - الباب الثالث و يتضمن منهجية البحث التقنيات والأدوات المستخدمة في البحث.
  - الباب الرابع يتحدث عن تطبيق المشروع المقترن و خطوات حل المشكلة.
  - الباب الخامس يتضمن الخاتمة التي تحتوي على النتائج والتوصيات والمرجع.

## 1.2 استرجاع المعلومات (Information Retrieval)

تعني عملية استرجاع المعلومات استرجاع الوثائق (documents) التي تطابق (matching) طلب المستخدم (query) وذلك حسب صلتها به.

وبذلك فإن عملية استرجاع المعلومات تحوي عدة عمليات يقوم بها نظام استرجاع المعلومات قبل البدء في مطابقة طلب المستخدم مع الوثائق التي يجب استرجاعها، وتلك العمليات هي :

### 1. التقطيع (Tokenization)

هي عملية تقطيع المحتوى إلى كلمات ذات معنى تسمى قطع (Tokens) وهي أصغر وحدة لها معنى .

### 2. إستبعاد الكلمات غير ذات المعنى (Removing Stopwords)

هي الكلمات التي ليس لها معنى مفيد .

### 3. تطبيع حروف بعض الكلمات (Normalization)

هي عملية توحيد الكلمات التي تم تقطيعها ووضعها في الصيغة الموحدة وذلك لزيادة مدى المطابقة بين مجموعة الكلمات المقطعة .

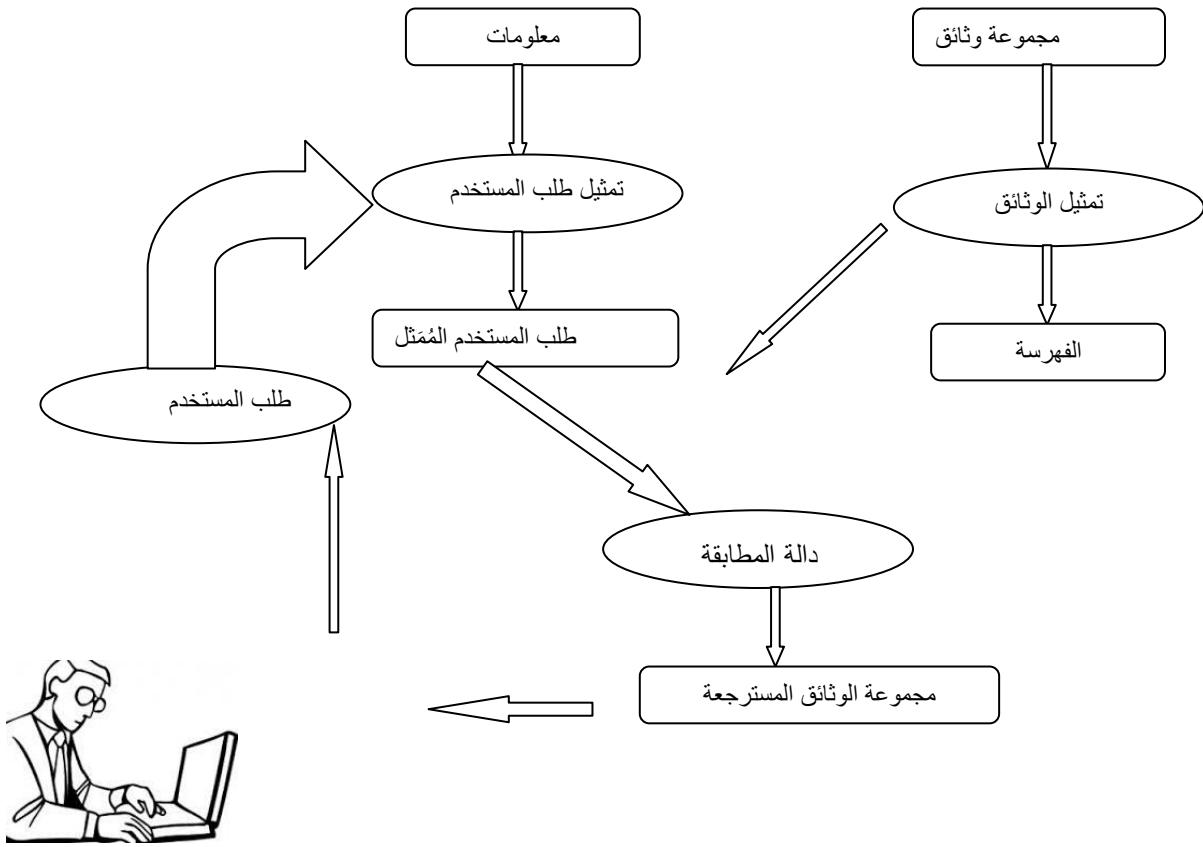
### 4. التشذيب أو تحويل الكلمات إلى أصل شترك فيه (Stemming)

هو عملية إزالة بواط الكلمة ولو احتجها وذلك لإنتاج الجذر أو الأصل .

### 5. فهرسة الوثائق بناءً على تشذيب الكلمات (Indexing)

هي العملية التي يتم فيها تمثيل الوثائق في شكل مجموعة من الكلمات المفتاحية .

وسيتم شرح هذه العمليات لاحقاً في هذا الفصل .



شكل (1.2): العمليات الأساسية في نظام إسترجاع المعلومات

ومن الرسم يتضح أن المشكلة الأساسية هي عملية مطابقة طلب المستخدم مع الوثائق التي يتم إسترجاعها.

## 1.1.2 العمليات الأساسية في نظام إسترجاع المعلومات

### 1.1.1.2 عملية التمثيل

يتم فيها وضع المحتوى في صيغة أو هيكلية محددة ويمكن تصنيف عملية التمثيل إلى :

#### أ- تمثيل الوثائق

حيث يتم وضع الوثائق المحفوظة في هيكلية تسمح بفهرستها والبحث فيها.

والفهرسة (Indexing) هي التي يتم فيها تمثيل الوثائق في شكل مجموعة من الكلمات المفتاحية (Keyword) والكلمات المفتاحية التي يتم تمثيلها في الفهرس تسمى المصطلح (Term) وهو الوحدة الأساسية المستخدمة في تمثيل الإستفسار ويمكن أن يكون المصطلح عبارة عن كلمة (Word) أو جملة (Phrase) أو جملة (Word) إعتماداً على ما يحتاجه من فهرسة الملفات، عملية إنتاج الفهرس أحياناً تمر بعدد من العمليات وذلك إعتماداً على اللغة مثلاً عملية التشذيب والقطيع (Tokenizing)، ونتائج عملية الفهرسة هو وصف منظم للملفات قابل للبحث في شكل مجموعة من المصطلحات (stemming).

#### ب- تمثيل طلب المستخدم

حيث يتم تمثيل الإستعلام (query) ووضعه في هيكلية محددة لاستخدامه في عملية البحث في الوثائق المفهرسة.

### 2.1.1.2 عملية المطابقة

يتم فيها مطابقة طلب المستخدم مع الوثائق المفهرسة وإرجاع الوثائق المطابقة وتدعى هذه الوثائق بالوثائق ذات الصلة، وترتبط تنازلياً حسب درجة صلتها بالموضوع. ومن أهم الصفات التي يجب أن تتحقق في هذه المرحلة هي الفاعلية، صحة النتائج، السرعة، قلة نسبة الخطأ، حسن الأداء في الإسترجاع، صغر حجم الفهرس، وعملية الإسترجاع لهذه الوثائق تسمى الإستدعاء أو الإرجاع (Recall) وتزداد هذه الخاصية بزيادة الوثائق المسترجعة. ومن خلال الوثائق المسترجعة يمكننا تحديد مدى نجاح نظام إسترجاع المعلومات من فشله وذلك بتحديد مدى الصلة بينها وبين طلب المستخدم وهذه تسمى الدقة (Precision) والدقة هي القدرة على تحليل جميع الكلمات المدخلة وإعطاء تقسيم صحيح واحد على الأقل، ولتحقيقها لابد من الاعتناء بأمرتين مهمين هما :

1. تحديد هل الوثيقة المسترجعة ذات صلة بموضوع طلب المستخدم أم لا.

2. تحديد ترتيب الوثائق ذات الصلة على حسب صلتها بالموضوع (Ranking)

وكفاءة عملية المطابقة تعتمد على عدة أمور:

1. حجم مجموعة الوثائق.
2. مواضع تلك الوثائق.
3. ثقافة المستخدم الذي يصبح طلبه للمعلومات.

## 2.1.2 نماذج أنظمة إسترجاع المعلومات

### 1.2.1.2 النموذج المنطقي( Boolean Model)

تصاغ فيه الإستعلامات من قبل مجموعة من المصطلحات (Term) مع العوامل المنطقية (And ,Or ,Not) وهذه العوامل المنطقية يتم التعامل معها في الإسترجاع بطريقة مشابهة لاستخدامها في الجداول الحقيقة التقليدية وبعد بناء الوثائق يتم تمثيلها في شكل مصطلحات (Term) فإذا كان المدخل مساوياً للواحد(1) يعني أن المصطلح مذكور في الوثيقة المعينة، وإذا كان مساوياً للصفر(0) فهذا يعني أنه لم يذكر في الوثيقة. النموذج المنطقي لا يأخذ في الاعتبار عدد مرات ظهور المصطلح (Term Frequency) في الوثيقة بل على حسب ظهور المصطلح في الوثيقة أو عدم ظهوره.

للنموذج المنطقي جوانب قصور وهي :

أولاً: لا يوفر وثائق لها ترتيب وزن (Ranked and Weighted) ويركز حول ما إذا كانت الوثيقة مطابقة للاستعلام أم لا.

ثانياً: الاستعلام في النموذج المنطقي معقد نسبياً.

ونلاحظ أنه يتم ترتيب النتائج زمنياً بدلاً من تقدير دقيق لدرجة أهميتها للاستعلام.

### 2.2.1.2 نموذج الإسترجاع المرتب(Ranked Retrieval Model)

هو نموذج إسترجاع يقوم على تقدير درجة أهمية الوثائق وتحديد أي منها هو الأفضل مطابقة للاستعلام ،في هذا السياق ثبت أن النماذج المرتبة أكثر فعالية من النموذج المنطقي، وممليي يصف لنا أمثلة النماذج المرتبة.

### 1. نموذج الناقلات(Vector Space Model)

هو نموذج يستمد إسمه من حقيقة أنه يمثل كل من الوثائق والإستعلامات في شكل ناقل ويتم تمثيل كل مصطلح في بين الوثيقة (similarity) الناقل في أبسط معانيه، ونجد أن الطريقة المتبعة في هذا النموذج لقياس التشابه والاستعلام هي قياس جيب تمام الزاوية بين متجه الوثيقة ومتوجه الاستعلام فإذا كانت هذه الزاوية صغيرة هذا يعني أن متجه الوثائق يعتبر ذا صلة أكبر بالنسبة لمتجه الاستعلام .

ا هو الوثيقة ويتم تمثيله إستناداً على تردد المصطلحات (j) هو المصطلح و (i) حيث أن ( $W_{ij}$ ) ويرمز للوزن هنا بالرمز

والذي يعرف بأنه عدد مرات ظهور المصطلح في كل وثيقة، أو إستناداً على (terms frequencies)، الذي يستخدم لتحديد أهمية المصطلح في Inverse Document Frequency (TF.IDF scheme) ، الذي يستخدم لمجموع الوثائق وعادة ما تستخدم الصيغة التالية :

$$\text{cosine}(dj, q) = \frac{(dj * q)}{\|dj\| * \|q\|} = \frac{\sum_{i=1}^{|V|} w_{ij} * w_{iq}}{\sum_{i=1}^{|V|} w_{ij}^2 * \sum_{i=1}^{|V|} w_{iq}^2}$$

حيث:-

$dj$ : متجه الوثيقة.

$q$ : متجه الاستعلام.

$w_{ij}$ : وزن المصطلح  $j$  في الوثيقة.

$w_{iq}$ : وزن المصطلح  $i$  في متجه الاستعلام.

هو عدد الابعاد في الناقلات وهذا هو العدد الاجمالي للكلمات الهامة.  $|V|$  :

## 2- نموذج الاسترجاع الإحتمالي (Probabilistic Retrieval Model) :

تقاس (Similarity) في فضاء الناقلات بطريقة عشوائية نوعاً ما مثلاً: النموذج يفترض أن الوثائق التي تكون قريبة في الزاوية من الاستعلام (أي الوثائق التي تكون الزاوية بينها وبين الاستعلام صغيره) تعتبر ذات صلة أكبر بالنسبة لمتجه الاستعلام. لكن في نموذج الاسترجاع الإحتمالي يتم استخدام طرق أكثر دقة بحيث يتم ترتيب الوثائق على حسب الإحتمالية المقدرة لهذه الوثيقة بأن تكون ذات صلة بالإستعلام وهذا يعتبر أساساً مبدأ نموذج الاسترجاع الإحتمالي الذي طور من قبل روبرتسون .

في نموذج الاسترجاع الإحتمالي على نظام إسترجاع المعلومات أن يقرر بأن الوثائق إما تتبع لمجموعة الوثائق ذات الصلة (relevant set) أو مجموعة الوثائق التي ليس لها صلة بالإستعلام (non relevant) . ولإتخاذ هذا القرار يفترض وجود مجموعة وثائق ذات صلة، ومجموعة وثائق ليست ذات صلة لهذا الاستعلام ومهمته هي حساب إحتمالية أن تتبع الوثيقة المعينة إلى مجموعة الوثائق ذات الصلة أو مقارنتها مع إحتمالية أن تكون الوثيقة تتبع لمجموعة الوثائق ليست ذات الصلة .

إذا كان:

D: الوثيقة

R: ذات صلة للوثيقة

NR : ليست ذات صلة للوثيقة

يمكن حساب الإحتمالية بإستخدام ما يسمى Bayes Rule

$$P(R|D) = p(D|R) * p(R) / p(D)$$

$$P(NR|D) = p(D|NR) * p(NR) / p(D)$$

P : هي إحتمالية (Probability) أن الوثيقة (Document) التي تم إسترجاعها بالنظام تتنمي للوثائق ذات الصلة (Relevant Document)

P(D|R) : بافتراض أن لدينا (R) Relevant Document فإن (P(D|R)) هي إحتمالية أن الوثيقة المسترجعة تتنمي للوثائق ذات الصلة .

P(R|D) هي إحتمالية المطابقة للوثيقة (D)

P(NR|D) الوثائق ليست ذات الصلة

وتعتبر الوثيقة D على أنها ذات صلة إذا كان :

$$P(R|D) > P(NR|D)$$

تنتمي إلى الوثائق ذات الصلة D كوزن لتحديد إمكانية أن تكون الوثيقة (P(R|D) و p(D|NR)) يتم إستخدام وتعتبر خوارزمية Best Match25 هي من أشهر خوارزميات نموذج الإسترجاع الإحتمالي المستخدم على نطاق واسع وهي الخوارزمية التي تم استخدامها في هذا البحث.

### 3.1.2 المراحل الأساسية في عملية تمثيل المعلومات

#### 1.3.1.2 التقطيع (Tokenization)

تعتبر المرحلة الأولى والأساسية في عملية تمثيل المعلومات ويتم فيها تقطيع المحتوى إلى كلمات ذات معنى تسمى قطع (Tokens) وهي أصغر وحدة لها معنى، ويتم أيضاً حذف بعض الحروف الخاصة كعلامات الترقيم والشرطيات وغيرها. الكلمات الناتجة من عملية التقطيع يمكن الاستفادة منها في عملية الفهرسة. هناك العديد من الإستراتيجيات التي يمكن أن تستخدم للتقطيع منها :

التقطيع بالفراغات

يتم اعتبار أي كلمة بعدها فراغ قطعة مستقلة، ومن عيوب هذه الطريقة أن الكلمات المقاطعة قد تكون اختصار أو علامة ترقيم أو ضمير أو كلمات موصولة مثل (عبدالله) أو مفصولة مثل (علاء الدين – هي كلمة واحدة ولكن مفصولة بفراغ) وهذا الإختلاف في القطع قد يؤدي إلى نتائج خاطئة، ويعتبر الفراغ بين الكلمات في اللغة العربية ميزة جيدة لا توجد في بقية اللغات مثل اللغات الآسيوية لذلك لا يمكن تطبيق هذا النوع من التقطيع على كل اللغات. ويتم اعتبار اي كلمة بعدها علامة (-) قطعة مستقلة، ومن عيوب هذه الطريقة أن بعض الكلمات قد تحتوي على هذه العلامة كجزء أساسي من الكلمة. فمثلاً كلمة "كريم" سيتم تقطيعها إلى قطعتين "ك" و "ريم"، إذا كان

المستخدم يقصد التشبيه بالغزال ريم، فهي تم تقطيعها بشكل جيد. لكن في حالة أنه كان يقصد صفة الكرم بهذه الطريقة تعتبر غير مجدية وستعطي نتائج خاطئة.

### 2.3.1.2 الكلمات المراد حذفها (Stopwords)

تعتبر المرحلة الثانية في عملية تمثيل المعلومات ويتم فيها حذف الكلمات التي لا تحمل معنى لذاتها والتي تعتبر كلمات ربط حروف الجر والضمائر وأسماء الموصول وغيرها . مثل لذلك: حروف الجر(Article)، by، وـ(كـ)، (an) مثل (an)، والضمائر مثل (أنت) فمثل هذه الكلمات لا يتم فهرستها. هناك كلمات لا يجب اعتبارها (Stopword) وقد تبين أنه ليس بالصحيح دائماً إزالة (a) فإنه سيتم حذف حرف(a) Vitamin a) وبالتالي حذفها مثل إذا قلنا ال (Stopwprds) وخاصة في اللغات الإعرابية مثل اللغة العربية وهي غالباً تكون مقيدة في تحديد أجزاء الكلام، ويتم حذف هذه الكلمات الزائدة قبل الترجمة.

أمثلة لبعض الكلمات الزائدة في اللغة العربية : من، إلى، هو،....، الكلمات المترجمة من اللغات الأخرى، والضمائر، وحروف الجر، العبارات مثل (السيد، العزيز،...).

### 3.3.1.2 التطبيع (normalization)

بعد تقسيم النص إلى مجموعة من الكلمات المتقطعة (Tokens) في هذه المرحلة يتم توحيد تلك الكلمات المتقطعة ووضعها في صيغة موحدة (Canonical Form) وذلك لزيادة مدى المطابقة بين مجموعة الكلمات المتقطعة، ومن طرق التطبيع التعديلات الإملائية أو مايسى بـ(الإستبدال) وهذه العملية تعتمد على اللغة فمثلاً في اللغة العربية يتم إستبدال (ي) بـ(ى)، وإستبدال (أ،إ) بـ(ا)، وأيضاً من العمليات التي تتم في هذه المرحلة عملية التعديل في الكلمات كتحويلها إلى حروف صغيره (lower case) .

### 4.3.1.2 التسذيب (stemming)

هو عملية إزالة بواي الكلمة ولو احتجها وذلك لإنتاج الجذر كما في بعض الحالات أو النابعة (الأصل) كما في حالات أخرى، وهذه الطريقة تجمع كل الكلمات التي لها نفس الأصل وتمتلك بعض العلاقات الدلالية بحيث يتم فيها تخطيط وتحويل كل الصيغ المختلفة للكلمة إلى صيغ مشتركة وشكل موحد لتلك الكلمات التي لها نفس الجذر ومثال لذلك إستخدام التسذيب في اللغة الإنجليزية حيث يتم إسترجاع كل الوثائق(documents) التي تحوي ”play,plays,player,playing“ ونثال الصيغ الناتجة يجب أن تكون ملائمة لعملية الفهرسة والبحث. وكمثال من اللغة العربية جذر كلمة ”طفيليات“ هو ” طفل“ وجذر كلمة أطفال هي أيضاً ” طفل“ ولكن الكلمتان تختلفان في المعنى .

تبُرَز أهمية هذه العملية في إجراءات التصنيف وبناء الفهارس وبحث وإسترجاع المعلومات وذلك لكونها تقلل من تأثير اختلاف الأنماط والأشكال للكلمات العربية، كما أنه يساهم في تقليل الحجم المطلوب لخزن الكلمات في حالة تخزينها في الفهارس بأشكالها المختلفة والموجودة في النصوص الأصلية.

للتشذيب تأثير كبير على الاستدعاء (recall) بمعنى أن عدد الوثائق المسترجعة كثير، وإذا كانت تلك الوثائق المسترجعة بعد عملية التشذيب ذات صله عاليه بطلب المستخدم (query) هذا يعني نجاح ودقة تلك العملية (precision). وسنقوم بتوضيح الإستدعاء والدقة لاحقاً.

### طرق التشذيب

هناك العديد من طرق التشذيب لكن أغلبها يعتمد اعتماداً كلياً على اللغة التي يتم التعامل معها و فيما يلي سيتم توضيح طريقين من طرق التشذيب :

#### أ- التشذيب الخفيف (Light Stemming)

في هذا النوع يتم تجريد الكلمات من السوابق والواحد بحيث تستخدم الكلمات المجردة لفهرسة الوثائق غالباً ما تكون الكلمات التي تجرب من نفس الكلمة لها نفس المعنى و بذلك فإن هذه الطريقة تقوم باسترجاع عدد الوثائق ذات الصلة بطلب المستخدم و نتيجة لذلك فإن هذا النوع من التشذيب يجد من نطاق البحث بصورة مؤثرة.

خصائصه :

- غير فعال في حالة وجود كلمتين لهما نفس المعنى ويختلفان في الميزان الصرفي مما يقلل من إحتمالية المطابقة و هذه المشكلة تعرف بالقصور في التشذيب (under-stemming).
- لم يساعد في التعرف على أجزاء الكلام (part of speech) أي معرفة نوع الكلمة في مجموعة الوثائق هل هي (إسم، فعل، صفة حال (الظروف الزمانية، الظروف المكانية)) فمثلاً الإسم والفعل اللذان لهما نفس الصيغة (single form) عند التشذيب يتم إسترجاع الوثائق التي وردت في كليهما.

#### ب- التشذيب بإيجاد جذر الكلمة (Root-based Stemming)

في هذا النوع يتم إستخراج جذور الكلمات الموجودة في المستند وذلك بإجراء بعض العمليات اللغوية عليها ومن ثم يتم استخدام الجذور المستخرجة ككلمات دلالية (Indexing Terms) في الفهرس. التشذيب بإيجاد جذر الكلمة قائم على الجذر حيث يقوم

بعمليات لغوية حتى يستخرج أصل الكلمة (root)، وبهتم هذا النوع من التشذيب بإستخراج جذور الكلمات الموجودة في الوثيقة وإستخدامها ككلمات دلالية (Indexing Terms)، وبذلك فإن جميع الكلمات الواردة في الوثيقة الواحدة والتي لها نفس الجذر ستفهرس بنفس الكلمة بالرغم من أنه ليس بالضرورة أن يكون لها نفس المعنى مما يؤدي إلى تقليل حجم الفهرس.

ومن خصائصه :

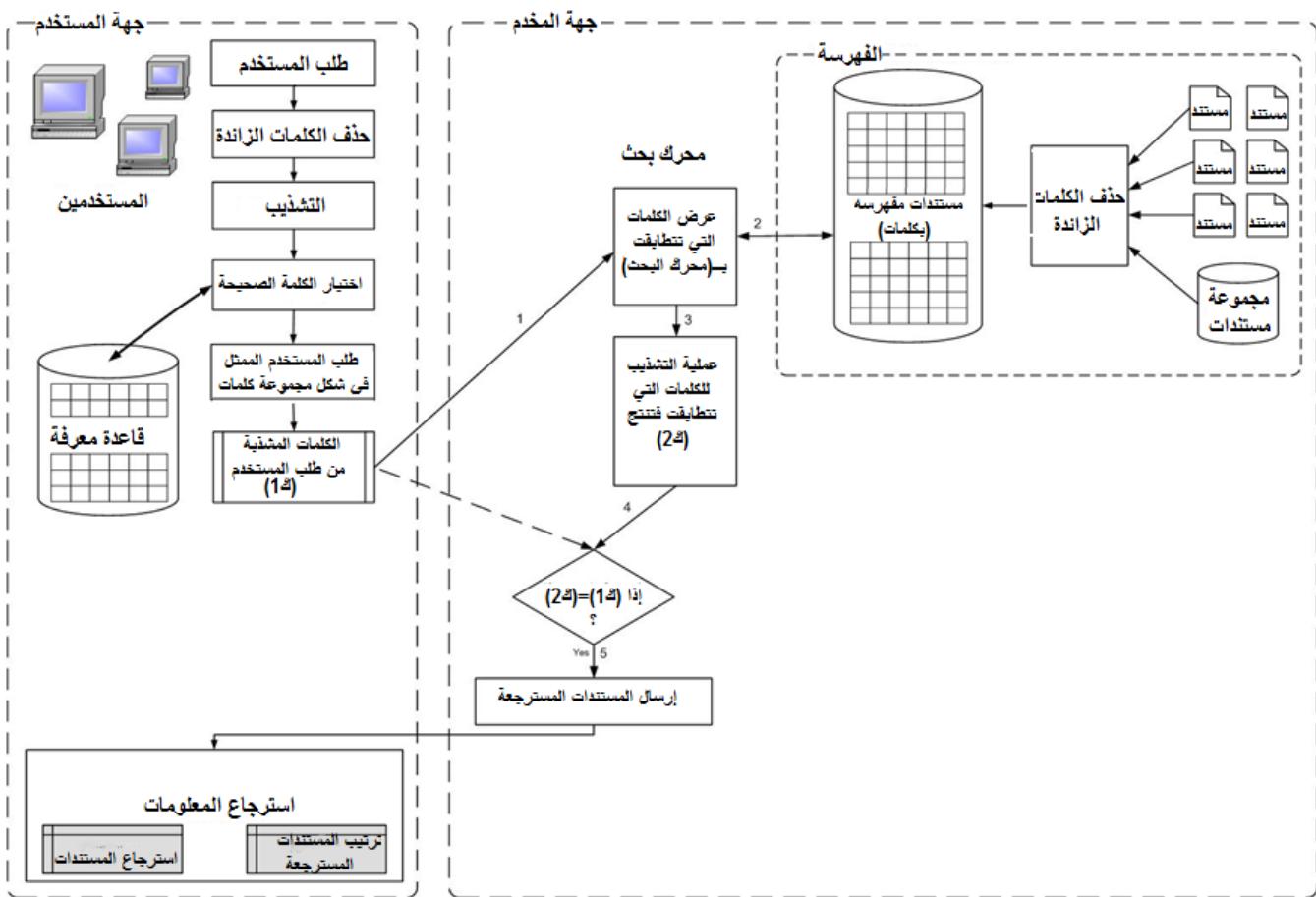
- أن هذه الطريقة تسترجع للمستخدم جميع الوثائق التي تحتوي على أي صيغة صرفية للكلمات الواردة في طلب المستخدم مما يزيد من إمكانية عثور المستخدم على المعلومة المطلوبة .

بـ- في حالة مجموعات الوثائق الضخمة جداً والمتتجدة بإستمرار كصفحات الإنترنت فإن هذه الطريقة غير مجده لأنها توسيع كثيراً من نطاق البحث ولن تصل بالمستخدم إلى المعلومة المطلوبة.

تـ- يقوم بتنقلي حجم الفهرس.

ثـ- يركز على إستدعاء أكبر كمية من الوثائق التي تشترك كلماتها في الأصل مع كلمات طلب المستخدم ولكن بعض الكلمات لها نفس الأصل وليس لها علاقة ببعضها البعض من ناحية المعنى مما أدى إلى كثرة في الإستدعاء وإسترجاع الوثائق التي يريدها المستخدم والتي لا يريدها مثلاً كلمة "يذهبان" و"ذهب"(gold) لهما نفس الجذر وهو "ذهب" ولكنهما يختلفان في المعنى، وبالتالي تؤدي هذه الطريقة من الدقة، و هذه المشكلة تعرف بالتشذيب أكثر من اللازم (over-stemming).

كل التفاصيل التي ذكرت في النقاط السابقة يمكن إجمالها في الشكل



: يوضح المراحل الأساسية في تمثيل المعلومات(2.2) الشكل

يمكن قياس أداء أنظمة إسترجاع المعلومات بطرق مختلفة اعتماداً على مهمة الإسترجاع وذلك على حسب الوثائق المسترجعة سواء أنها كانت ذات صلة أو ليست ذات صلة بطلب المستخدم لتقييم تلك الوثائق. ومن هذه الطرق:

- الإستدعاء (Recall)-1

هي نسبة تعبّر عن عدد الوثائق المسترجعة ذات صلة بالإستعلام من مجموعة الوثائق الكلية

$$Recall = \frac{\text{number of retrieved relevant documents}}{\text{number of relevant documents in the collection}}$$

- الدقة (Precition) - 2

هي نسبة تعبّر عن عدد الوثائق المسترجعة ذات الصلة من مجموعة الوثائق المسترجعة

$$Precition = \frac{\text{number of retrieved relevant documents}}{\text{number of retrieved documents}}$$

ولقياس جودة ترتيب الوثائق نستخدم مقياس (Discounted Cumulative Gain) DCG ، وهذا المقياس غالباً

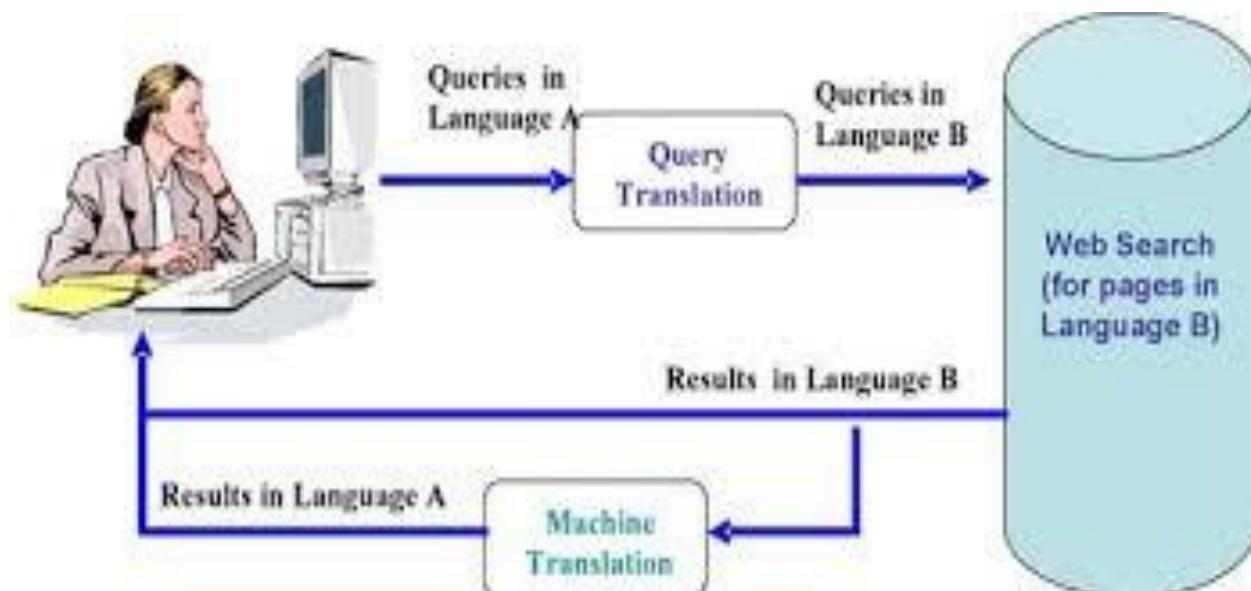
ما يستخدم لقياس فعالية خوارزميات محركات البحث أو التطبيقات المشابهة لها .

## 2.2 الدراسات السابقة

في أنظمة إسترجاع المعلومات ذات اللغة الواحدة (Monolingual IR) يتم تمثيل كل من الإستعلامات ومجموعة الوثائق بلغة واحدة، وتكون النتيجة المسترجعة بنفس اللغة ولكن طبيعة المعلومات المتاحة على الشبكة العنكبوتية أنها توجد بلغات مختلفة، فقد يحتاج المستخدم إلى كتابة الإستعلام بلغة (مختلفة عن لغة الوثيقة) ويريد النتيجة بلغة أخرى. أو كتابة الإستعلام نفسه بأكثر من لغة ومن هنا ظهر مفهوم CLIR و MLIR وفيما يلي توضيح لكل منهما .

### 1.2.2 أنظمة الإسترجاع ثنائية اللغة (CLIR)

فيها يقوم المستخدم بكتابه الإستعلام بلغة واحدة مختلفة عن لغة الوثيقة ربما لأن المستخدم لا يجيد لغة الوثيقة ولكن تكون النتيجة بلغة مختلفة عن لغة الإستعلام . وتسمى اللغة المكتوب بها الإستعلام باللغة المصدر ، وتسمى لغة الوثائق المسترجعة باللغة الثانوية (Languge Target) . حيث تتم ترجمة الإستعلام من اللغة المصدر إلى اللغة الثانوية بإستخدام أحد طرق الترجمة والبحث بهذا الإستعلام في مجموعة الوثائق لإسترجاع النتائج.



شكل (3.2) : مفهوم أنظمة الإسترجاع ثنائية اللغة

ونظراً لأن الوثيقة الواحدة على الشبكة العنكبوتية قد تكون مكتوبة بأكثر من لغة خصوصاً الوثائق المكتوبة بلغة غير الإنجليزية أو أن المستخدم يريد البحث عن بعض المصطلحات التي لا يستطيع التعبير عنها بلغته، ينتج عن ذلك كتابة المستخدم للإستعلام بأكثر من لغة ومن هنا ظهر مفهوم أنظمة إسترجاع المعلومات متعددة اللغات (MLIR) .

## 2.2.2 أنظمة إسترجاع المعلومات متعددة اللغات (Multilingual Information Retrieval)

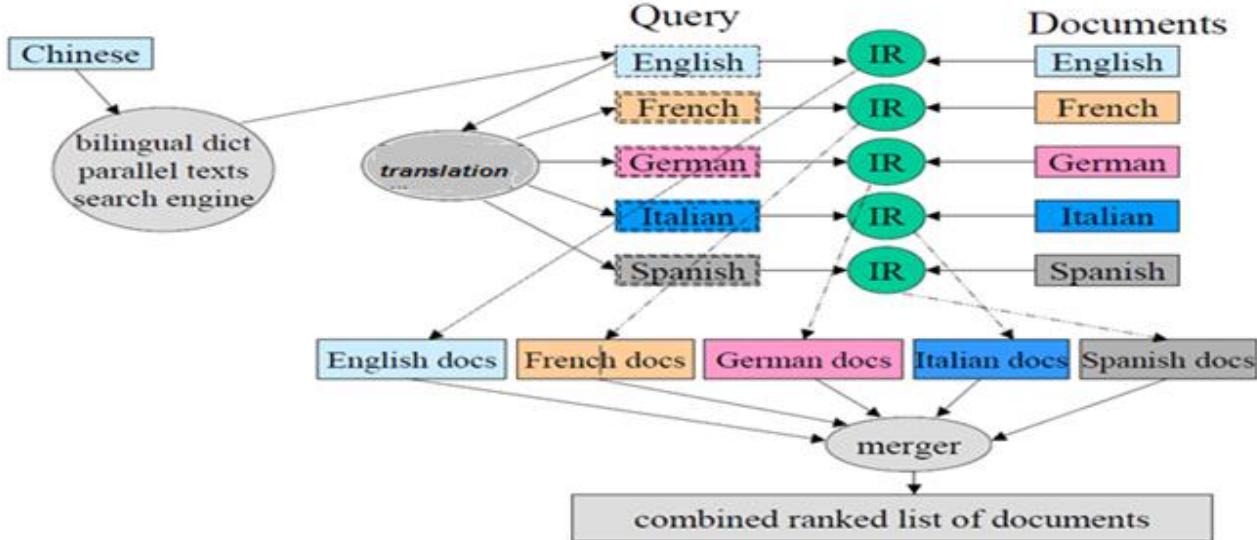
MLIR هو مهمة البحث عن الوثائق ذات الصلة في مجموعة من الوثائق سواء كانت بلغة واحدة أو متعددة اللغات لاستعلام مكتوب بأكثر من لغة، ويتم فيها ترجمة الاستعلام إلى مجموعة من اللغات والبحث بهذه الاستعلامات المترجمة فيمجموعات الوثائق (البحث بكل استعلام مترجم في مجموعة الوثائق التي لها نفس لغته) ثم دمج الوثائق المسترجعة من عملية البحث بكل استعلام وعرض قائمة وثائق موحدة الترتيب (ranked) المستخدم بغض النظر عن اللغة.

### الاستعلام متعدد اللغات

هو استعلام مكتوب بأكثر من لغة، مثلاً "مفهوم ال deadlock" ومعنى هذا الاستعلام بالعربية هو مفهوم لإफال وهذا الاستعلام مكتوب باللغتين العربية والإنجليزية غالباً ما تكون الأجزاء الإنجليزية في الاستعلام هي عبارة عن مصطلحات.

### الوثائق متعددة اللغات

هي وثائق مكتوبة بأكثر من لغة ، وفي هذه الوثيقة تكون هناك لغة أساسية ولغات أخرى ثانوية، غالباً ما تكون اللغة الثانوية هي اللغة الإنجليزية .



شكل(4.2) : أنظمة إسترجاع المعلومات ذات اللغات المتعددة(MLIR)

### 3.2.2 الطرق المتبعة في فهرسة الوثائق

كما علمنا سابقاً أن عملية الفهرسة تعني تمثيل الوثائق في شكل مجموعة من الكلمات المفتاحية ونذكر الكلمات المفتاحية التي يتم تمثيلها في الفهرس تسمى المصطلح (term) وهو الوحدة الأساسية المستخدمة في تمثيل الإستعلام والوثيقة.

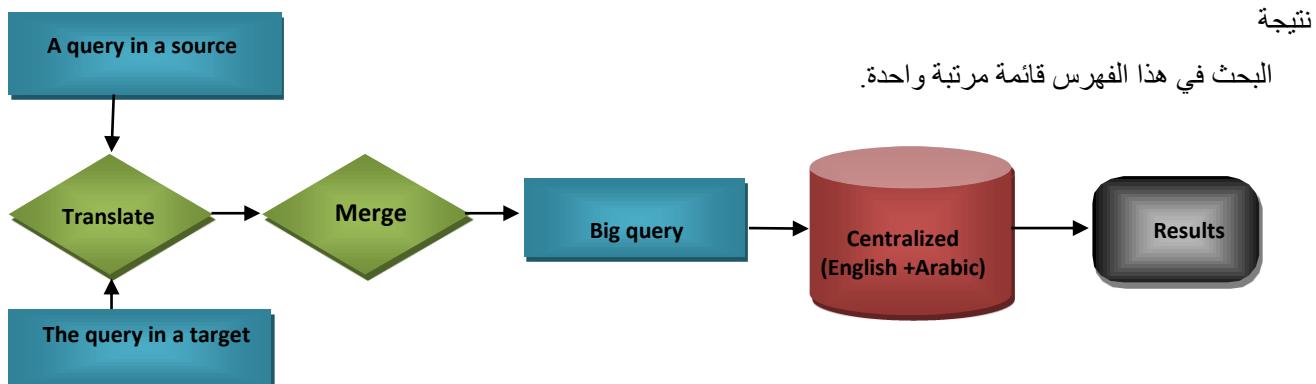
في البداية كانت أنظمة إسترجاع المعلومات تقوم بفهرسة الوثائق بطريقة مركزية وذلك بعمل فهرس واحد لجميع الوثائق سواء كانت أحادية أو متعددة اللغات، وطريقة أخرى للفهرسة هي الفهرسة الموزعة أي لكل لغة فهرس منفرد، وطريقة ثالثة تقوم بالدمج بين الفهرسة المركزية والفهرسة الموزعة لتلافي عيوب كل منها وسيتم ذكر هذه العيوب لاحقاً.

وفيما يلي شرح لطرق الفهرسة كما تم تصنيفها مسبقاً إلى:

#### 1. طريقة الفهرسة المركزية (Centralized Indexing Approach)

في هذه الطريقة يتم بناء فهرس مركزي وحيد لكل الوثائق بإختلاف لغاتها، أي تتم فهرسة الوثائق العربية والإنجليزية وذات اللغات المتعددة في ذات الفهرس. وعندما يقوم المستخدم بكتابه إستعلام معين بغض النظر عن لغته سواء كان بلغة وحيدة أو بأكثر من لغة فإنه يتم البحث بهذا الإستعلام في الفهرس المركزي الذي تم بناءه لكل الوثائق بلغاتها المختلفة، وبعد عملية البحث يتم إسترجاع النتائج المطابقة لطلب المستخدم.

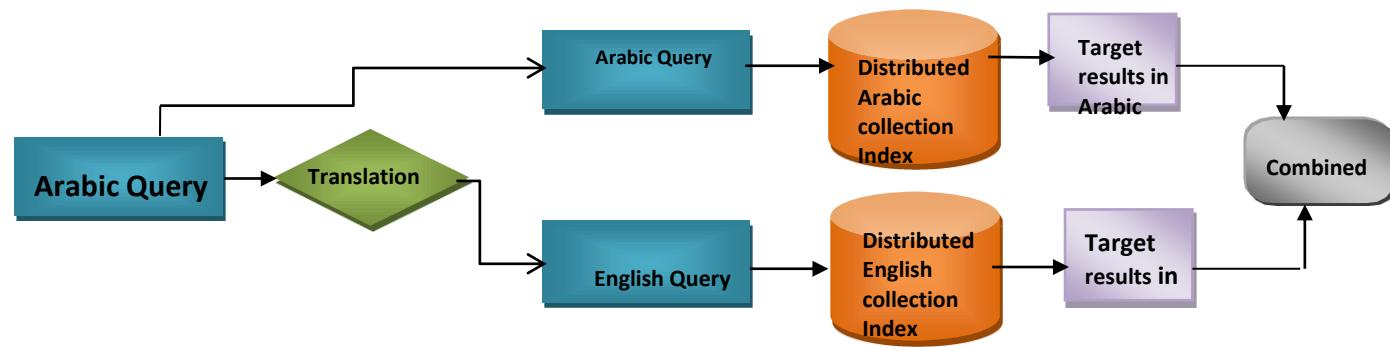
مثال لذلك إذا كان المستخدم يبحث عن (مفهوم ال deadlock) في فهرس مركزي يضم وثائق عربية ووثائق إنجليزية ووثائق ذات لغتين (عربي-إنجليزي)، يتم ترجمة الجزء العربي من الإستعلام إلى اللغة الإنجليزية وترجمة الجزء الإنجليزي إلى اللغة العربية ودمج الأجزاء المترجمة في إستعلام واحد ليصبح (مفهوم الاقفال deadlock concept)، ويتم البحث بهذا الاستعلام في الفهرس المركزي وبالنسبة لمجموعة الوثائق العربية يتم تجاهل وزن الجزء الانجليزي من الإستعلام (أي وزنه يساوي صفر)، وبالنسبة لمجموعة الوثائق الإنجليزية وفي هذه الحالة يتم تجاهل وزن الجزء العربي من الإستعلام (أي وزنه يساوي صفر)، وبالنسبة للوثائق ذات اللغتين لن يتم تجاهل وزن أي من الجزء العربي والجزء الإنجليزي من الإستعلام، وتكون نتائج البحث في هذا الفهرس قائمة مرتبة واحدة.



شكل(5.2) طريقة الفهرسة المركزية

## 2. طريقة الفهرسة الموزعة (*Distributed Indexing Approach*)

في هذه الطريقة تتم فهرسة الوثائق بطريقة موزعة أي لكل لغة فهرس منفرد، ويعتمد هذا المنهج على توزيع الوثائق وفقاً للغاتها، تتم فهرسة الوثائق العربية في فهرس عربي أحادي (monolingual arabic) وفهرسة الوثائق الإنجليزية في فهرس أحادي آخر (monolingual English) والوثائق متعددة اللغات يتم فهرسة محتوياتها على حسب لغاتها. مثل لذلك إذا قام المستخدم بكتابية إستعلام "مفهوم العولمة" يتم ترجمة هذا الإستعلام إلى اللغة الإنجليزية "globalization concept" ومن ثم يتم البحث بالإستعلام العربي ("مفهوم العولمة") في الفهرس العربي وإسترجاع النتيجة، ومن ثم البحث بالإستعلام الإنجليزي (globalization concept) في الفهرس الإنجليزي ودمج تلك النتائج.

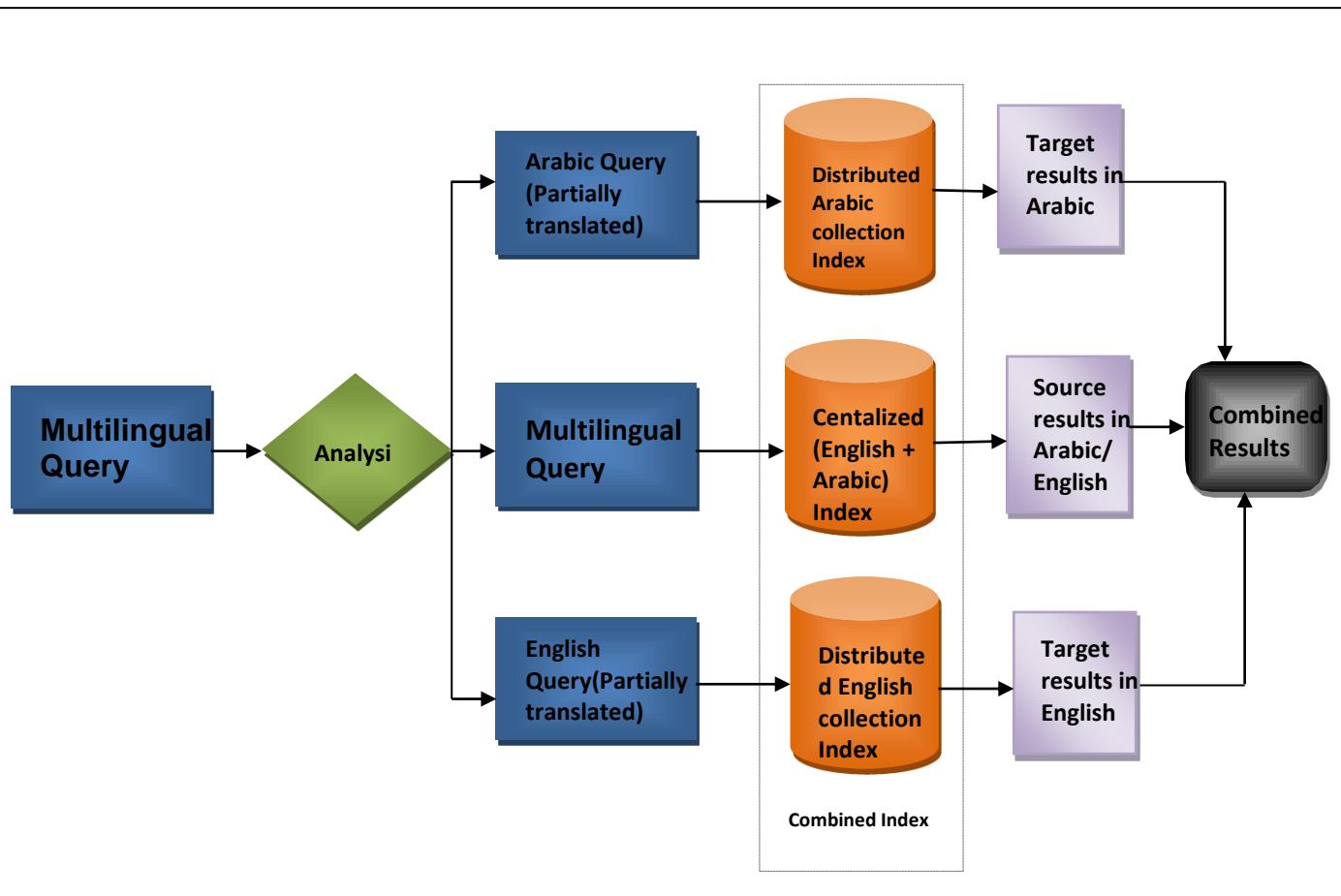


شكل(6.2) يوضح الفهرسة الموزعة

## 3. الطريقة المركزية الموزعة (*Hybrid approach of Indexing*)

تدمج هذه الطريقة بين الفهرسة المركزية والفهرسة الموزعة للإستفاده من مميزات كل منها وفي هذه الطريقة تتم فهرسة الوثائق أحادية اللغة كما تتم فهرستها في الطريقة الموزعة وذلك ببناء فهرس منفرد لكل لغة، وفهرسة الوثائق متعددة اللغات بالطريقة المركزية وذلك ببناء فهرس مركزي لكل من الوثائق متعددة اللغات (عربي-إنجليزي).

مثال لذلك إذا كان المستخدم يبحث عن (مفهوم الـ inheritance) يقوم محرك البحث بتحويل الإستعلام إلى اللغة العربية (monolingual arabic) مرة بحيث يصبح (مفهوم الوراثة) ويتم البحث به في الفهرس العربي، ومرة أخرى إلى الإنجليزية (monolingual English) (ويصبح inheritance concept) بحيث يتم البحث به في الفهرس الإنجليزي. واستعلام ثالث بدمج الإستعلام باللغتين العربية والإنجليزية (مفهوم الوراثة inheritance concept)، ويتم البحث به في الفهرس المركزي متعدد اللغات بحيث يتم تطبيق وزنه والبحث به في ذلك الفهرس ومن ثم دمج النتائج المسترجعة من كلا من الفهارس.



شكل(5.2) طريقة الفهرسة المركزية

## 4.2.2 خوارزميات دمج النتائج (Algorithms)

### 1. خوارزمية القيمة الخام (Raw Score)

في هذه الخوارزمية يتم دمج الوثائق من قائمة النتائج المسترجعة من مجموعة الوثائق العربية مع قائمة النتائج المسترجعة من مجموعة الوثائق الإنجليزية وقائمة الوثائق المسترجعة من مجموعة الوثائق متعددة اللغات، وذلك بأخذ وثيقة من العربية ووثيقة من الإنجليزية وأخرى من متعددة اللغات ومن ثم تقوم بعرضها للمستخدم.

### 2. خوارزمية تطبيق الوزن (Normalized score merging)

تعتمد هذه الخوارزمية على أكبر وزن في قائمة الوثائق المسترجعة وأيضاً على أقل وزن في تلك القائمة بحيث يتم تطبيق وزن(score) كل وثيقة بطرح وزنها من أكبر وزن في القائمة مقسوماً على المدى(أكبر وزن - أقل وزن) ومن ثم دمج النتائج المسترجعة من كل فهرس بناءً على هذا الوزن المطبع، فمن عملية البحث يتم إسترجاع ثلاث قوائم الأولى تحتوي على وثائق عربية والثانية تحتوي على وثائق إنجليزية والثالثة تحتوي على وثائق باللغتين(عربي-إنجليزي) وتقوم بدمج تلك الوثائق بأخذ الوثيقة الأولى من كل قائمة من القوائم المطبعة ومقارنتها أوزانها ووضع الوثيقة ذات الوزن الأعلى في بداية القائمة التي سيتم عرضها للمستخدم.

### 3. خوارزمية CORI (Collection Retrieval Inference Network)

لكي نقوم بتطبيع الوزن قمنا بإستخدام هذه الخوارزمية لأنها تعتبر اكثراً دقة من سابقاتها وفيها يتم تطبيق الوزن بناءً على طول الوثائق المسترجعة من كل مجموعة. بحيث يكون الوزن النهائي المطبع لا ي وثيقه عبارة عن وزن الوثيقة مضروباً في وزن المجموعة التي تنتمي له، ونستخدم ذلك الوزن المطبع لدمج تلك النتائج .  
الفكرة الأساسية لهذه الخوارزمية هي زيادة وزن الوثائق التي لها وزن أعلى من الوزن المتوسط، ويقل وزن الوثائق التي لها وزن أقل من الوزن المتوسط . ومن دواعي استخدامنا لهذه الخوارزمية انها بسيطة لأن مدخلاتها هي وزن الوثائق وطول النتائج المسترجعة . وايضاً هذه الخوارزمية لا تحتاج الى معلومات عن مجموعة الوثائق وبالتالي فإن الوسيط (Broker) الذي يتعامل مع هذه الخوارزمية لا يحتاج الى تخزين معلومات عن مجموعة الوثائق . ولكن اذا كانت التنبؤات عن مجموعة الوثائق مطلوبة في بيئة متغيرة او حركية (مثل الويب) فإن هذه المعلومات تحتاج الى التحديث بصورة دورية وهذا غير ممكن الا بوجود تعاون بين الوسيط وخوادم المجموعات (Collection Server).

في هذا الباب سيتم توضيح الطرق المتتبعة في عملية فهرسة الوثائق والخوارزميات المتتبعة في دمج النتائج المسترجعة من عملية البحث في مجموعة الفهارس والتقييمات التي تم استخدامها لتكوين بيئة التجربة .

## 1.3 الطرق المتبعة في فهرسة الوثائق

### 1.1.3 الطريقة المركزية

في البداية كانت أنظمة إسترجاع المعلومات تقوم بفهرسة الوثائق بطريقة مركزية وذلك بعمل فهرس واحد لجميع الوثائق سواء كانت أحادية أو متعددة اللغات . ويتم حساب الوزن للمصطلح الوارد في الإستعلام والوارد في الوثيقة كالتالي :

$$TF * \log(N/DF)$$

: بحيث أن :

TF : عدد مرات ظهور المصطلح في الوثيقة المعينة

DF : عدد الوثائق التي ورد فيها المصطلح

N: عدد كل الوثائق في المجموعة المعينة

ولكن لهذه الطريقة عدة عيوب وعيوبها الأساسي هو زيادة الوزن (Overweighting) ويعني أن أوزان الوثائق التي تنتهي إلى المجموعة التي تحتوي على وثائق بسيطة سيكون أعلى من أوزان الوثائق التي تنتهي إلى المجموعة التي تحتوي على وثائق أكثر. وذلك لأنه سيزيد عدد الوثائق (لأنه يتم وضعها في مجموعة واحدة أي أن عدد الوثائق N سيزيد ) ولكن عدد تكرار المصطلح او عدد مرات ظهوره (TF) و عدد الوثائق التي ورد فيها المصطلح(DF) ستنظل كما هي من دون تغيير مما يؤدي إلى زيادة أوزان المصطلحات التي تظهر في المجموعة التي تحتوي على وثائق بسيطة فمثلاً إذا كان عدد وثائق اللغة العربية (2,000 وثيقة ) وعدد وثائق اللغة الإنجليزية (18,000 وثيقة ) وتم دمج هذه الوثائق في مجموعة واحدة (مركزية ) سيصل عدد الوثائق الكلية إلى (20,000 وثيقة)، ولذلك ستكون هناك زيادة في وزن الوثائق العربية لأن المصطلحات الواردة في الوثائق العربية تجعلها مميزة أكثر من نظيراتها في اللغة الإنجليزية [1] .

وعند التحدث عن الإستعلامات و الوثائق متعددة اللغات سيكون هناك عيب آخر لهذه الطريقة وهو أن المصطلحات المتشابهة عبر اللغات تحسب أوزانها بطريقة منفصلة على أساس أنها مصطلحات مختلفة، فمثلاً إذا كان لدينا الإستعلام الآتي ("الوراثة Inheritance" حيث أن كل مصطلح هو ترجمه للأخر ) سيكون وزن المصطلح "وراثة" في الفهرس المركزي مختلف عن وزن المصطلح "Inheritance" وسيتم جمع الوزنين للمصطلحين أعلاه للبحث في الفهرس و عند البحث فإن الوثائق العربية سيوجدها المصطلح "وراثة" فقط وفي الوثائق الانجليزية سيوجده المصطلح "Inheritance" فقط ولكن في الوثائق متعددة اللغات (عربي-إنجليزي)سيوجد المصطلحان مما يزيد من وزن الوثائق متعددة اللغات وهذا هو السبب الذي يؤدي إلى هيمنة الوثائق متعددة اللغات على بداية القائمة المسترجعة . أحياناً يكون المصطلح مصحوباً بترجمته في الوثائق خاصة في الوثائق ذات اللغات غير الانجليزية مما ينتج عنه تكرار المصطلح "Co-occurring Terms" فمثلاً مصطلح "الإغفال " في الوثائق ذات اللغة العربية قد يكون مصحوباً بترجمته "deadlock" مما ينتج عنه تكرار المصطلح لكن بلغات مختلفة وهذه الخاصية موجودة على عدد معقول من الوثائق غير الإنجليزية على الشبكة العنكبوتية. وعندما يأتي الحديث عن الإستعلام ات متعددة اللغات في الفهارس المركزية فإن هذه المشكلة "Co-occurring Terms" تؤدي إلى زيادة أوزان الوثائق متعددة اللغات مما يجعلها تكسب أوزاناً إضافية هي ليست جزءاً حقيقياً من وزنها الأصلي فمثلاً إذا كان لدينا وثيقتين هما

و1، و2 . وكانت الوثيقة 1 وثيقة متعددة اللغات حيث فيها اللغة العربية هي اللغة الرئيسية بالإضافة إلى اللغة الإنجليزية كلغة ثانوية وكانت الوثيقة 2 هي وثيقة ذات لغة احادية وهي الانجليزية

"1: " تؤدي عملية التطبيع normalization لإنشاء مجموعة جداول tables ذات..."

" 2: " The process of normalization leads to the creation of tables, whose..."

ولكن الوثيقة 2 هي بالضبط ترجمة للوثيقة 1 فإذا كان الإستعلام كالتالي "التطبيع normalization" هذا الإستعلام يجعل الوثيقة متعددة اللغات 1 ذات صلة أعلى لذلك تكون في أعلى قائمة الوثائق المسترجعة لكن يمكن ان توجد وثائق اخرى ذات صلة اكبر بالإستعلام من هذه الوثيقة لكن تكرار المصطلحات في الوثيقة متعددة اللغات أدى إلى زيادة وزنها مما جعلها تعتلي قائمة الوثائق المسترجعة وهذه من المشاكل الرئيسية في الوثائق والإستعلامات متعددة اللغات في الطريقة المركزية [6] .

ولكن من مميزات الفهرسة المركزية أنها تقوم بعملية البحث والإسترجاع من فهرس واحد بغض النظر عن لغة الوثيقة وبالتالي لا تتطلب هذه الطريقة عملية دمج لتلك الوثائق وهذه الميزة مفيدة جداً للوثائق ذات اللغات المتعددة وذلك لأن عملية الإسترجاع من الفهرس الواحد ذات كفاءة أعلى من عملية الإسترجاع من الفهارات المتعددة [1] .

### 2.1.3 الطريقة الموزعة

ولحل مشكلة زيادة الوزن في طريقة الفهرسة المركزية تمت فهرسة الوثائق بطريقة موزعة أي لكل لغة فهرس منفرد ، ويعتمد هذا المنهج على توزيع الوثائق وفقاً للغاتها [2].

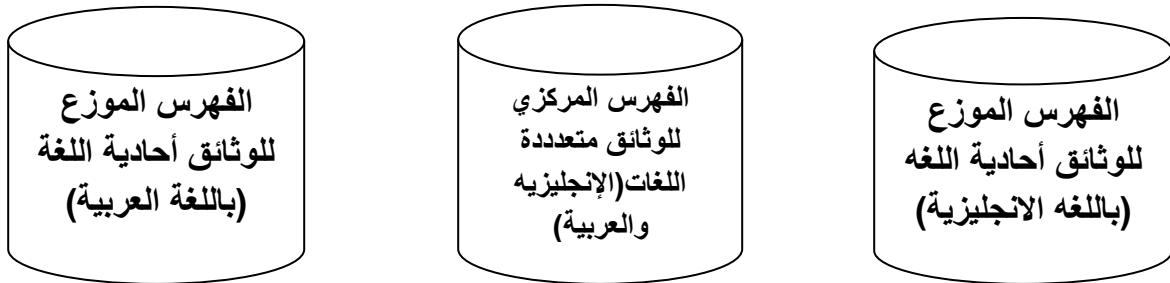
ولكن عند التحدث عن هذه الطريقة لن تكون مشكلة الوثائق المتعددة فقط محصورة على الوزن وإنما تمتد إلى كيف ستتم فهرسة الوثائق متعددة اللغات ؟ فكما ذكر سابقاً أن هذه الطريقة تعتمد على تقسيم الوثائق وفهرسة محتوياتها على حسب اللغة [8]. ولكن عند إجراء عملية التقسيم هذه على الوثائق متعددة اللغات ( على حسب اللغات الواردة فيها ) هذا يجعل المعلومات والمصطلحات الموجودة في الوثيقة متعددة اللغات تفقد قيمتها ومعناها ،بالإضافة إلى ذلك إذا تم تقسيم الوثائق متعددة اللغات وفهرستها على عدد من المجموعات الفرعية (sub-collections) إعتماداً على اللغات المكونة منها تلك الوثيقة ، فإن ذلك سيؤدي إلى نقصان وزن الوثيقة متعددة اللغات "underweighting" في كل مجموعة فرعية لأن الوزن يتم حسابه فقط من جزء الوثيقة الموجود في المجموعة الفرعية وليس من الوثيقة كاملة وهذا يجعل الوثائق متعددة اللغات غير قابلة للمقارنة مع الوثائق احادية اللغات(لأن الوثيقة احادية اللغة لن يتم تقسيمها إلى وثائق فرعية وبالتالي لا تفقد معناها ولا ينقص وزنها) مما يجعل الوثائق متعددة اللغات حتى وإن كانت ذات صلة أعلى بالإستعلام ان تكون في نهاية القائمة المسترجعة ، وهذا هو الفرق بين الطريقة الموزعة والطريقة المركزية حيث أنه في الطريقة المركزية تتم فهرسة الوثيقة متعددة اللغات بإفتراض أنها لغة أحادية مما يؤدي إلى زيادة اوزان الوثائق متعددة اللغات وفي الطريقة الموزعة يتم فهرسة الوثيقة متعددة اللغات بتقسيمها إلى وثائق فرعية حسب اللغات المكونة لها مما يؤدي إلى نقصان اوزانها [2].

هناك فائدته كبيرة في الفهارات الموزعة أحادية اللغة وهي أن تكون عملية الإسترجاع منها ذات كفاءة عالية وذلك نتيجة للمطابقه بين الوثيقة والإستعلام الذي تم ترجمته، خاصة مع الجزء الواحد المترجم للإستعلام ذو اللغات المتعددة.

### 3.1.3 الطريقة المركزية الموزعة

بعد التعرف على الطريقتين أعلاه و معرفة عيوب كل منهما كان لابد من استخدام طريقة للحد من سلبيات كل منها و تمكن المستخدم من البحث باستخدام الإستعلام ات متعددة اللغات في الوثائق متعددة اللغات وان تكون نتيجه البحث تضم الوثائق ذات الصلة الأعلى بالإستعلام بغض النظر عن اللغة المستخدمة في كتابته او امكانية المستخدم في التعبير عن احتياجاته . لذلك تم استخدام هذه الطريقة التي تدمج بين الفهرسة المركزية والفهرسة الموزعة لبناء الفهرس في هذا البحث للحد من سلبيات كلا من الطريقتين السابقتين والإستفاده من فوائدهما.

و هذا النهج الهجين من كلتا الآليتين يساعد في التعامل مع الوثائق أحادية اللغة و متعددة اللغات بصورة فعالة ووفقاً لذلك من المحتمل أن يكون ذلك النهج الهجين أكثر ملائمه من بقية الطرق وذلك ببناء فهرس أحادي اللغة للوثائق أحادية اللغة (monolingual ) وهذا يكون مثل الطريقة الموزعة وبناء فهرس اخر للوثائق ذات اللغات المتعددة (الطريقة المركزية) كما هو مبين في الشكل التالي(باعتبار اللغتين الإنجليزية و العربية):



شكل (1.3) يوضح استخدام الفهرسة التي تدمج بين الموزعة والمركزية

وهذه الطريقة المقترحة تقلل من مشكلة زيادة الوزن(Overweighting) الموجودة في الفهرسة المركزية إلى أدنى مستوى وذلك لأن الوثائق أحادية اللغة لن تتم فهرستها في الفهرس المركزي(أي لن يتم دمجها في فهرس واحد مع الوثائق متعددة اللغات ) وبالتالي لن يزيد عدد الوثائق الكلي الموجود في المجموعة وهذا يحد من مشكلة زيادة الوزن . إلى جانب حل مشكلة زيادة الوزن في الطريقة المركزية من ناحية أخرى هذه الطريقة تضم فوائد الفهرسة المركزية التي تقوم بعملية البحث والإسترجاع من فهرس واحد بغض النظر عن لغة الوثيقة كما ذكر سابقاً وبالتالي يمكن بناء فهرس مركزي وحيد يحوي الوثائق ذات اللغات المتعددة أما الوثائق ذات اللغة الواحدة يتم فهرستها بالطريقة الموزعة.

أيضاً هذه الطريقة المقترحة تحد من مشكلة نقصان الوزن(underweighting) بالنسبة للوثائق متعددة اللغات في الطريقة الموزعة لانه لا يتم تقسيم الوثيقة متعددة اللغات إلى مجموعة من الوثائق الفرعية أي لا يتم فصل محتواها على حسب اللغة لذلك لن تضيع المعلومات الموجودة في تلك الوثائق ولن تفقد قيمتها وعلى عكس الطريقة الموزعة هذا النوع من الفهرسة يجعل أوزان الوثائق ذات اللغات المتعددة مكافئة لغيرها من الوثائق الأحادية(لأنه لا يتم تقسيم الوثيقة متعددة اللغات).

فمثلاً إذا كان لدينا ثلاثة فهارس أحدهم للغة العربية يحتوي على مجموعة الوثائق ذات اللغة العربية فقط وفهرس اللغة الانجليزية يحتوي على مجموعة الوثائق ذات اللغة الانجليزية فقط و الفهرس ذو اللغتين للوثائق التي تحتوي على اللغتين معاً فإذا كان المستخدم يبحث عن (مفهوم deadlock) فإن عملية البحث ستتم كالتالي: يقوم محرك البحث

بتحويل الإستعلام الى اللغة العربية (monolingual arabic) مرة وذلك بترجمة الأجزاء الإنجليزية في الإستعلام بحيث يصبح (مفهوم الإقفال) ويتم به البحث في فهرس الوثائق العربية (ويتم تجاهل وزن الجزء الانجليزي من الإستعلام أي أن وزنه يساوي صفر) ويتم تحويل الإستعلام الى اللغة الانجليزية (monolingual English) وذلك بترجمة الأجزاء العربية في الإستعلام بحيث يصبح (deadlock concept) ويتم البحث به في فهرس الوثائق الإنجلizية (ويتم تجاهل وزن الجزء العربي من الإستعلام أي أن وزنه يساوي صفر) ، ويتم دمج الإستعلامين أعلاه(العربي،الإنجليزي) ليصبح الإستعلام (مفهوم الإقفال deadlock concept)، ويستخدم هذا الإستعلام للبحث في مجموعة الوثائق ذات اللغتين بحيث تنتج من عملية البحث في الثلاثة فهارس ثلاث قوائم تحتوي على الوثائق التي تعتبر ذات صلة أعلى بالإستعلام بحيث تكون هذه القوائم مرتبة ترتيباً تنازلياً اعتماداً على أوزان الوثائق في كل قائمة وهنا تكمن المشكلة في كيفيه دمج النتائج المسترجعة من الفهارس الثلاثة في قائمة واحدة تظهر للمستخدم بحيث تكون فيها الوثائق الأكثر صلة بالإستعلام في أعلى القائمه المسترجعة .

## 2.3 خوارزميات دمج النتائج (Algorithms)

### 1.2.3 خوارزمية القيمة الخام (Raw Score)

في هذه الخوارزمية لا يتم تطبيق وزن الوثائق او تعديله وإنما يتم دمج القوائم المسترجعة من الفهارس الثلاثة (عربي ،إنجليزي،متعدد اللغات ) في قائمة واحدة كما ذكر في الباب السابق إعتماداً على أوزانها الأصلية دون تعديل[3] فمثلاً إذا أجرينا عملية البحث باستخدام الإستعلام الآتي "مفهوم الـbinary tree" وكانت الوثائق المسترجعة من الفهرس ذو اللغة العربية هي :

"في علم الحاسوب شجرة البحث الثنائية أو أحياناً يطبق عليها شجرة البحث المرتبة ...": Da1

"هي بنية معلوماتية حيث ان كل رأس فيها اثنين من الابناء على الاكثر غالباً مميزين بـ اب ايمن ...": Da2

" هي الشجرة التي يكون لكل عقدة فيها ابن وحيد فقط . يكون ارتفاع هذه الشجرة عادةً...": Da3

بالأوزان الآتية على التوالي: 0.318662405 ، 0.590278973 ، 0.67558967

والوثائق المسترجعة من الفهرس ذو اللغة الانجليزية هي:

the basic concepts of binary trees, and "De1:  
... "then

De2: " The binary tree is a fundamental data  
."...used in computer science structure

In **computer science**, a binary tree is a "De3:  
..."**tree data structure** in which each node has

بالأوزان الآتية على التوالي : 0.7761199  
، 0.63523394 ، 0.7066437 ،

والوثائق المسترجعة من الفهرس متعدد اللغات هي :

Dm1: "مفهوم الشجرة الثنائية هي عبارة عن هيكل بيانات...": binary tree

Dm2: "الشجرة الثنائية تتكون من مجموعة من العقد ... nodes"

Dm3: "مفهوم الشجرة الثنائية (binary tree concept) : هي جزء اساسي من هيكل البيانات..."

بالأوزان الآتية على التوالي :

0.0364784 ، 0.6044762 ، 0.7046437

فإنه في هذه الخوارزمية سيتم دمج الوثائق المسترجعة في القوائم الثلاث بمقارنة وزن الوثيقة الأولى من كل قائمة ووضع الوثيقة ذات الوزن الأعلى في بداية القائمة التي سيتم عرضها للمستخدم ،أي أنه سيتم مقارنة اوزان الوثائق

الثلاث، وزن الوثيقة Da1 في القائمه المسترجعة من الفهرس ذو اللغة العربية مع وزن الوثيقة De1 في القائمه المسترجعة من الفهرس ذو اللغة الإنجليزية ومع وزن الوثيقة Dm1 في القائمه المسترجعة من الفهرس ذو اللغتين (عربي-إنجليزي) ويتم وضع الوثيقة التي تحتوي على وزن أكبر (وهي De1) في أعلى القائمه التي سيتم عرضها للمستخدم والإنتقال إلى الوثيقة التي تليها في القائمه ومقارنتها مع الوثيقتين من الفهرسين الآخرين إلى ان نتوصل إلى نهاية القوائم الثلاث . وستكون القائمه النهائية عند دمج القوائم الثلاث هي

*De1,De2,Dm1,Da2,De3 ,Dm2 ,Da2, Da3,Dm3*

ولكن من عيوب هذه الخوارزميه انه يمكن أن تكون أوزان الوثائق في القوائم المختلفة غير قابلة للمقارنة مع بعضها البعض فمثلاً يمكن أن تكون أوزان الوثائق في القائمه المسترجعة من الفهرس العربي كبيرة جداً وأوزان الوثائق في القائمه المسترجعة من الفهرس الإنجليزي صغيره جداً وبالتالي ستظهر الوثائق العربية أولاً مع العلم بأنه يمكن أن تكون الوثائق الإنجليزية ذات صلة أكبر بالإستعلام من الوثائق العربية.

### 2.2.3 خوارزمية تطبيق الوزن (Min Max)

بما أنه يمكن أن تكون أوزان الوثائق في القوائم المختلفة غير قابلة للمقارنة عندها يكون من المعقول أن يتم تطبيق الوزن أو تعديله لكل قائمه من القوائم المسترجعة من الفهارس الثلاثة كلاً على حده قبل أن يتم دمجها في قائمه واحدة ومن إحدى الطرق المستخدمة هي قسمة أوزان الوثائق في كل قائمه على أكبر وزن (max score ) وهو وزن أول وثيقة في تلك القائمه المسترجعة المرتبة مما يجعل أوزان تلك الوثائق محصورة بين 0 و 1 الوثيقة ذات أعلى وزن سيطبع وزنها للواحد وآخر وثيقة في تلك القائمه سيطبع وزنها ويكون قريباً من الصفر. ومن ثم يتم إعادة ترتيب القوائم بعد تعديلها ومن ثم دمجها في قائمه واحدة تظهر للمستخدم [4].

وأيضاً في هذه الخوارزميه يتم استخدام المعادله التالية لتطبيق الوزن:

$$\text{normalized\_score} = \frac{\text{original\_score} - \text{min\_score}}{\text{max\_score} - \text{min\_score}}$$

حيث أن:

Original\_score: هي الوزن الأصلي للوثيقة

min\_score: هي أصغر وزن في القائمه المسترجعة (وزن آخر وثيقة في القائمه)

Max\_score: هو أكبر وزن في القائمه المسترجعة (وزن أول وثيقة في القائمه)

فباعتبار القائمه المسترجعة من الفهرس العربي min\_score و Max\_score هي 0.67558967 ، 0.318662405 وعند تطبيق المعادله أعلاه سيكون الوزن الجديد المطبع لهذه القائمه لكل من الوثائق الثلاث Da1, Da2, Da3 هو 0.761، 0.761 على التوالي

و عند تطبيق المعادلة على القوائم الثلاث وتطبيع أوزان الوثائق فيها ودمجها في قائمه واحدة كما ذكر في الباب السابق فإن القائمه النهائية التي سيتم إظهارها للمستخدم هي :

*De1,Da1,Dm1,Dm2,Da2 ,Da2 ,De3,Da3,Dm3*

### 3.2.3 خوارزمية CORI

كما ذكر سابقاً هذه طريقة أخرى لتطبيع الوزن وذلك باستخدام أوزان الوثائق بالإضافة إلى الوزن القائمة المطبع (weighted score) الذي يتم استخراجه من القائمة التي تنتهي لها الوثيقة [5].

وهذه الطريقة تعتبر أكثر دقة من سابقاتها ومن أفضل خوارزميات دمج النتائج وفيها يتم تطبيق الوزن بناء على طول الوثائق المسترجعة من كل مجموعة بحيث يكون الوزن النهائي المطبع لأي وثيقة عبارة عن وزن الوثيقة مضروباً في وزن المجموعة التي تنتهي لها، ونستخدم ذلك الوزن المطبع لدمج تلك النتائج [7].

والفكرة الأساسية لهذه الخوارزمية هي زيادة وزن الوثائق التي لها وزن أعلى من الوزن المتوسط ، ويقل وزن الوثائق التي لها وزن أقل من الوزن المتوسط . ومن دواعي إستخدامنا لهذه الخوارزمية أنها بسيطة لأن مدخلاتها هي وزن الوثائق وطول النتائج المسترجعة . وأيضاً هذه الخوارزمية لا تحتاج إلى معلومات عن مجموعة الوثائق وبالتالي فإن الوسيط (Broker) الذي يتعامل مع هذه الخوارزمية لا يحتاج إلى تخزين معلومات عن مجموعة الوثائق . ولكن إذا كانت التنبؤات عن مجموعة الوثائق مطلوبة في بيئه متغيرة او حركية (مثل الويب) فإن هذه المعلومات تحتاج إلى التحديث بصورة دورية وهذا غير ممكن إلا بوجود تعاون بين الوسيط وخوادم المجموعات (Collection Server).

وهذه الخوارزمية تقوم بحساب وزن الوثائق بناءً على طول النتائج المسترجعة (عدد الوثائق المسترجعة) لكل مجموعة ، والوزن لكل مجموعة يقوم على إفتراض أن كل الوثائق المسترجعة هي أعلى صلة بالاستعلام وتسخدم هذه الطريقة لتطبيع الوزن.

$$\text{normalized}_{\text{score}} = \text{original}_{\text{score}} * \left[ 1 + \frac{S_i - \text{avg}_S}{\text{avg}_S} \right]$$

حيث:

$\text{avg}_S$  : متوسط وزن الوثائق في المجموعة الواحدة

$S_i$  : هو وزن مجموعة الوثائق ذات اللغة المراد تطبيق وثائقها.

وبعد أن تقوم خوارزمية ال CORI بتعديل أوزان الوثائق عند تطبيق معادلة تطبيع الوزن أعلاه يتم دمج القوائم الثلاث في قائمة واحدة بمقارنة وزن الوثيقة الاولى من كل قائمة ووضع الوثيقة ذات الوزن الأعلى في بدايه القائمه التي سيتم عرضها للمستخدم ويكون ترتيب القائمه الناتجه كالتالي :

*Dm1,Dm2,Da1,Da2,De1,De2,De3, Da3,dm3*

### 3.3 التقنيات المستخدمة

#### 1.3.3 لغة الجافا (Java)

هي عبارة عن لغة برمجة ابتكرها جيمس جوسلينج (James Gosling) في عام 1992 أثناء عمله في مختبرات شركة صن (Sun Microsystems) وذلك لاستخدامها بمثابة العقل المفكرة المستخدم لتشغيل الأجهزة التطبيقية الذكية مثل التلفاز التفاعلي.

مميزاتها :

- ❖ إحدى لغات الكائنية التوجه (Object Oriented)، آمنة وجيدة الأداء.
- ❖ إضافة الحركة والصوت إلى صفحات الويب، كتابة الألعاب والبرامج المساعدة، وإنشاء برامج ذات واجهة مستخدم رسومية.
- ❖ تصميم برمجيات تستفيد من كل مميزات الأنترنت، توفر لغة الجافا بيئة تفاعلية عبر الشبكة العنكبوتية وبالتالي تستعمل لكتابه برامج تعليمية للإنترنت عبر برمجيات المحاكاة الحاسوبية للتجارب العلمية وبرمجيات الفصول الإفتراضية للتعليم الإلكتروني والتعليم عن بعد. لا تحصر فاعلية الجافا في الشبكة العنكبوتية فقط بل تمكنا من إنشاء برامج للاستعمال الشخصي والمهني حيث يوجد العديد من البرمجيات التي تسهل عملية كتابة الأوامر ك [9] Eclipse و netbeans.

في المشروع محل الدراسة قمنا بإستخدام هذه اللغة لأنها تدعم الكثير من المكتبات (API) التي تساعدنا في تطبيق النظام بصورة أكثر كفاءة.

#### 2.3.3 تقنية (Java Server Page) JSP

هي تقنية صممت لتساعد مطوري البرمجيات في إنشاء صفحات ويب متغيرة المحتوى (динاميكية). تم إصدارها عام 1999 بواسطة شركة (Sun Microsystems) وهي تشبه لغة PHP إلا أنها تستخدم لغة البرمجة جافا لتنفيذ ونشر محتوى الصفحات المكتوب بهذه اللغة لابد أن يكون خادم الويب (web server) متوافق مع لغة سيرفات Jetty أو Apache Tomcat (servlet).

مميزاتها :

- ❖ لا تعتمد على منصة بعينها (Multi Platform).
- ❖ يوجد العديد من المكونات التي يمكن إعادة استخدامها داخل صفحة JSP ولعل أشهرها (JavaBeans).
- ❖ نسبة لإعتمادها على لغة الجافا فهي تمكن من الاستفادة من قوة ومزایا الجافا [10].

في المشروع محل الدراسة تم إستخدام هذه التقنية لأن النظام وكما ذكرنا سابقاً تم تطويره بلغة الجافا ولأن لغة JSP تُمكِّن من التعامل مع برامج مكتوبة بلغة الجافا بكل سهولة ، قمنا بإستخدامها لبناء واجهة ويب للنظام.

### 3.3.3 خادم ويب تومكات (Apache Tomcat Server)

هو خادم ويب وحاوية سيرفلت (Servlet) مفتوح المصدر أصدرته مؤسسة أبashi للبرمجيات (Apache Software Foundation)، كما يوفر أيضاً هذا الخادم العديد من المميزات التي تهيئ منصة عمل لتطوير تطبيقات و خدمات الويب.

مميزاته : المرونة (flexibility)، الإستقرار (Stability). [11]

في المشروع محل الدراسة قمنا بإستخدام هذا الخادم نسبة لأن لغة الويب المستخدمة في صفحاتنا هي JSP وهو يتبع لنا إمكانية ترجمة هذه اللغة بصورة مجانية.

### 4.3.3 محل جيريوكو لوثائق HTML (Parser)

هو مكتبة جافا تتيح التحليل والتلاعب في أجزاء من وثيقة HTML، و أيضاً الأجزاء التي بجانب الخادم (server-side tags)، كما أنه يمكن عمل استنساخ حرفياً لأي جزء في وثيقة HTML غير معترف به أو غير صالح.

مميزاته :

- ❖ التعامل مع الوثائق الكبيرة في شكل أجزاء متسلسلة مما يحافظ على مساحة الذاكرة.
- ❖ متطلباته من الموارد ذات نسبة معقولة مقارنة مع الم حللات الأخرى.
- ❖ يسمح بإخلاص كل النص الموجود بالوثيقة المحددة.
- ❖ يسمح بحذف المسافات غير المرغوب فيها بين الكلمات.
- ❖ يسمح بإخلاص النص الموجود في الوثيقة المحددة بصورة مناسبة لتغذية محركات البحث مثل Apache Lucene . [12]

في المشروع محل الدراسة تم إستخدام هذا المحل في إخلاص النص العربي(النص فقط) من صفحات المخزنة كوثائق لفهرستها و البحث فيها.

### 5.3.3 محرك البحث لوسين (Lucene)

هو مكتبة مكتوبة بلغة جافا بواسطة مؤسسة أباشي للبرمجيات (Apache software foundation) تُستخدم للبحث عن المعلومات داخل الوثائق صممت إضافةً إمكانات البحث للتطبيقات الخاصة بالمستخدم.

مميزاته :

- ❖ الأداء العالي و إمكانية التوسيع.
- ❖ مجاني مفتوح المصدر.
- ❖ بالإضافة لإمكانيات البحث توفر هذه المكتبة أيضاً بعض الإضافات أو الطائف الجديدة مثل: تظليل النص و التدقيق الإملائي.

طريقة عمله:

هناك بعض العمليات تتم على الوثائق المراد البحث فيها و عمليات أيضاً تتم على طلب المستخدم.

العمليات التي تتم على مجموعة الوثائق هي:

أ- استخلاص المحتوى (Acquire content)

يتم استخلاص محتوى الوثائق المراد فهرستها للبحث فيها.

ب- بناء وثيقة جديدة (Build document)

يتم بناء وثيقة جديدة مبنية في صورة مفتاح و قيمة لتكون مهيكلة بصورة تسمح بالبحث فيها.

ت- تحليل الوثيقة (Analyze document)

يتم تطبيق و تشذيب الكلمات الموجودة بالوثيقة لرفع درجة كفاءة البحث.

ث- فهرسة الوثيقة (Index document)

تتم إضافة الوثيق المحللة للفهرس.

العمليات التي تتم على طلب المستخدم هي:

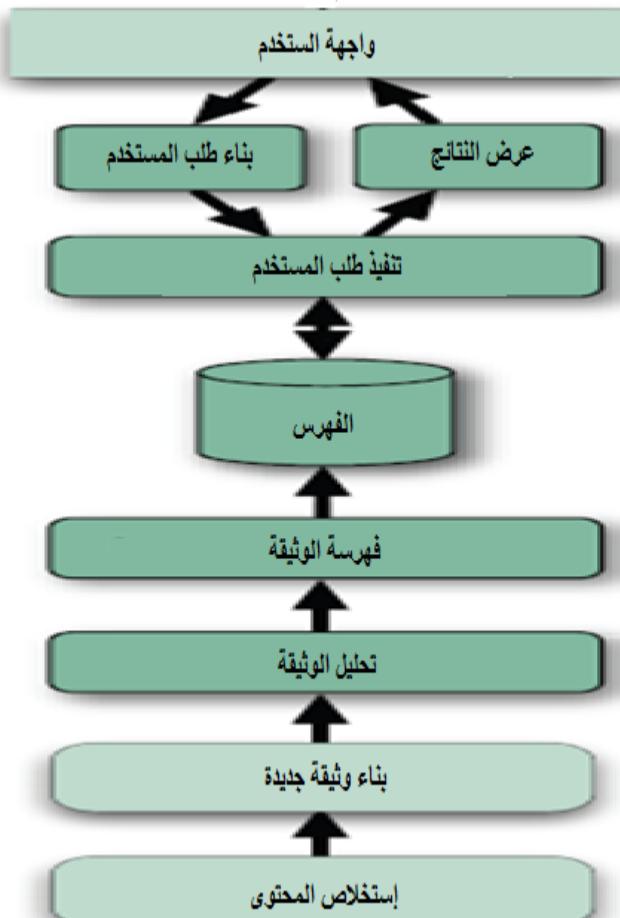
نـ بناء طلب المستخدم (Build query)

يتم تكوين طلب المستخدم باستخدام الحزمة QueryParser التي تقوم بتحويل النص الذي أدخله المستخدم إلى طلب ذو صيغة محددة متوافقة مع صيغة البحث.

## ii. تنفيذ طلب المستخدم (Run query)

يتم البحث في الفهرس عن طلب المستخدم لإسترجاع الوثائق الأكثر ملاءمة.

الشكل (1.3) يوضح مجمل الخطوات والعمليات التي يقوم بها لوسين[13].



شكل (2.3) العمليات التي تتم في محرك بحث لوسين

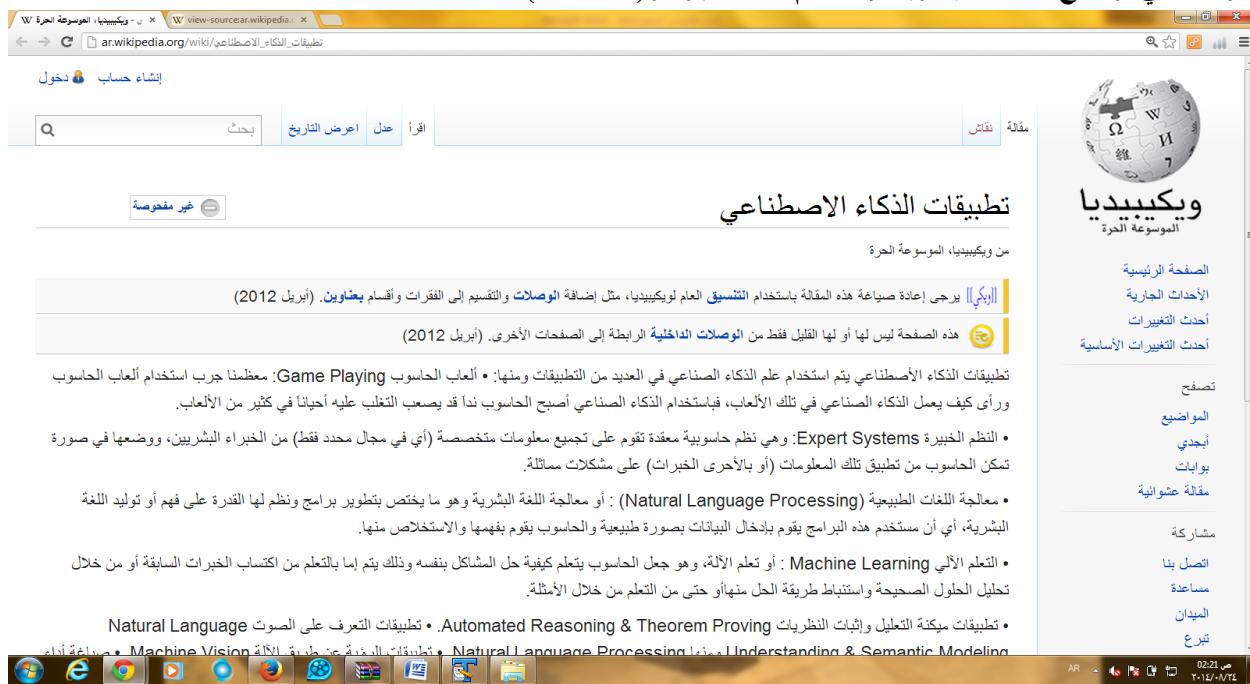
في النظام محل الدراسة تم استخدام محرك بحث لوسين في إضافة إمكانية البحث في الوثائق التي تمت فهرستها بالطريقة قيد الدراسة.

## 1.4 خطوات تطبيق المشروع

تم تكوين بيئة الإختبار لهذا المشروع باستخدام محرك البحث لوسين (Lucene) (النسخة 4.6.0) والذي يستخدم للبحث عن المعلومات داخل الوثائق، وإستخلاص النصوص من الصفحات تم إستخدام محل الجريko (Jericho) وحيث تم في هذه التجربة إستخدام عدد 20 الف وثيقة مختلفة المواضيع مثل وثائق تتحدث عن الذكاء الإصطناعي وأخرى عن حالة الجمود (deadlock) وغيرها من المواضيع. ثم عملية تطبيع النصوص العربية والإنجليزية ومن ثم تشذيبها وتكون الوثائق الجديدة وفهرستها ببناء فهرس بالطريقة المقترنة وهي التي تدمج بين الطريقة المركزية والموزعة كما تم توضيحها سابقاً، وترجمة الأجزاء العربية والإنجليزية والبحث بها في ذلك الفهرس، وأخيراً يتم دمج النتائج المسترجعة من عملية البحث بخوارزميات الدمج التي تم توضيحها سابقاً، وفيما يلي شرح لهذه الخطوات:

### أولاً: تحديد مجلد أو أكثر لفهرسة الملفات الموجودة بداخله

تتم قراءة الملفات الموجودة بداخل المجلد ملف، ملف، واستخلاص النص فقط من الملفات ذات الإمتداد (html) أو (htm)، بإستخدام محل جريko بحيث يتم تجاهل الصور والنصوص غير العربية الموجودة بداخل ملفات (HTML). يتم وضع النص المستخلص في شكل سلسلة من الحروف (string) وإعتبارها مدخل لعملية التطبيع. وفيما يلي توضيح لصفحة قبل وبعد إستخدام محل الجريko (Jericho)



الشكل (1.4): الصفحة قبل إستخدام محل الجريko

ولتنقية مثل هذه الصفحة قمنا بإستخلاص النص فقط (العربي والإنجليزي) وإزالة كل الصور وعلامات الترقيم والأحرف بأي لغة خلاف اللغتين الإنجليزية والعربية وبما أن الصفحة كما موضح في النموذج أعلاه خليط بين اللغتين العربية والإنجليزية فعندما تكون الفقرة باللغة العربية نقوم بتغيير ال <tag> لهذه الفقرة إلى العربية ليصبح <tag lang = "ar"> وفيما يلي توضيح لذلك :

تطبيقات الذكاء الاصطناعي من ويكيبيديا الموسوعة الحرة غير مفحوصة اذهب إلى تصفح ابحث يرجى إعادة صياغة هذه المقالة باستخدام التنسيق العام لويكيبيديا مثل إضافة الوصلات وأقسام بعناوين أبولو هذه الصفحة ليس لها أو لها القليل فقط من الوصلات الداخلية الرابطة إلى المفهومات الأخرى أبولو تطبيقات الذكاء الصناعي في العديد من التطبيقات ومنها ألعاب الحاسوب

معظمنا جرب استخدام ألعاب الحاسوب ورأى كيف يعمل الذكاء الصناعي في تلك الألعاب في باستخدام الذكاء الصناعي `<p lang='ar'>game playing</p>`

المنابع أصبحت قد يصعب التغلب عليه أحياناً في كثير من الألعاب النجم الخيرية وهي نظم حاسوبية مقدمة تقوم على تجميع معلومات متخصصة أي في مجال محدد فقط من الخبراء البشريين

ووضعها في صورة تمكن الحاسوب من تطبيق تلك المعلومات أو بالأحرى الخبراء على مشكلات مماثلة معالجة اللغات أو معالجة اللغة الطبيعية وهو ما يختص بتطوير برنامج `<p lang='ar'>natural language processing</p>`

ونظم لها القدرة على فهم أو توليد اللغة البشرية أي أن مستخدم هذه البرامج يقوم بادخال البيانات بصورة طبيعية أو تعلم الآلة وهو جعل `<p lang='ar'>machine learning</p>` والجهاز يقوم بهما والاستخلاص منها التعلم الآلي الحاسوب يتعلم كيفية حل المشاكل بنفسه وذلك يتم إما بالتعلم من اكتساب الخبرات السابقة أو من خلال تحليل الحلول `<p lang='ar'>الصحيحة واستنباط طريقة الحل منها أو حتى من التعلم من خلال الأسئلة تطبيقات ميكنة التعلم وإثبات النظريات automated reasoning theorem proving</p>`

تطبيقات التعرف على الصوت `<p lang='ar'>natural language understanding semantic modeling</p>` ومنها `<p lang='ar'>natural language processing</p>` `<p lang='ar'>modelling human performance</p>` `<p lang='ar'>machine vision</p>` `<p lang='ar'>الروبوتية عن طريق الآلة</p>` `<p lang='ar'>planning robotics</p>` `<p lang='ar'>languages environments for ai</p>` `<p lang='ar'>الحوسبة الظاهرية والمعالجة المتوازية parallel distributed processing pdp</p>` `<p lang='ar'>emergent computation</p>` `<p lang='ar'>heuristic classification</p>` `<p lang='ar'>ai philosophy</p>` `<p lang='ar'>الفلسفة والذكاء الاصطناعي</p>`

فعلاً عند استخدام هذا العمل لتطوير الأنظمة الحديثة يتم تخزين المabytes داخل الحاسوب لتكون قاعدة بيانات له مثل ما تخزن المعلومات داخل العقل البشري من خلال التعلم والخبرات اليومية التي يكتسبها ثم يتم بعد ذلك تطوير برنامج خاص ليستطيع الحاسوب استخدامها في التعامل مع هذه البيانات واستخدامها بطريقة منطقية في حل المشكلات الازلية لمعنى القرار وقد نجح العلماء حتى الآن في تطوير بعض النماذج المغيرة من نظم الذكاء الاصطناعي ومنها أجهزة الروبوت

## الشكل (2.4) : عملية تعريف tag

تطبيقات الذكاء الاصطناعي من ويكيبيديا الموسوعة الحرة غير مفحوصة اذهب إلى تصفح ابحث يرجى إعادة صياغة هذه المقالة باستخدام التنسيق العام لويكيبيديا مثل إضافة الوصلات وأقسام بعناوين أبولو هذه الصفحة ليس لها أو لها القليل فقط من الوصلات الداخلية الرابطة إلى المفهومات الأخرى أبولو تطبيقات الذكاء الصناعي يتم استخدام علم الذكاء الصناعي في العديد من التطبيقات ومنها ألعاب الحاسوب

معظمنا جرب استخدام ألعاب الحاسوب ورأى كيف يعمل الذكاء الصناعي في تلك الألعاب في باستخدام الذكاء الصناعي `<p lang='ar'>game playing</p>`

المنابع أصبحت قد يصعب التغلب عليه أحياناً في كثير من الألعاب النجم الخيرية وهي نظم حاسوبية مقدمة تقوم على تجميع معلومات متخصصة أي في مجال محدد فقط من الخبراء البشريين ووضعها في صورة تتمكن الحاسوب من تطبيق تلك المعلومات أو بالأحرى الخبراء على مشكلات مماثلة معالجة اللغات الطبيعية أو معالجة اللغة البشرية وهو ما يختص بتطوير برنامج ونظم لها القدرة على فهم أو توليد اللغة البشرية أي أن مستخدم هذه البرامج يقوم بادخال البيانات بصورة طبيعية والجهاز يقوم بهما والاستخلاص منها التعلم الآلي

أو تعلم الآلة وهو جعل `<p lang='ar'>machine learning</p>` والجهاز يقوم بهما والاستخلاص منها التعلم الآلي الحاسوب يتعلم كيفية حل المشاكل بنفسه وذلك يتم إما بالتعلم من اكتساب الخبرات السابقة أو من خلال تحليل الحلول `<p lang='ar'>الصحيحة واستنباط طريقة الحل منها أو حتى من التعلم من خلال الأسئلة تطبيقات ميكنة التعلم وإثبات النظريات automated reasoning theorem proving</p>`

تطبيقات التعرف على الصوت `<p lang='ar'>natural language understanding semantic modeling</p>` ومنها `<p lang='ar'>natural language processing</p>` `<p lang='ar'>modelling human performance</p>` `<p lang='ar'>machine vision</p>` `<p lang='ar'>الروبوتية عن طريق الآلة</p>` `<p lang='ar'>planning robotics</p>` `<p lang='ar'>languages environments for ai</p>` `<p lang='ar'>الحوسبة الظاهرية والمعالجة المتوازية parallel distributed processing pdp</p>` `<p lang='ar'>emergent computation</p>` `<p lang='ar'>heuristic classification</p>` `<p lang='ar'>ai philosophy</p>` `<p lang='ar'>الفلسفة والذكاء الاصطناعي</p>`

أو تعلم الآلة وهو جعل الحاسوب يتعلم كيفية حل المشاكل بنفسه وذلك يتم إما بالتعلم من اكتساب الخبرات السابقة أو من خلال تحليل الحلول الصحيحة واستنباط طريقة الحل منها حتى من التعلم من خلال الأسئلة تطبيقات ميكنة التعلم وإثبات النظريات `<p lang='ar'>automated reasoning theorem proving</p>`

تطبيقات التعرف على الصوت `<p lang='ar'>natural language understanding semantic modeling</p>`

## الشكل (3.4) : الصفحة بعد استخدام الجريجو

### ثانياً: تطبيق محتوى ملفات (HTML)

يتم تطبيق النص المستخلص و ذلك بتوحيد الحروف - مثلاً: قلب كل (ا)، (ا) الى (ا)- وإزالة التشكيل ما عدا الشدة باعتبارها حرف من حروف الكلمة الأصلية ليكون النص بعد ذلك جاهز للتشذيب.

### ثالثاً: تشذيب النص العربي والإنجليزي (Stemming)

بعد تطبيق النص يتم تشذيب النص لرفع كفاءة البحث والذي يقوم بازالة بوادي ولواحق الكلمات وإرجاعها الى أصلها وله أنواع فالتشذيب الكلمات العربية تم استخدام ( light10 stemmer ) وهذا المشذب مستخدم على نطاق واسع في إسترجاع المعلومات العربية.

الواحد	السوابق	نوع المشذب
هـ، انـ، اـتـ، وـنـ، يـنـ، يـهـ، هـ، يـةـ، يـ	الـ، وـالـ، بـالـ، كـالـ، فـالـ، لـلـ، وـ	Light 10

جدول(1.4): السوابق والواحد في 10 Light

وبذلك تكون بحاجة للتعرف على الأسماء والأفعال في سلسلة الحروف المراد تشذيبها.

### رابعاً: يتم تكوين الوثيقة الجديدة المراد فهرستها وإضافتها إلى الفهرس

بعد تطبيق و تشذيب النص المكتوب في الوثيقة لابد من تكوين وثيقة جديدة يتم حفظ المحتوى فيها في شكل مفتاح و قيمة لتصبح الوثيقة قابلة للبحث كما هو موضح في الجدول:

القيمة	الحقل
مسار الوثيقة الأصلي	المسار
رقم الوثيقة (كل وثيقة رقم)	رقم الوثيقة
تاريخ آخر تعديل تم على الوثيقة	التعديل
مسار الوثيقة مضاف الى التاريخ الذي تمت فيه إضافة الوثيقة	تعريف المستخدم
عنوان الموجود في الإشارة <title>	عنوان الوثيقة
النص الموجود في الإشارة <body>	النص

جدول (2.4) : مكونات وثيقة لوسين

ومن ثم تتم إضافة الوثيقة الجديدة في الفهرس الذي تم تحديد موقعه من قبل المستخدم. بذلك تكون كل الملفات التي تحتوي على نصوص عربية وإنجليزية في المجلد الذي حده المستخدم فُهرست في فهارس تم حفظها في المكان الذي حده المستخدم أيضاً.

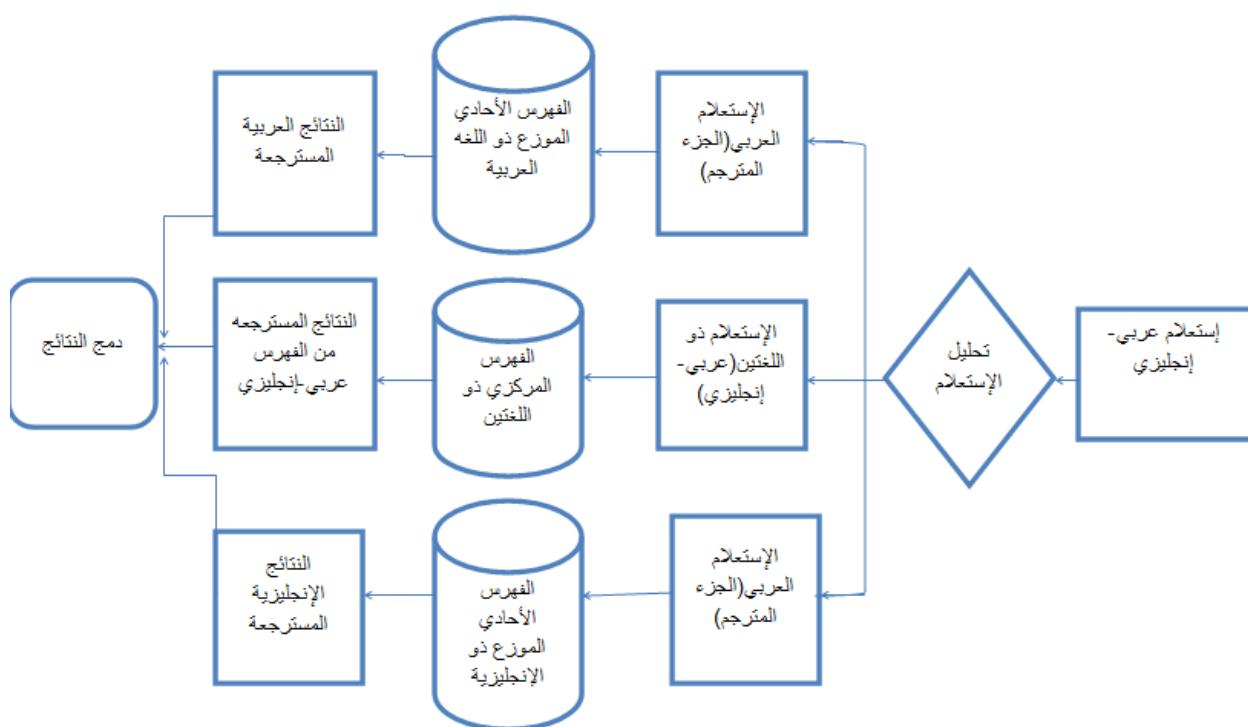
### خامساً: ترجمة الأجزاء العربية والإنجليزية الواردة في طلب المستخدم

يتم إدخال الإستعلام عربي-إنجليزي وترجمة الجزء العربي من الإستعلام عربي-إنجليزي الى اللغة الإنجليزية ودمج الجزء المترجم الى اللغة الإنجليزي من الإستعلام الأصلي سينتج إستعلام عربي والبحث بهذا الإستعلام في الفهرس العربي أحادي اللغة وأيضاً ترجمة الجزء الإنجليزي من الإستعلام عربي-إنجليزي الى اللغة العربية ودمج الجزء المترجم إلى الجزء العربي من الإستعلام الأصلي لينتاج الإستعلام الإنجليزي والبحث بهذا الإستعلام في الفهرس الإنجليزي أحادي اللغة وأخيراً دمج الإستعلام العربي أحادي اللغة مع الإستعلام الإنجليزي

أحادي اللغة لتكوين إستعلام عربي-إنجليزي جديد للبحث به في الفهرس ذو اللغتين العربية والإنجليزية، ثم إسترجاع النتائج في ثلات قوائم ، القائمة الأولى تحوي وثائق عربية والثانية وثائق إنجلizية والثالثة وثائق عربي-إنجليزي .

#### سادساً : تطبيق خوارزميات الدمج على النتائج المسترجعة

بعد عملية البحث يتم إسترجاع ثلات قوائم ، القائمة الأولى يتم إسترجاع نتائجها من الفهرس الأحادي ذو اللغة العربية والقائمة الثانية من الفهرس ذو اللغة الإنجلizية وقائمة ثالثة مسترجعة من الفهرس عربي-إنجليزي . ومن ثم يتم دمج تلك النتائج المسترجعة بإستخدام خوارزميات الدمج التي تم توضيحها سابقا.



الشكل(4.4): عملية البحث و دمج النتائج في الفهارس المركزية الموزعة

## 1.1.5 التصميم التجاري

في هذا الإختبار سيتم قياس كفاءة طريقة الفهرسة التي تم بناءها وتطبيق خوارزميات لدمج النتائج ومقارنتها نتائج تلك الخوارزميات لمعرفة أيهم أحسن كفاءه في عملية الإسترجاع، سيتم قياس الكفاءة عن طريق معيار DCG (Discounted Cumulative Gain) الذي يتم حسابه بناءً على مدى ارتباط الوثيقة المسترجعة بطلب المستخدم، بالإضافة إلى ترتيب الوثيقة ضمن مجموعة الوثائق المسترجعة.

سيتم تطبيق التجربة على مجموعة وثائق HTML تحتوى على نصوص عربية وأخرى تحتوى على نصوص إنجليزية ووثائق تحتوى على نصوص متعددة اللغات (عربية-إنجليزية) مفهرسة بالطريقة المركزية الموزعة ، ليتم البحث فيها وإسترجاع قوائم نتائج من تلك الفهارس ودمج تلك النتائج في قائمة واحدة مرتبة للمستخدم بإستخدام الخوارزميات الثلاث التي تم توضيحها سابقاً و مقارنتها مع بعضها البعض لتحديد أيهما أحسن كفاءة.

## 2.1.5 الهدف من التجربة

معرفة مدى تأثير الطريقة الجديدة المستخدمة في الفهرسة على كفاءة إسترجاع المعلومات، و مقارنة كفاءة الخوارزميات الثلاث التي تم تطبيقها لدمج النتائج المسترجعة.

## 3.1.5 منهجية التجربة

تم إستخدام محرك بحث لوسين في فهرسة الوثائق بطريقة تسمح بالبحث فيها، فأولاً تمت تنقية هذه الوثائق من الصور، و الرموز، و كل العناصر ما عدا النصوص العربية والإنجليزية و بعد ذلك تم تطبيق تلك النصوص التي تكون محتويات كل وثيقة لتصبح جاهزةً لعملية التشذيب. ثالثاً تم تشذيب محتويات هذه الوثائق بإستخدام لایت10، ومن ثم عملية الفهرسة.

تم بناء فهرس مركزي موزع للبحث فيه، ومن ثم تم جمع ثمانية إستعلامات (quires) من مستخدمين مختلفين، وفي هذه الإستعلامات تمت ترجمة الأجزاء العربية والإنجليزية وتم تكوين الإستعلام العربي والإستعلام الانجليزي والإستعلام متعدد اللغات (عربي-إنجليزي)، و إستخدام تلك الإستعلامات للبحث في ذلك الفهرس الذي تم بناءه.

من ضمن الوثائق المسترجعة تم إختيار أول عشر وثائق، و تم تقييم كل وثيقة برقم في المدى (0—5) على حسب مدى إرتباطها بطلب المستخدم ، ليتم بعد ذلك حساب معيار DCG (Discounted Cumulative Gain)

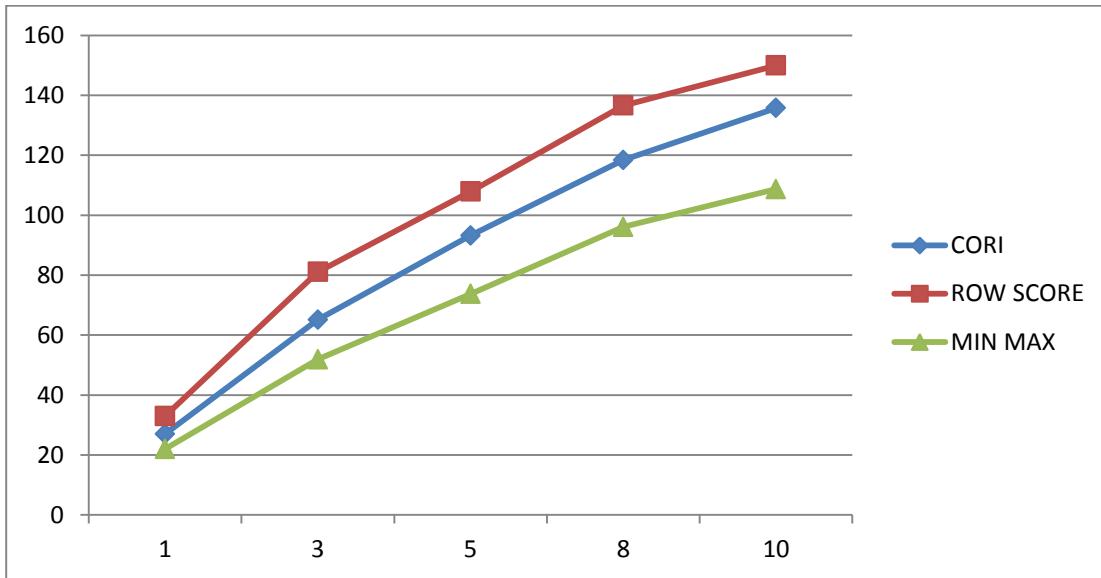
وهو معيار لκفاءة إسترجاع المعلومات في محركات البحث، حيث يتم حسابه لكل طلب في كل خوارزمية من خوارزميات الدمج الثلاثة، لتتم من خلاله مقارنة κفاءة.

## 4.1.5 النتائج والمناقشة

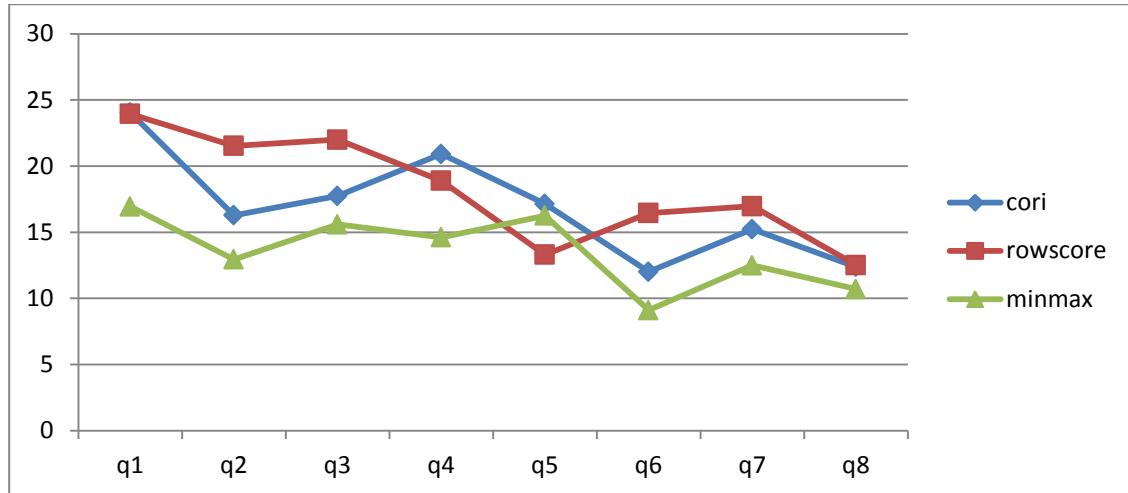
وبعد إجراء عملية البحث بإستخدام الثمانية إستعلامات في الفهرس المركزي الموزع وإسترجاع قائمة نتائج مرتبة من كل فهرس ومن ثم تطبيق خوارزميات المجموع الثلاثة عليها فكان متوسط معيار DCG لكل من الخوارزميات الثلاثة كالتالي:

- متوسط DCG بإستخدام خوارزمية CORI = 16.97
- متوسط DCG بإستخدام خوارزمية RawScore = 18.74
- متوسط DCG بإستخدام خوارزمية NormalizedScore = 13.59

بالنظر للشكل (1.5) يمكننا مقارنة متوسطات DCG بصورة أوضح، وينتجي لنا أن خوارزمية Raw Score تفوقت على الآخريات بحصولها على أكبر متوسط من DCG على مستوى عدد الوثائق في قائمة النتائج المسترجعة وذلك على مستوى وثيقه واحده وعلى مستوى ثلث وثائق وعلى مستوى خمس وثائق وعلى مستوى ثمانى وثائق وعلى مستوى عشر وثائق.



الشكل(1.5):متوسطات معيار DCG للخوارزميات الثلاث



الشكل(2.5):نتائج الخوارزميات للثلاث استعلامات

كما ذكر سابقاً أن خوارزمية CORI هي الأفضل وكان ذلك بناء على استخدام إستعلام أحادي اللغة ولكن بالنظر للشكلين (2.5)، (1.5) أعلاه وجد أن خوارزمية Raw Score تعطي نتائج أفضل وذلك لأن عملية البحث هنا تتم بإستخدام إستعلام ذو لغتين(عربي\_إنجليزي) مما أدى إلى زيادة طول الإستعلام الذي يتم البحث به في الفهرس متعدد اللغات والذي بدوره أثر على وزن الوثائق وبالتالي أثر على ترتيب الوثائق حسب الأهمية.

### 5.1.5 جوانب القصور في نتيجة التجربة

- تم إجراء التجربة على عدد غير كاف من الوثائق العربية.
- من المفترض أن يتم نوع من اختبارات الأهمية على نتائج طلبات المستخدمين، لكن اختبارات الأهمية يجب أن تتم على الأقل على 31 طلب مستخدم، وفي تجربتنا هذه إستخدمنا 8 طلبات فقط، لذلك لم يكن من الممكن إجراء اختبار الأهمية.

### 6.1.5 الصعوبات التي واجهت البحث

عدم القدرة على الحصول على Google translator API حيث تم إغلاقها من قبل الشركة مؤخراً لذا تم بناء قاموس يحتوي على أكثر الكلمات استخداماً في عملية البحث لاستخدامه في عملية ترجمة الاستعلام.

## 2.5 التوصيات

- ايجاد طريقة لتعديل معادلة تطبيع الوزن في خوارزمية cori وذلك بإدخال معامل طول طلب المستخدم لإعطاء نتائج أفضل
- تطبيق الفهرسة و البحث على مجموعة أكبر من الوثائق للحصول على نتائج أفضل عند مقارنة خوارزميات دمج النتائج.
- إجراء عملية البحث عن طريق طلبات مستخدمين تتجاوز الثلاثينيات، للتمكن من إجراء اختبار الأهمية.

## الخاتمة

بعد توفيق من الله سبحانه وتعالى تم استكمال هذا البحث الذي يستخدم الطريقة المركزية الموزعة لفهرسة الوثائق و التي تمكننا من فهرسة الوثائق متعددة اللغات أو أحادية اللغة بطريقة فعالة تزيد من كفاءة الإسترجاع في انظمة إسترجاع المعلومات، وأيضاً تمكن المستخدم من كتابة الاستعلام باللغة التي يراها مناسبه أو التعبير عن مصطلحاته باللغة التي يعرفها. ومن ثم تم تطبيق خوارزميات دمج النتائج المسترجعه من عملية البحث بهذا الاستعلام وتكوين قائمة واحدة تكون فيها النتائج مرتبه على حسب صلتها بالإستعلام وأهميتها بالنسبة للمستخدم بحيث تكون الوثائق ذات الصلة الاعلى بالإستعلام و ذات الأهمية الاكبر بالنسبة للمستخدم في بداية القائمة الناتجة.

