# **ResultsMerginginaDistributedSearchSystem**

V.V.KLUEV TheCoreandInformationTechnologyCenter TheUniversityofAizu TsurugaIkki -machiAizu -WakamatsuCity Fukushima965 -8580 JAPAN

Abstract: Mechanismsofresultsmerginginadistributedsearchsystemhavebeendiscussedinthispaper. Such mechanisms' aimsaretoselectthemostrelevantdocumentsretrievedby differentservers and putthemonthe topofthelistreturned to the endu ser. Inour calculation, we use the clustering technique to divider etrieved results into several homogeneous groups and ametric conthe base of the vector space model to arrange items inside each group. The mainthrust of our approachis both the proposed clustering metric and the ranking metric. Our mechanisms have been implemented in the OASIS system, asystem which be long stothe distributed search systems in the Internet.

Key-Words: -ResultsMerging,SearchEngine,DistributedSystem,Metasearch,Ran king

# **1Introduction**

ItiswellknownthattheInternetcontinuesto expandrapidly.Therearedozensofpowerful searchtoolsonthenetavailablewithoutcost. Nevertheless,searchingforinformationisstill inefficient.Howisitpossibletofindapp ropriate informationeasilyandquickly?Thismainquestion isstillunanswered.Thereareaplentyof explanationsonthenethowtosearch.Oneofthem canbefoundat[3].Thecurrenttrendindeveloping newtoolsistogetmoreandmoreinstruments designedonthebaseofthedistributedarchitecture. Themainadvantagesofadistributedsystem againstasystemusingthecentralizedarchitecture areasfollows:

- Administrationofasystemshouldbe easier.
- Theindexofawholesystemshouldbe largerwh encomparedtootherapproaches.
- Resultsofasearchhavetobemore accurate.

Thekeyproblemsofdistributedsearcheshavebeen discussed in the following citations [2,45,8,10, 12,13,15]. Problems pointed to are reflected here:

- 1) Howcanasearcha ccuratelydeterminea smallnumberofpotentiallyuseful collections(localsearchengine)toinvoke foreachuserquery?
- 2) Howcananenginesearchinsidethiskind ofcollectionsconsistingoftopicidentical documents for relevant ones?

Theanswerstothe sequestionsdependonthe natureofthesystem.Distributedsystemscanbe classifiedasfollows:

- 1) Generalsystems. Thespacehasbeen dividedonthebasegeographical principle. Forexample, one server can index documents from sources located incentral Europe; the responsibility of the second one is Japan, etc. The classical example is the HARVEST system [10].
- Metasystems. Theseformaseparatetype inthedistributedclasssystems. They usuallysupportunifiedaccesstomultiple localsearchengines. T hesesystemsdonot maintaintheirownindexonwebpages. Particularsystemscanbeusedas extensionstothemostpopularpowerful searchengines. Forexample, SearchPadis ac lientsidemetasearchtoolthatsearches allmajoronlinesearchenginesato nce usingdefinitionspresetbytheuserthat allowforincrediblyexactmatches [2].
- Topicspecificsystems.Eachlocalserver hascollecteddocumentsrelativetoa specialtopic.Forexample,thetopicof searchmaybealgorithms,sightseeing,etc. OASISisanexampleofthiskindof system[1].

Acommonanddifficulttaskforeachkindof systemsisresultsmerging.Theaimofthisprocess is to combine results returned from selected collections(localsearchengines)intoafinalsingle rankedlist.Th emajorobstaclestowardresolving thistaskareasfollows:

- The different local servers can return the same documents because theservers' databases intersect. Duplicates have to be eliminated.
- Metricsbeingusedbydifferentservers cannotbecorrelated, evenserversusethe samerankingmethod.Thefullsetof returneddocumentshastobeevaluated somehowbeforepresentingtheresultsto theenduser.

Inthispaper,wepresentanewtechniquetoresolve theaforementionedtasks.Thistechniquehas shownpromisingresults.Ithasbeenappliedtothe OASISsystem[1].

Therestofthepaperhasbeenorganizedasfollows. Section2describessomewell -knowntechniques forresultsmerging.Section3presentsour approachindetail.Section4discussesthed etailsof evaluationswecarriedoutinourtestbeds.Thefinal sectiongivescommentsonexperimentalresults.

# 2RelatedWork

Thetaskofmergingresultswasdiscussedin severalcitations[2,4,5,6,7,8].Proposalswith differentsolutionscanbefo undthere.Inthis section,wewillreviewthemshortly. Inrespondtoaqueryeachserver(collection) returnsalistofdocumentsretrieved.Thislist usuallyconsistsofthefollowingitems:

- URL;
- Thetitleofthedocument;
- Ashortdescriptionofthedoc ument(the aimofthisitemistohelptheuserto estimatethedocument);
- Adocumentscorecalculatedaccordingthe measureappliedontheserver;
- Auxiliaryinformation:specificationofthe documentformat,dateofindexing, documentsize,etc.

### 2.1Dupl icateRecognition

Thesimplestvariantoftheduplicaterecognition taskisdetectionofthesameURLs.Thisisnot easilyaccomplishedifdocumentsarelocatedon differentservers("mirrors")andpresentedin differentformats(forexampleps,pdf,html or text).Thereisnotadequateinformationtomakea decision.Inanycase,averysimplemethod proposed in [9] can be applied. If names of documentfilesdifferonlyinsuffixesdocuments canbeconsideredasthesame.

### 2.2RankingStrategy

Currentappr oachestocombineresultsreturnedby differentsearchenginesintoasinglerankedlistcan bedistinguishedinthefollowingway:

- 1) Rowscoremerging;
- 2) Roundrobinmerging;
- 3) Mergingonthebaseofcollectionscores;
- 4) Mergingonthebaseofglobalsimilarities
- 5) Clustering.

Authorsofthefirstapproach[6]assumethat documentcollectionshavebeenindexedusingthe samemodelanddocumentscoresobtainedasa resultofthesearcharecomparableacrossall collections. Thesescoresarethenusedtomerge searchresultsintoasinglelist. Thebasicassumptionforthesecondapproach lookslikethis:Eachdocumentcollectionconsists ofapproximatelythesamenumberofrelevant itemsandtheyareequallydistributedwithinresults lists[4].Fromthishypothesis afinalresultslistcan becreatedinaroundrobinmanner.

Thesystemsusingthethirdmergingstrategy computeascoreforeachcollectionandforagiven query.Thiscollectionscoreisthenutilizedto modifythescoreattachedtoeachdocument.The detaileddescriptionandnecessaryformulascanbe foundin[10].

Asnotedin[5],threeaforementionedapproaches areallheuristicsandtheycannotguaranteethe selectionofallpotentiallyusefuldocumentsfora givenquery.

Anapproachbasedonag lobalsimilaritythreshold aimstoretrieveallpotentiallyusefuldocuments fromeachselectedcollection.Theproblemisthat documents(fullcontent)mayneedtobefetchedto thesearchsystemtoenablethecomputationof theirglobalsimilarities.Th eadvantageisthathigh qualitymergingcanbeachieved[5].

Theclusteringstrategyhasbeenappliedtothe OASISsystem.Itsaimistosplitresultsinto varioussets(clusters).Thedocumentsinsideeach grouparesimilartoeachotherbytopicsthey are relatedto.Onerepresentativefromeachclusterhas beenincludedintoafinalresultslist.Theclustering isexecutedbythecompetitivelearningneural network[1].Thedisadvantageofthissolutionis that huge document profiles are sent though the Internet.Moreover,LSI(LatentSemanticIndexing) similarcalculationisrunningbeforeclusteringto reducethesizeofdocumentprofiles.Asaresult,it decreasestheperformanceofthesystem.Collected statisticsshow:Aneuralnetworkdoesnotpro duce wellsplitclusters.

Themainproblemininformationretrievalisthe lackofapowerfullanguagemodel.Fromthis,most ofapplicationmethodsareempiricalandbasedon heuristics.Asaresult,inourcase,itispractically impossibletocompareth eaforementioned approachesbecausetheyweredesignedfor differentsystemsandweretestedusingdifferent datasets.

# 3MechanismsUsed

#### 3.1CommonNote

Thebasicassumptionsforourapproachareas follows. Usersusually interact with a search system severaltimestofindinformationtheylookfor. For eachchunkofinformationneedtheymaysubmit several queries, often by a process of query refinement. Theaimofeachiterationistonarrow thesearchtopicanddiscardunnecessary documents. Thebes tstrategy for a systemistos plit returnedresultsintothevarioussetsanddeliverto theuserthemostimportantrepresentativesfrom these sets without losing relevant data. The systemshouldminimize the volume information sent thoughtheInternetd uringeveryiteration.Thekey ideasfromtheapproachesmentionedinSection2 havebeenutilizedhere.Ourproposedmechanisms takeintoaccount:

- Thetitleoftheeveryreturneddocument;
- Aselectedpartofthedocumenttextto browsebyauser;
- Sizeof thedocumentinbytes;
- Scoreforeachselectedcollectionforquery propagation;
- Termfrequencyofeachquerytermin everyreturneddocument;

Thefirstthreecomponentsareastandarddefacto forbrowsinginthesearchengineworld.Theusage ofthela sttwowillnotsignificantlyincreasethe volumeofinformationbeingsentthoughthe Internet.

#### 3.2ClusteringScheme

Itisknownthatclusteringisacommonwayto divide retrieved documents into several sets. A

numberofstudies[13,14]havereporte don experimentsrelatedtoanautomatictextclustering technique.Researchers[12]testedWard's clusteringandobtainedpromisingresults.We adoptedatechniqueusedbythem.Themain clusteringalgorithmisverysimple,anditis practicallyidentical toWard'stechnique.

*DescriptionoftheAlgorithm:* LetXbethesetof documentstobeclustered;letCbethesetof clustersandNbethepredefinednumberofclusters (thenumberofreturnedresultstobepresenttothe user).Thealgorithmstartswi thaseparatecluster foreachdocumentfromX.Ineachstep,thetwo mostsimilarclustersC <sub>i</sub>andC <sub>j</sub>aredeterminedand mergedintoanewclusterC <sub>new</sub>.Thealgorithm terminateswhenapredefinednumberofclustersN containingalldocumentsfromXhasbee nformed.

Metricsandthresholdsusedinthesecalculations arekeycomponents.Ourmetricsisbasedonthe widelyusedvectorspacemetric.Inpursuanceof [11]weestimateanumberofwordsinthe documentwritteninEnglishas:

$$WN = \frac{CN}{5}$$

where CN is the size of the document in bytes. Allaforementioned components from document description returned by the collection have been used in our distancemetric. This description was considered as a document. We modified the standard tf \* idf formulas to compute term weights [16]: the final term weight is the product of tf \* idf weight and the collection score from the retrieval source of the document:

$$w_final_{ij} = w_{ij} * col_score_j$$

where

- $w_{ij}$  is tf \* idf weightofterm *i* indocument *j*
- *col\_score*<sub>j</sub> is acollectionscore(a collectionselectionsubsystem calculates thisscore).Notethatthecollectionscore is less that 1.

#### 3.3Ranking

Tomakeafinaldecisionaboutdocumentrelev ance scores, we modified the formula proposed in [8].

$$score(D,Q) = \frac{nb_q + c_1 * nb_occ + c_2 * doc_score + c_3 * col_score}{\sum_{i} score(D_i,Q)} (*)$$

Weagreewithauthors[8]thatthenumberofquery termsincludedineachdocumentandfrequencies ofquerytermsindocumentsaregoodrelevant indicators.Thefollowingcomp onentsshouldbe takenintoaccountaswell:a)thecollectionscores andb)thedocumentscoresassignedbythe collectionfromwhichthedocumentwasretrieved. Onthebasisofthisassumption,wecalculatethe documentscoreaccordingformula(\*).Para meters foreachdocumentDareasfollows:

- *nb<sub>q</sub>* is the number of query terms presented in D;
- *nb\_occ* isthetotalnumberofoccurrences ofquerytermsinD;
- doc\_score isadocumentscoreassigned bythecolle ctionfromwhichthedocument wasretrieved;
- col\_score isacollectionscore(a collectionselectionsubsystemcalculates thisscore);
- $c_1, c_2$  and  $c_3$  are constants, which are set to 0.100.

Aforementionedform ulaisappliedtoeachcluster centroid.Documentspresentedtotheuserare arrangedaccordingtotheirscores.

### **4ResultsoftheTests**

Theproposedapproachhasbeentestedusingthe OASISsystem.Thesetestsarediscussedinthis section.Thefollowing configurationofthesystem wasdedicatedforourtests:Threelocalservers(the installationplacebeingAizuWakamatsuCity)and oneremoteserver(atKoriyamaCity).Thedistance betweenthesecitiesisabout60km.Thetesttopic specificcollections [9,17]consistingofthereal Internetdatawereinvolvedinourexperiments. Table1describesadistributionofdocument collectionsinstalledonservers.Theaimofthese testsistocomparethepreviouslyimplemented methodswithproposedonesinthis paper.

Weassumethatthefollowingareagoodaccuracy indicatorofthemergingtechnique:thenumberof actuallyrelevantdocumentsonthetopofthelist presentedtotheuser,andthenumberofactually relevantdocumentsineachcompiledclusteras well.Theseparametersareveryimportantfromthe user's point of view. We used a set of short queries, similartothosesubmittedtothesearchengine. Theyconsistofone,twoandthreewords.These queriesreflecttherealsearchprocessontheWeb. Weconductedourtestsasfollows:

- 1) Aquerywassubmittedtothesystem;
- 2) Thesystemfetchedresultsofthesearch;
- 3) Afterdeletingtheduplicatesreturned resultswereclustered;
- 4) Itemsinsideeachclusterwerearranged accordingtorelevancescore;
- 5) Centroidso f10clustersweresubmittedto theuser;
- 6) Theusercouldseeuptothreedocuments fromeverycluster;
- 7) Theevaluationoftheresultsaccuracywas doneonlybyhumaninspection.

Table1Locationofthecollections							
Servers	Collections	Numberof					
		documents					
Aizu:1	ProgrammingLanguages	7659					
Aizu:2	Algorithms	7775					
Aizu:3	Travelogues	226					
	Linux&Unix	488					
	InformationRetrieval	202					
	ResearchGroups	811					
	Physics	467					
	CardGames	798					
	Museum	444					
	Monitors	70					
Koriyama	ProgrammingLanguages	445					
-	Cars	427					

Table2sho wsthemainresultsofourtests.

Thesystemretrievedresultsforthesamequeryset usingtheoldermergingmethodandthesystem presentlybeingdiscussed.Inourtests,document setsbeforeamergingprocesswerethesamefor bothmethods.Thenumber ofretrieveddocuments byeachserverwasbetween0and20.Aswecan see,columns2 –4presentthelengthofqueriesin words;columns5and7showthenumberof relevantitemsamongthetopthreedocumentsfrom thelistpresentedtotheuser.Columns6 and8point toaquoteofrelevantdocumentsamongthetop five.Ingeneral,thenewmethodproducesbetter results.Theamountofinformationsentviathenet was reduced significantly without losing the

Topic	Querylength			Searchresults						
	1W 2W 3		3W	Oldmethod		Newmethod				
				Top3	Top5	Top3	Top5			
1	2	3	4	5	6	7	8			
ProgrammingLanguages	5	1	4	2	3	3	3			
Algorithms	4	2	1	2	2	2	3			
Cars	1		2	1	2	2	3			
Travelogues			2	0	1	1	1			
Linux&Unix		2		0	1	0	1			
InformationRetrieval		2	1	0	1	0	2			
ResearchGroups	1			0	0	0	1			
Physics		2		1	0	0	1			
CardGames	2	1	2	2	2	1	2			
Museums	4	3	2	3	3	3	3			
Monitors		1		0	1	0	0			

Table2.Resultsofthesearch

accuracyofthesearch. Thesecondadvantageisin reducing calculations to merge the results. Our new

 $O(N^2)$  time[15]. clusteringalgori thmrunsin Usingthisalgorithmwediscardedthetime consumingclusteringbyneuralnetworks[1]and timeconsumingLSIcalculationaswell.Thenew methodshowedthepoorerresultsthantheolder oneinthecaseofsea rchingforinformationrelated topicsofMonitors,CardGamesandPhysics.Both methodsproducedthesameresultsfortopicsof MuseumsandLinux&Unix.Thedocumentsthat aretheclosesttothecentroidswerealsoanalyzed. Thenumberofrelevanthitsfrom themisthesame onaverageforbothmethods(2 -3usingeach methods).

# **5**Conclusion

Thispaperintroducesanewtechniquetomerge searchresults indistributed systems. It utilizes the ideaofpreliminaryclusteringandthenarranging itemsinsideea chcluster.Clusterrepresentatives canbesubmittedtotheuserintheusualsearch enginemanner.Allcalculationsarebasedon standarddefactoinformation:thedocumenttitle, itssize,aselectedpartofthetexttobrowsebya calculations, this technique user,etc.Fromthese canbeappliedtometasearchsystemsaswell. ThistechniquewasincorporatedintotheOASIS system.Resultsshownduringourpreliminarytests arepromising. They are qualitatively at least as goodasresultsproducedbythe timeconsuming neuralnetworkclusteringandLSIcalculations.

#### References:

- OASIS: DistributedSearchSysteminthe Internet, EditedbyA.Patel, L.Petrosjan, W. Rosenstiel, St.Petersburg: St.Petersburg StateUniversityPublishedPress, 1999, -614 p.(ISBN5 -7997-0138-0)
- KrishnaBharat,SearchPad:Explicit
  CaptureofSearchContexttoSupportWeb
  Search, InProceedingsof9 thInternational
  WorldWideWebConference ,Amsterdam, May15 –19,2000.
- [3] <u>http://searchwebmanagement.techtarget.co</u> m/sDefinition/0,,sid27\_gci212955,00.html
- [4] VoorheesE.M.,GuptaN.K.,andJonson LairdB.,LearningCollectionFusion Strategies, *inProceedingsoftheACM SIGIR*'95,1995,pp .172 –179.
- [5] ZonghuanWu,WeiyiMeng,ClementYu, ZhuogangLi,TowardsaHighly -Scalableand EffectiveMetasearchEngine, InProceedings of10thInternationalWorldWideWeb Conference,2001
- [6] KworkK.L.,GrunfeldL.,andLewisD.D., TREC-3Ad -hoc:RoutingR etrievaland ThresholdingExperimentsUsingPIRCS, *In ProceedingsofTREC -3,1995,pp.247 -255.*
- [7] L.Gravano,C.Chang,H.Garcia -Molina, andA.Paepcke. *STARTS:StanfordProposal forInternetMeta -Searching*.ACMSIGMOD Conference,Tucson,May1997,pp.207 -218.
- YvesRasolofo,FaizaAbbaci,Jacques
  Savoy:ApproachestoCollectionSelection
  andResultsMergingforDistributed
  InformationRetrieval, CIKM2001, pp.191 198

- [9] V.Kluev,CompilingDocumentCollection fromtheInternet, *ACMSIGIRForum*,Vol. 34,Num ber2,Fall2000,pp.9 –14.
- [10] CallanJ.DistributedInformationRetrieval. InAdvancesinInformationRetrieval,Edited byW.B.Croft ,KluwerAcademicPublisher, 2000,pp.127 –150.
- [11] ReadRace <u>http://</u> www.epixmed.com/hunsocma/ReadRace/
- [12] RoussinovD.G.,ChenH,Document ClusteringforElectronicMeetings:An ExperimentalComparisonofTwo Techniques, *DecisionSupportSystems*,27 (1999),pp.67 -79.
- [13] MartinEster,Hans -PeterKriegel,Jourg Sander,Miha elWimmer,andXiaoweiXu, IncrementalClusteringforMininginaData WarehousingEnvironment, *inProceedingsof the24* <sup>th</sup>VLDBConference ,USA,1998.
- [14] ChristopherD.ManningandHinrich Schutze, *FoundationsofStatisticalNatural LanguageProcessing* ,TheM TIPress, England,2000,680p.
- [15] LarkeyL.S.,ConnelM.E.,andCallanJ, CollectionSelectionandResultsMerging withTopicallyOrganizedU.S.Patentsand TRECData,inProceedingsofCIKM'200, pp.282 –289.
- [16] G.Salton, Automatictextprocessing:the transformation,analysis,andretrievalof Informationbycomputer .Reading, MA:Addison-Wesley,1989.
- [17] V.Kluev,SourceSelectioninaDistributed SearchSystem, inV.Kluev,N.Mastorakis, editors,TopicsinAppliedandTheoretical MathematicsandComputerScience,WSEAS Press,2001,pp.293 –298.