# Multi-document Summarization via Information Extraction

Michael White and Tanya Korelsky
CoGenTex, Inc., Ithaca, NY
*mike,tanya@cogentex.com*

Claire Cardie, Vincent Ng, David Pierce, and Kiri Wagstaff
Department of Computer Science, Cornell University, Ithaca, NY
*cardie,yung,pierce,wkiri@cs.cornell.edu*

## 1 Introduction

Although recent years has seen increased and successful research efforts in the areas of single-document summarization, multi-document summarization, and information extraction, very few investigations have explored the potential of merging summarization and information extraction techniques. This paper presents and evaluates the initial version of RIPTIDES, a system that combines information extraction (IE), extraction-based summarization, and natural language generation to support user-directed multi-document summarization. We hypothesize that IE-supported summarization will enable the generation of more accurate and targeted summaries in specific domains than is possible with current domain-independent techniques.

In the sections below, we describe the initial implementation and evaluation of the RIPTIDES IE-supported summarization system. We conclude with a brief discussion of related and ongoing work.

## 2 System Design

Figure 1 depicts the IE-supported summarization system. The system first requires that the user select (1) a set of documents in which to search for information, and (2) one or more scenario templates (extraction domains) to activate. The user optionally provides filters and preferences on the scenario template slots, specifying what information s/he wants to be reported in the summary. RIPTIDES next applies its Information Extraction subsystem to generate a database of extracted events for the selected domain and then invokes the Summarizer to generate a natural language summary of the extracted information subject to the user's constraints. In the paragraphs below, we describe the IE system and the Summarizer in turn.

**IE system .** In the RIPTIDES evaluation of Section 3, we assume "perfect" output templates from the IE system. As a result, we focus here on scenario template design considerations that directly affect the quality of the summaries produced rather than the architecture of the IE system, [1] which is of less importance for the purposes of this paper. (We anticipate that the final paper will include both a description of the IE system and an evaluation of the Summarizer using its output.)

The domain selected for the initial system and its evaluation is natural disasters. A top-level natural disasters scenario template contains: document-level information (e.g. *docno*, *date-time*); zero or more *agent* elements denoting each *person*, *group*, and *organization* in the text; and zero or more *disaster* elements. *Agent* elements encode standard information for named entities (e.g. *name*, *position*, *geo-political unit* ). For the most part, *disaster* elements also contain standard event-related fields (e.g. *type*,

---

[1] In brief, the RIPTIDES IE system uses a traditional architecture (Cardie, 1997): a preprocessor finds sentences and tokens; a parser identifies syntactic structure; syntactico-semantic pattern-matching identifies text fragments for extraction; coreference resolution guides template creation. We are investigating the use of weakly supervised learning techniques (e.g. Riloff and Jones, 1999; Thompson et al., 1999) for the automatic construction of each IE system component.

*number*, *date*, *time*, *location*, *damage* sub -elements). The final product of t he RIPTIDES system, however, is not a set of scenario templates, but a user -directed multi -document summary. This difference in goals influenced a number of template design issues. First, disaster elements must distinguish different reports or views of t he same event (from multiple sources) —the system creates a separate *disaster* event for each such account —and should include the *reporting agent* , *date*, *time*, and *location* whenever possible. In addition, *damage* elements (i.e. *human* and *physical effects* ) are best grouped according to the reporting event. Finally, a slight broadening of the IE task was necessary in that extracted text was not constrained to noun phrases. In particular, adjectival and adverbial phrases that encode *reporter confidence*, and s entences and clauses denoting *relief effort* progress appear beneficial for creating informed summaries. The disaster elements extracted for each text are provided as input to the summarization component.
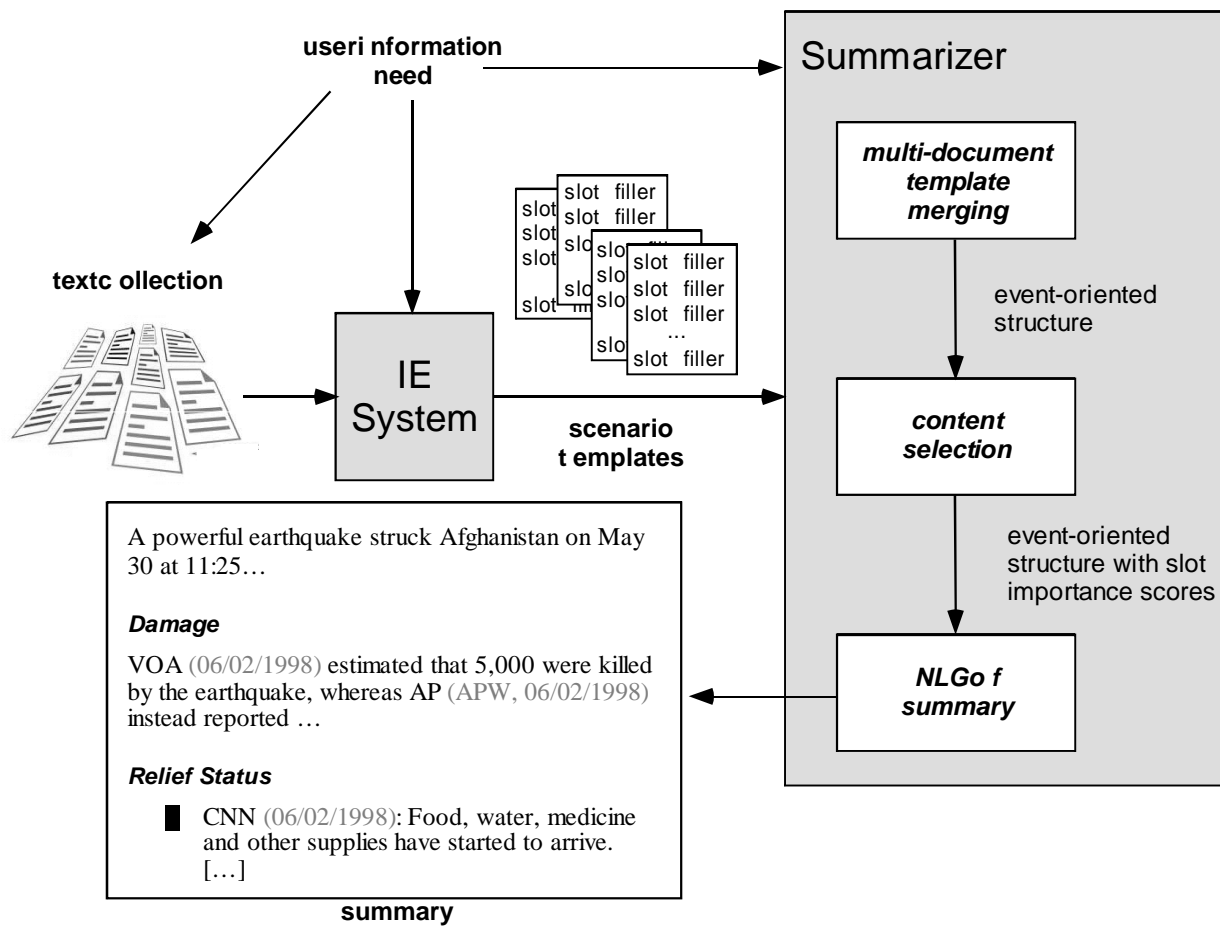


**Figure1**

**The Summarizer.** A sample summary generat ed by the initial version of RIPTIDES is shown in Figure 2 below.[2] The Summarizer produces each summary in three main stages. In the first stage, the output

---

[2] The sample summary was produced for the second evaluation test case (see Section 3) —a summary of all test articles th at emphasizes factual information. For space reasons, some of the summary is not shown. In this example, preference is given to the basic facts (the first paragraph), then to the overall damage information (vs. damage reports in specific locations), then to relief information (which has a slight preference over the remaining facts). Note that the relatively poor coherence of the second paragraph stems primarily from the overly simplistic heuristics for comparing damage reports; in the current version of t he system, more robust rules are used (see section 4). Also

templatesaremergedintoanevent    -orientedstructure,whilekeepingtrackofsourceinformation.    The
mergeoperationcurrentlyreliesonsimpleheuristicstogroupextractedfactsthatarecomparable;for
example,duringthisphasedamagereportsaregroupedaccordingtowhethertheypertaintotheeventasa
whole,orinsteadtodamageinthesame    particularlocation.Heuristicsarealsousedinthisstageto
determinethemostrelevantdamagereports,takingintoaccountspecificity,recencyandnewssource.
Withreliefslots,word  -overlapclusteringisusedtogroupslotsfromdifferentdocument    sintoclusters
thatarelikelytoreportsimilarcontent.Inthesecondstage,abaseimportancescoreisfirstassignedto
eachslotbasedonacombinationofdocumentposition,documentrecencyandgroup/clustermembership.
Thebaseimportancescores   arethenadjustedaccordingtouser   -specifiedslotpreferencesandmatching
criteria.Theadjustedscoresareusedtoselectthemostimportantslotstoincludeinthesummary,subject
totheuser  -specifiedwordlimit.Inthethirdandfinalstage,thes     ummaryisgeneratedfromtheresulting
contentpoolusingacombinationoftop   -down,schema -liketextbuildingrulesandsurface   -oriented
revisions.(Atpresent,however,extractedreliefsentences/clausesaresimplylistedindocumentorder.)

---

## EarthquakestrikesnorthernAfghanistan

ApowerfulearthquakestruckAfghanistanonMay30at11:25.Theearthquakewascenteredinaremotepartofthe
countryandhadamagnitudeof6.9ontheRichterscale.

### *Damage*

VOA (06/02/1998)estimatedthat5,000wereki   lledbytheearthquake,whereasAP    (APW,06/02/1998) instead
reportedanywherefrom2,000to5,000peopledead.CNN        (06/02/1998)insteadreportedupto4,000peopledied,
whileI (PRI,06/01/1998) estimatedseveralthousandpeoplemayhavedied.I       (PRI,0 6/01/1998)estimatedthat
thousandswerelefthomeless.[…]

### *Relief*

#### Status

- CNN (06/02/1998):Food,water,medicineandothersupplieshavestartedtoarrive.[…]

#### Problems/Obstacles

- VOA (06/03/1998):BadweatherinAfghanistanishamperingeffortstoreach     victimsoflastweek's
devastatingearthquake.[…]

### *FurtherDetails*

Heavyaftershocksshooknorthernafghanistan.Landslidesormudslidesalsohitthearea.[…]

**Figure2**

---

TheSummarizerisimplementedusingtheApacheimplementationofXSLT(Apache,2        000)and
CoGenTex'sExemplarsFramework(WhiteandCaldwell,1998;White,2001).      TheApacheXSLT
implementationprovidedaconvenientwaytorapidlydevelopaprototypeimplementationofthefirsttwo
processingstagesusingaseriesofXMLtransformation      s.Inthefirststepofthethirdsummary
generationstage,thetextbuildingcomponentoftheExemplarsFrameworkconstructsa"roughdraft"of

---

notethatthereferenceto"I"inthisparagraphshouldhavebeentothereporterTonyKahn;ourheuristicwastouse
theinitialreferencetoasource,and"I"happenedtobetheinitialreferen        ceinthearticle.

thesummarytext. [3]Inthisroughdraftversion,XMLmarkupisusedtopartiallyencodetherhetorical, referential,semanticandmorpho -syntactic structureofthetext.Inthesecondgenerationstep,the Exemplarstextpolishingcomponentmakesuseofthismarkuptotriggersurface -orientedrevisionrules thatsmooththetextintoamorepolishedform.Adisti nguishingfeatureofourtextpolishingapproachis theuseofabootstrappingtooltopartiallyautomatetheacquisitionofapplication -specificrevisionrules fromexamples;cf.White(2001)fordetails.

## 3EvaluationandInitialResults

Toevaluate thei nitialversionof theIE -supportedsummarizationsystem,weusedTopic89fromthe TDT2collection —25textsonthe1990Afghanistanearthquake.Eachdocumentwasannotated manuallywiththenaturaldisasterscenariotemplatesthatcomprisethedesiredo utputoftheIEsystem.In addition,treebank -stylesyntacticstructureannotationswereaddedautomaticallyusingtheCharniak (1999)parser.Finally,MUC -stylenounphrasecoreferenceannotationsweresuppliedmanually.All annotationsareinXML.

Next,theTopic89textsweresplitintoadevelopmentcorpusandatestcorpus.Thedevelopmentcorpus wasusedtobuildthesummarizationsystem;theevaluationsummariesweregeneratedfromthetest corpus.SummariesgeneratedbytheRIPTIDESsystemwer ecomparedtoasimple,sentence -extraction multi-documentsummarizerthatreliesonlyondocumentposition,recency,andwordoverlapclustering. Inaddition,theRIPTIDESandBaselinesystemsummarieswerecomparedagainstthesummariesoftwo humanaut hors.Allsummariesweregradedwithrespecttocontent,organization,andreadabilityonan A-Fscalebyfourgraduatestudents/professionals,allofwhomwereunfamiliarwiththisproject.

Eachsystemandauthorwasaskedtogeneratefoursummaries ofdifferentlengthsandemphases:(1)a 100-wordsummaryoftheMay30andMay31articles;(2)a400 -wordsummaryofalltestarticles, emphasizingspecific,factualinformation;(3)a200 -wordsummaryofalltestarticles,focusingonthe damagecaused bythequake,andexcludinginformationaboutreliefefforts,and(4)a200 -wordsummary ofalltestarticles,focusingonthereliefefforts,andhighlightingtheRedCross'sroleintheseefforts.

| RIPTIDES | Baseline | Person1 | Person2 |
|---|---|---|---|
| C | D/D+ | A- | B+ |

**Table 1**

| | RIPTIDES | Baseline |
|---|---|---|
| Overall | 1.92±0.53 | 1.16±0.48 |
| Content | 2.15±0.96 | 1.77±1.11 |
| Organization | 1.99±1.02 | 0.48±0.49 |
| Readability | 1.81±0.71 | 1.19±0.95 |

**Table 2**

TheresultsareshowninTables1and2. Table1providestheoverallgradeforeachsystemorauthor averagedacrossallgradersandsummaries,whereeachassignedgradehasfirstbeenconvertedtoa number(A=4.0,B=3.0,C=2.0,D=1.0,F=0.0)andtheaverageconvertedbacktoalettergrade. Table2

---

[3]TheExemplarstextbuilderemploysaprocessingmodelthatissimilartoXSLT;theprimarydifferencebetween thetwoisthattheExemplarstextbuildingrulesaremoreobject -orientedthanXSLTtemplates,enablinggreater rulesophist icationandreuseacrossvaryingcontexts(cf.WhiteandCaldwell,1998).

shows the mean and standard deviations of the overall, content, organization, and readability scores for the RIPTIDES and the Baseline system averaged across all graders and summaries.

Given the amount of development effort that went into the initial version of the system, we were not surprised that our summarizer fared poorly when compared against manually written summaries, receiving an average grade of C, vs. A- and B+ for the human authors; nevertheless, the initial RIPTIDES system scored almost a full grade ahead of the baseline summarizer, which received a D/D+. The difference in the overall scores was significant, as were the scores for organization and readability (though not content). The most notable improvement was in organization, which was not surprising given that the Baseline system just listed extracted sentences in document order.

The comments of the evaluators helped to identify the most important problems to focus on in ongoing work. These problems include the need for better event description merging, more refined comparison of differences in reported numbers, improved rhetorical structuring of relief information, temporal expression normalization, and sentence reduction. With progress in these areas, we hope to achieve scores within one grade of human performance.

## 4 Related and Ongoing Work

The RIPTIDES system is most similar to the SUMMONS system of Radev and McKeown (1998), which summarized the results of MUC-4 IE systems in the terrorism domain. In comparison to SUMMONS, the RIPTIDES system appears to be designed to more completely summarize larger input document sets; in particular, we believe our system will scale to handle the hundreds of news articles we have collected about the recent earthquakes in Central America and India, whereas SUMMONS was more of an exploratory prototype that was never run on more than a handful of documents. Another important difference is that SUMMONS sidestepped the problem of comparing reported numbers of varying specificity (e.g. "several thousand" vs. "anywhere from 2000 to 5000" vs. "up to 4000" vs. "5000"), whereas we have recently implemented more robust rules for doing so. In our approach, a range encompassing the current reports across available news sources is constructed, and any lower, less specific or incomparable estimates (e.g. "more than half the region's residents") are noted (space permitting).[4]

In its treatment of relief information, the RIPTIDES system is also similar to, though simpler than, the domain-independent multi-document summarizers of Goldstein et al. (2000) and Radev et al. (2000) in the way it clusters sentences across documents to help determine which sentences are central to the collection, as well as to reduce redundancy amongst sentences included in the summary. It is also similar in spirit to MultiGen (Barzilay et al., 2001), though much less ambitious in its approach.

In ongoing work, we are in the process of refining our algorithm for summarizing differences in reported numbers and improving our treatment of relief information. At the conference, we plan on showing output from the current version of the summarizer, using the actual results of the IE system. For the final version of the paper, we plan on repeating our evaluation with the improved system, and will include the updated results in the final version of the paper.

---

[4] Less specific estimates such as "hundreds" are considered lower than more specific numbers such as "5000" when they are lower by more than a factor of 10.

# References

TheApacheXMLProject.2001."XalanJava." http://xml.apache.org/.

Barzilay,R.,N.ElhadadandK.McKeown.2001."Sen tenceOrderinginMultidocument Summarization."ToappearintheProceedingsofHLT2001.

Cardie,C.1997." EmpiricalMethodsinInformationExtraction." *AIMagazine* 18(4):65 -79.

Charniak,E.1999."Amaximum -entropy-inspiredparser."BrownUnivers ityTechnicalReportCS99 -12.

Goldstein,J.,Mittal,V.,Carbonell,J.,andKantrowitz,M.2000."Multi -documentsummarizationby sentenceextraction."In *ProceedingsoftheANLP/NAACLWorkshoponAutomaticSummarization,* Seattle,WA.

Radev,D.R.and McKeown,K.R.1998."Generatingnaturallanguagesummariesfrommultipleon -line sources." *ComputationalLinguistics* 24(3):469 -500.

Radev,D.R.,Jing,H.,andBudzikowska,M.2000."Summarizationofmultipledocuments:clustering, sentenceextracti on,andevaluation."In *ProceedingsoftheANLP/NAACLWorkshopon Summarization,*Seattle,WA.

Riloff,E.andJones,R.1999."LearningDictionariesforInformationExtractionbyMulti -Level Bootstrapping."In *ProceedingsoftheSixteenthNationalCon ferenceonArtificialIntelligence* , Orlando,FL,pp.474 -479.

Thompson,C.A.,Califf,M.A.,andMooney,R.J.1999."Activelearningfornaturallanguageparsing andinformationextraction."In *ProceedingsoftheSixteenthInternationalConference onMachine Learning*,Bled,Slovenia,pp.406 -414.

White,M.andCaldwell,T.1998."EXEMPLARS:APractical,ExtensibleFrameworkforDynamic TextGeneration."In *ProceedingsoftheNinthInternationalWorkshoponNaturalLanguage Generation*,Niagara -on-the-Lake,Canada,pp.266 -275.

White,M.2001."TextPolishing:Surface -OrientedSmoothingofGeneratedTextsviaMarkup -Based RevisionRules."Inpreparation.