

NATSUM: Narrative abstractive summarization through cross-document timeline generation

Cristina Barros, Elena Lloret, Estela Saquete, Borja Navarro-Colorado

*Department of Software and Computing Systems, University of Alicante
Apdo. de Correos 99 E-03080, Alicante, Spain
{cbarros,elloret,stela,borja}@dlsi.ua.es*

Abstract

A new approach to narrative abstractive summarization (NATSUM) is presented in this paper. NATSUM is centered on generating a narrative chronologically ordered summary about a target entity from several news documents related to the same topic. To achieve this, first, our system creates a cross-document timeline where a time point contains all the event mentions that refer to the same event. This timeline is enriched with all the arguments of the events that are extracted from different documents. Secondly, using natural language generation techniques, one sentence for each event is produced using the arguments involved in the event. Specifically, a hybrid surface realization approach is used, based on over-generation and ranking techniques. The evaluation demonstrates that NATSUM performed better than extractive summarization approaches and competitive abstractive baselines, improving the F1-measure at least by 50%, when a real scenario is simulated.

Keywords: Narrative summarization, Abstractive summarization, Timeline Generation, Temporal Information Processing, Natural Language Generation

1. Introduction

Managing and processing the over-abundance of information and its heterogeneity is an enormous challenge for human beings in the digital era. Therefore, the application of Human Language Technologies (HLT) is necessary to facilitate access to and use of this information. For example, every day, online newspapers generate countless digital texts (news) about the

7 same facts. In this context, a summary is useful to support humans in the
8 analysis and processing of information [1]. Text summarization can provide
9 appropriate mechanisms to automatically condense the key information that
10 is spread over different documents (e.g. news) [2].

11 To provide users with easy and optimal access to all this information,
12 summaries must provide a coherent and natural structure. In this sense,
13 narrative structure is the most natural and friendly text structure for human
14 beings [3]. As human beings, we tend to organize the flux of happening in
15 narrative structures, where a narrative structure is the arrangement of a set
16 of events about one or more entities following a time order (that could be
17 natural chronological order —from past to future— or artificial order —with
18 time jumps—). Each event is a fact that occurs in the (real or imaginary)
19 world at a specific moment with a specific structure (the event structure) [4],
20 and denotes processes, activities, states, achievements or accomplishments
21 [5]. Furthermore, an event involves participants [6] and other components
22 that complete the event such as time, place, instruments, patients, etc¹.

23 Depending on how a summary is produced, a distinction can be made
24 between *extractive* and *abstractive* summaries. *Extractive* summaries are
25 produced by directly selecting the most significant sentences of a document
26 and copying them verbatim into the output. *Abstractive* summaries are more
27 challenging, since they include new or different vocabulary, linguistic expres-
28 sions or concepts that do not originally appear in the input documents, but
29 that paraphrase the most relevant information of the input. When the sum-
30 mary is intended to narrate or describe a series of events that happened at
31 a specific time, *extractive* summarization approaches will lose the tempo-
32 ral connections appearing in the text, that can lead to dangling references,

¹From a linguistic point of view, the participants and components of an event are called “arguments” and “modifiers”. An event mention is formed by an event head (normally a verb, but not always), a set of arguments and optional complements. The arguments are those elements of the event structure that complete the meaning of the verb (as, for example, the person that carries out the specific action expressed by the verb, the person or object that receives the action, the instrument used to perform the action, etc.). The modifiers are the remaining optional elements of the event structure (the place where the action occurs, the time, etc.). In this paper, the word “argument” is used as a linguistic term to refer to the elements of the event structure [4]. Given that there is no common typology of arguments in the linguistic literature, we follow the proposal of PropBank project [7] to nominate arguments with numbers from A0 (the argument closest to the verb) to A4 (the most external argument), and AM for the remaining modifiers.

33 and thus the resulting text may be ambiguous or difficult to understand.
34 For instance, an *extractive* summarization system could select the sentence
35 “*Terrorists provoked the blast*” from the text shown in Example 1 without
36 providing any additional information about other relevant information, such
37 as *when?* or *where?*. However, using an *abstractive* summarization approach,
38 the relevant information (e.g., *who? what?, when?, where?,...*) could be
39 fused together, leading to the generation of one or more new sentences. Fol-
40 lowing the same text fragment given as example (Example 1), the sentence
41 “*On Friday, terrorists exploded bombs in the U.S embassy in the Kenyan and*
42 *Tanzanian capitals.*” could be generated.

43 (1) Suspected bombs [exploded *event*] outside the U.S. embassies in the Kenyan and
44 Tanzanian capitals [Friday *time*]. Terrorists provoked the [blast *event*]

45 However, although *abstractive* summarization would be more appropriate
46 than *extractive* summarization, the detection and resolution of temporal in-
47 formation is of crucial importance to anchor the event to a precise date. This
48 avoids reader misunderstanding, (e.g. instead of “*On Friday*”, it would be
49 more appropriate for ordering purposes to reformulate the expression as “*On*
50 *the 7th of August 1998*”). In this way, the final summary would be clearer,
51 containing all the relevant information within a coherent and cohesive text,
52 thereby removing any possible ambiguity.

53 The main objective of this paper is to develop an *abstractive* summariza-
54 tion approach that generates narrative summaries based on a natural time
55 ordering of events from a set of documents (news in this case) that deal with
56 the same real events. Hereafter we will refer to it as the acronym NATSUM
57 (Narrative Abstractive Timeline Summarization). This system has two main
58 components: (i) a cross-document timeline generation module that extracts
59 events related to the same entity from several texts (cross-document) and
60 the time slot in which each event occurs, arranging them in a timeline; and
61 (ii) an *abstractive* summarization module that transforms these time-ordered
62 events into a single text with a time-based chronological narrative structure.

63 The task of extracting events involving a particular target entity among
64 different documents and ordering them chronologically is known as Cross-
65 document Timeline Extraction [8]. Timeline Extraction comprises the ac-
66 complishment of three stages. The first step involves determining whether
67 the events extracted from the different documents are related to the target
68 topic or entity. From this first cluster of events, a *temporal information pro-*
69 *cessing* is required in order to extract the temporal expressions and the tem-

70 poral relationships established between these events, determining thus which
71 events happened at the same time. Finally, *cross-document event coreference*
72 is needed in order to cluster all the mentions that occur at the same time
73 and actually refer to the same event, regardless of the words used to express
74 them. The previous Example 1 contains two event mentions² that refer to
75 the same event.

76 For the creation of the narrative abstractive summary, a single sentence
77 for all the events mentions referring to the same event is generated. This
78 sentence includes all the information related to this event as well as the time
79 it occurred. In this way, the abstractive summaries will be generated over the
80 structured knowledge previously obtained from an enriched timeline³. This
81 implies an advance on classical timeline extraction as it involves the addition
82 of all the arguments related to the event. Also, there is an improvement in
83 automatic narrative summarization as the temporal information (temporal
84 expressions, events and temporal relationships) is considered in the summary
85 generation process.

86 The paper is organized as follows. Section 2 contains a detailed back-
87 ground study of the different relevant research fields, involving Automatic
88 Timeline Generation, Abstractive Summarization and Natural Language Gen-
89 eration. Section 3 describes the architecture of our proposed system NAT-
90 SUM. Following this, Section 4 presents the main experiments conducted
91 together with the evaluation methodology. Section 5 reports on the results
92 obtained and a discussion of the findings. Furthermore, Section 6 reports ad-
93 ditional experiments and evaluation to assess NATSUM’s performance within
94 the similar task of timeline summarization and compare its results to the state
95 of the art. Finally, Section 7 highlights the main conclusions of this research
96 and outlines some potential areas of future work.

97 **2. Background**

98 Considering that our proposal is generating narrative abstractive sum-
99 maries based on timeline knowledge, both research issues are tackled in this
100 section.

²Event mention is a reference to an event, that is, the different forms to refer to the same event.

³We propose summarization focused on a target entity because we are using the timelines defined in Semeval2015 Task 4, which defined timelines related to a target entity.

101 *2.1. Automatic Timelines*

102 Recently, SemEval-2015 [9] included a task that tried to combine temporal
103 information processing and event coreference to obtain a timeline of events
104 related to a specific given entity, from a set of documents [8]. They proposed
105 two different tracks on the basis of the data used as input. Track A, for
106 which they provided only raw text sources, and Track B, for which they also
107 made gold event mentions available.

108 Track A had two participants: WHUNLP team, that processed the texts
109 with Stanford CoreNLP⁴ [10] and applied a rule-based approach to extract
110 target entities and their predicates and also performs temporal reasoning⁵
111 and the SPINOZAVU [11] system, that is based on a pipeline, developed
112 in the NewsReader project, and addressed entity resolution, event detection,
113 event-participant linking, coreference resolution, factuality profiling and tem-
114 poral relation processing, first at document level, and then at cross-document
115 level, in order to obtain timelines.

116 Track B had also two participants: Heildeltoul team approach [12] that
117 uses the HeideTime tool for temporal information processing, and the Stan-
118 dord CoreNLP for event coreference resolution. A cosine similarity matching
119 function and a distance measure are used to select which sentences and events
120 are relevant for the target entity. Finally, GPLSIUA team [13], that uses
121 the OPENER language analysis toolchain⁶ for entity detection, the TIPSem
122 tool [14, 15] for temporal processing and a topic modeling algorithm over
123 WikiNews corpus to detect event coreference.

124 Outside SemEval-2015 competition, the work presented by Laparra et
125 al. 2017 [16] developed three deterministic algorithms for timeline extrac-
126 tion based on two main ideas: a) addressing implicit temporal relations at
127 document level, and b) leveraging several multilingual resources to obtain
128 a single, interoperable, semantic representation of events across documents
129 and across languages.

130 The novelty of our proposal is going further with the timeline extrac-
131 tion task, including all the participants in the events, and combining this
132 technique with a summarization approach to generate narrative and ordered
133 texts related to a specific topic.

⁴<http://stanfordnlp.github.io/CoreNLP/>

⁵No bibliography is available apart from the general paper of SemEval 2015 Task 4

⁶<http://www.opener-project.eu/webservices>

134 *2.2. Abstractive Summarization and Natural Language Generation*

135 As it was stated in the previous section, abstractive summarization is
136 far more challenging than extractive summarization, since it requires under-
137 standing the information expressed in one or several documents and com-
138 press, fuse, integrate, enrich or generalize it to create a new text (i.e., sum-
139 mary) that contains the key aspects of the input documents. For generat-
140 ing high quality abstractive summaries, the integration of Natural Language
141 Generation (NLG) techniques are crucial to be able to paraphrase the infor-
142 mation expressed in the original sentences.

143 NLG tasks are commonly viewed as a pipeline of three broad stages: doc-
144 ument planning (also known as macroplanning), microplanning and surface
145 realization [17]. In the document planning stage, the system must decide
146 what information should be included in the text and how to organize it into
147 a coherent structure, leading to a document/text plan. From this document
148 plan, in the microplanning stage, a discourse plan will be generated, where
149 appropriate words and references will be brought together into sentences.
150 Finally, the surface realization stage generates the final text with the infor-
151 mation and structure selected. Each of the stages described has different
152 goals and tasks to complete. In some research they are dealt with one at
153 a time, or they focus on one task in particular. As examples of the latter,
154 some popular tools developed in the context of NLG include SimpleNLG [18],
155 which prioritizes the realization stage, or more specialized tools such as AI-
156 GRE [19], whose focus lies on the referring expression generation task. There
157 have been some attempts to address the whole process as well, mostly using
158 machine learning techniques. For instance, Duma and Klein [20] proposed
159 that automatic template acquisition, and learning the content selection, out-
160 put structure and the lexical choices to display take place simultaneously
161 in a single process. Konstas and Lapata [21] analyzed several mechanisms
162 for mapping database information (weather forecast records) into natural
163 language sentences. These included the use of probabilistic grammars, the
164 detection of patterns in input records and the learning of rhetorical relations
165 to provide document plans from these records.

166 As regards the techniques used for automatic language generation, since
167 this is not a trivial task, NLG systems have used either statistical or knowl-
168 edge-based approaches. The underlying idea of statistical approaches is based
169 on the probability of certain words appearing together and/or in proximity,
170 studying the creation of a sentence on the basis of a set of words [22, 23].
171 In contrast, knowledge-based approaches use linguistic theories, e.g., rhetor-

172 ical structure theory, to generate the text [24]. The fundamental difference
173 between these approaches is the type of data used. Knowledge-based ap-
174 proaches use linguistic information (morphological, lexical, syntactic, seman-
175 tic), together with rules and pre-defined templates. Statistical approaches use
176 probabilistic information extracted from a text corpus. It is also important
177 to note that rule-based knowledge approaches are oriented to a specific do-
178 main and language. Consequently, their adaptation to a different domain or
179 language is extremely difficult and costly. In this sense, statistical approaches
180 offer an advantage, since they are more versatile for application across dif-
181 ferent domains or languages, as long as the probabilities are learned from
182 the appropriate corpora. Languages models (LM) can be considered one
183 of the most-used mechanisms from the statistical perspective in HLT [25].
184 To obtain knowledge from a corpus on frequency and probability of word
185 appearance — the fundamental idea behind LMs — several techniques can
186 be applied: maximum likelihood [26] and support vector machines [27] have
187 been widely used, for example.

188 In contrast to the NLG techniques for tackling abstractive summarization,
189 other techniques employing neural networks models have emerged in recent
190 years. For instance, See et al. [28] present a hybrid pointer-generator archi-
191 tecture with coverage for multi-sentence abstractive summarization. Chen
192 and Bansal [29] propose a fast summarization model that generates a concise
193 overall summary by selecting and rewriting salient sentences abstractively.
194 These types of models tend to contain redundant and/or repeated informa-
195 tion in the summary. In addition to these techniques, there are others that,
196 in some way are a middle-ground between abstractive and extractive tech-
197 niques. Examples of these types of techniques can be found in Cordeiro et
198 al. [30] where a methodology for learning sentence reduction is presented;
199 or in Valizadeh and Brazdil [31], where a summary is generated by selecting
200 the sentences which satisfy actor-object relationships.

201 Our summarization approach is completely abstractive, focusing only on
202 the surface realization stage, since the cross-document timeline generation
203 will be used as a document plan. Moreover, different from the state of art, to
204 generate a sentence, our approach will combine a statistical model together
205 with semantic information, thus resulting in an hybrid surface realization
206 method.

207 *2.3. Narrative structures extraction*

208 To the best of our knowledge, we are not aware of any previous work that
209 attempts to generate narrative abstractive summaries using timeline infor-
210 mation and NLG techniques. However, some previous proposals exist that
211 attempt to extract event-based narrative structures from texts. Chambers
212 and Juravsky [32, 33] extract narrative chains that define a partially ordered
213 sets of events that share a common actor (an entity person). The relation-
214 ship between events is, in this case, time relations. Our approach is based
215 on these narrative chains. Similar approaches are used by [34], [35] or [36]
216 to create narrative chains, but their work is focused on the extraction of
217 common sense knowledge for a complete understanding of narrative texts.
218 All these proposals extract the narrative chains from only one text. Our
219 approach is, however, cross-document. We extract a single timeline of events
220 (as a narrative chain) from several texts that talk about the same entity and
221 about the same events.

222 Regarding timelines, a task close to our proposal is timeline summariza-
223 tion. According to [37], given a query (such as “BP oil spil”), timeline sum-
224 marization needs to (i) extract the most important events for the query and
225 their corresponding dates, and (ii) obtain concise daily summaries for each
226 selected date ([38] [39] [40] [41] [42] [43] [44]). Formally, a timeline is a se-
227 quence $(d_1, s_1), \dots, (d_k, s_k)$ where the d_i are dates and the s_i are summaries
228 for the dates d_i , given a query q and an associated corpus C_q that contains
229 documents relevant to the query. The task of timeline summarization is to
230 generate a timeline s_q based on the documents in C_q . The number of dates
231 in the generated timeline, as well as the length of the daily summaries, are
232 typically controlled by the user. However, the aim of our proposal is to gen-
233 erate narrative summaries and not timelines, whereby timelines are used to
234 generate the narrative structure, which means that the input of the summa-
235 rization module is a target oriented timeline and not a set of documents, as
236 in TS approaches.

237 The next section presents how the summary generation is performed,
238 based on the arrangement of events along a timeline.

239 3. Narrative Abstractive Timeline Summarization System (NAT- 240 SUM): Design and Development

241 The task we address consists of producing an abstractive multi-document
242 summary that narrates the most relevant events⁷ together with the date they
243 occurred and when a specific target entity is involved. In this way, as shown
244 in Figure 1, given as an input a target entity and a set of documents related
245 to that target, the proposed system has to i) determine which events hap-
246 pened and when, choosing only the most relevant ones related to the target
247 entity, building a timeline, which is used to ii) generate the final abstractive
248 summary as output.

249 Therefore, the architecture of NATSUM comprises two different modules
250 and it uses a set of news documents and a target entity as input. The two
251 modules of the architecture are as follows:

- 252 • Enriched Timeline extraction: This module structures all the informa-
253 tion related to a specific topic/target entity in a timeline. All the event
254 mentions happening at the same time and referring to the same event
255 are grouped together on the timeline. This module is an improved
256 extension of the system presented in [45] .
- 257 • Abstractive summarization: This module is responsible for generating
258 a chronological abstractive summary based on NLG techniques given
259 an enriched timeline as input. Specifically, it employs a hybrid surface
260 realization approach, based on over-generation and ranking techniques.

261 The integration of both modules as a pipeline results in the generation of
262 a narrative abstractive summary. The proposed architecture is graphically
263 depicted in Figure 2. In the following sections, the development of each of
264 the aforementioned modules is explained in more detail.

265 3.1. Enriched Timeline extraction

266 As previously explained, given a set of documents and a set of target
267 entities, the original task of Cross-Document Timeline Extraction consists of
268 building an event timeline for a target entity from a set of documents [46].

⁷According to TimeML temporal annotation schema “events” is something that hap-
pens or occurs. Events can be punctual or last for a period of time. They also consider as
events those predicates describing states or circumstances in which something obtains or
holds true.

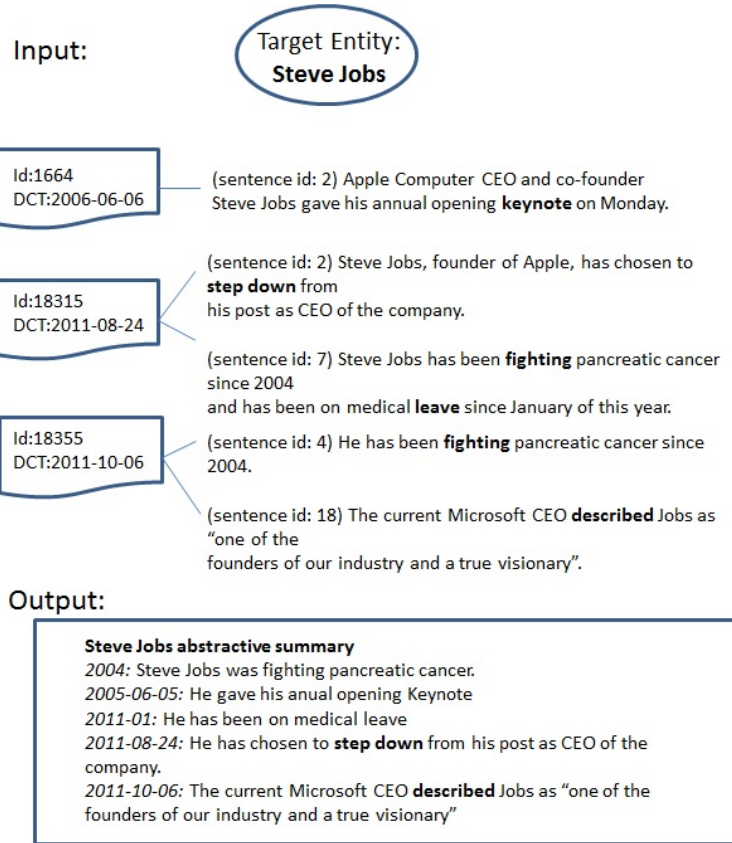


Figure 1: Example of input/output of the proposed system (NATSUM)

269 Theoretically, the main idea of our approach is that two events $e1$ and
 270 $e2$ will be coreferent if they are not only temporal compatible ($e1_t = e2_t$)⁸
 271 but also if they refer to the same facts (semantic compatibility: $e1_s \simeq e2_s$)⁹:

$$\text{coref}(e1, e2) \rightarrow (e1_t = e2_t) \wedge (e1_s \simeq e2_s) \quad (1)$$

272 Our proposal extends the approach by enriching the event clusters with

⁸ ei_t : Temporal information of the event i

⁹ ei_s : Semantic information of the event i

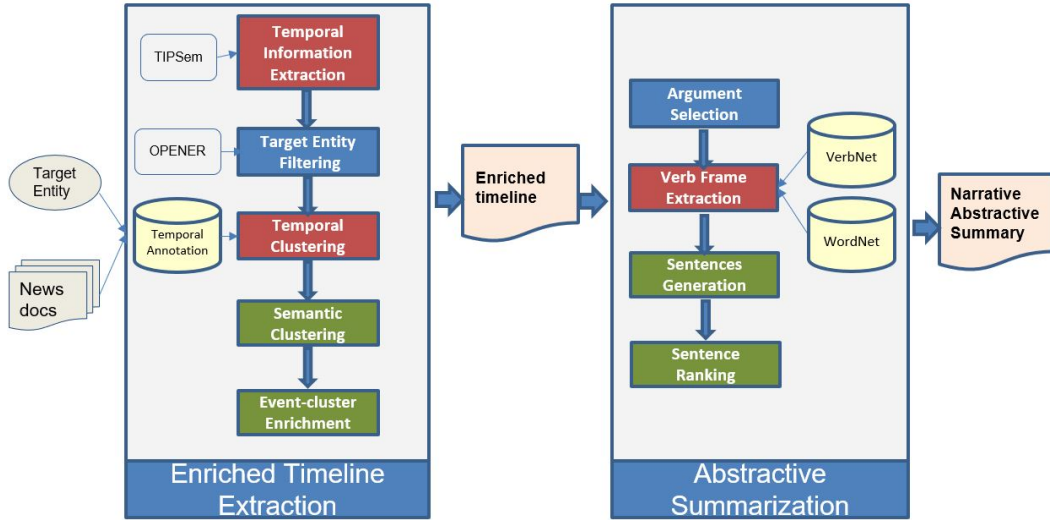


Figure 2: Architecture for our Narrative Abstractive Timeline Summarization system (NATSUM).

273 all the arguments extracted from these events in the different documents
 274 where they are presented. The steps of this module are:

- 275 • Temporal clustering¹⁰: by using the temporal information annotated
 276 by a temporal information processing system, the temporal relations
 277 between the events are processed and the events can be ordered and
 278 anchored to the timeline.
- 279 • Semantic clustering: the events are grouped together using event type
 280 information and distributional semantic knowledge.
- 281 • Event cluster enrichment: for each cluster of events, all the arguments
 282 related to the events in the cluster are added to the cluster.

283 3.1.1. Temporal Information Extraction

284 The input is a set of plain texts, and, therefore, the events in those
 285 texts must be automatically extracted. Furthermore, considering that the

¹⁰Temporal clustering in this context refers to Temporal Compatible Grouping, meaning that all the events happening at the same time are grouped together in a cluster. It is not the same concept as clustering in Machine Learning

286 final aim is building a timeline, temporal expressions and temporal links
287 between events and times are required. Therefore, plain texts need to be
288 annotated with all the temporal information. Several efforts have been made
289 to define standard ways to represent temporal information in texts. The
290 main objective of this representation is to make temporal information explicit
291 through standard annotation schemes. TimeML[47] is the most standardized
292 schema and it annotates not only events and temporal expressions, but also
293 temporal relations, known as links [48]. In this annotation schema, event is
294 used as a cover term to identify *something that can be said to obtain or hold*
295 *true, to happen or to occur*. This notion can also be referred to as eventuality
296 including all types of actions (punctuals or duratives) and states as well
297 (section 1, NewsReader Guidelines¹¹). Besides, according the task definition
298 of Semeval 2015 —task 4, not all events can be part of a TimeLine, amongst
299 others, counter-factual events will not appear in a TimeLine. Example (2)
300 shows a sentence annotated with TimeML temporal expressions (TIMEX3),
301 events (EVENT), and the links between them (TLINK).

```
302 (2) John <EVENT eid="e1">came</EVENT> on <TIMEX3 tid="t1">Monday</TIMEX3>  
303 <TLINK eventInstanceID="e1" relatedToTime="t1" relType="IS_INCLUDED" />
```

304 In our case, the first step is performing Temporal Information Extraction
305 and Processing, and TIPSem system (Temporal Information Processing using
306 Semantics) [14, 15]¹² is used for this purpose. TIPSem is able to automati-
307 cally annotate all the temporal information according to TimeML standard
308 annotation scheme [47], which means annotating all the temporal expressions
309 (TIMEX3), events (EVENT) and links (TLINKS) between them.

310 3.1.2. Target Entity Filtering

311 Considering that not all the events are necessary to build the timeline, but
312 only the ones related to a target entity, a Target Entity Filtering needs to be
313 performed in order to discard those events that are annotated but not related
314 to the given entity. The Target Entity Filtering requires resolving name entity
315 recognition and entity coreference resolution, and OPENER¹³ web services
316 are used for this purpose. To determine whether an event should be part of
317 the timeline, this module chooses: a) the events in which a target entity (or a

¹¹<http://www.newsreader-project.eu/files/2013/01/NWR-2014-2.pdf>

¹²<http://gplsi.dlsi.ua.es/demos/TIMEE/>

¹³<http://www.opener-project.eu/webservices>

318 target entity coreference) explicitly participates in a *has_participant* relation
319 with the semantic role A0 (i.e. agent) or A1 (i.e. patient), as defined in the
320 Propbank Project [7], and b) in case of nominal events, since the information
321 of A0 or A1 is not obtained, this module chooses this type of event if the
322 target entity is contained in the sentence. For example, for the target entity
323 “*Steve Jobs*” and the nominal event “*keynote*”, this event should be chosen
324 due to the sentence in which appears: “*Steve Jobs gave his annual opening*
325 *keynote on Monday*”.

326 Otherwise, the event is discarded.

327 3.1.3. Temporal Clustering

328 Considering the premise that two events referring to the same event hap-
329 pen at the same time, and using the temporal annotation of the input texts
330 (TimeML annotation schema¹⁴), the temporal clustering algorithm performs
331 two steps:

- 332 • *Within-document temporal clustering*: For each document, the tem-
333 poral information of each event is extracted. Each event is anchored
334 to a time anchor¹⁵ when a temporal SIMULTANEOUS/ BEGIN/ IN-
335 CLUDES link exists between this event and a temporal expression.
336 After this, two events are grouped together if they are temporally com-
337 patible. This means that: a) two events are anchored to the same
338 time anchor, or b) two events have a temporal SIMULTANEOUS link
339 between them.

340 Example 3 shows two events temporally compatible and grouped to-
341 gether.

342 (3) a. The <EVENT eid="e1"> meeting </EVENT> was
343 <TIMEX3 tid="t1" value="2014-03-22"> yesterday </TIMEX3>.

¹⁴<http://www.timeml.org/>

¹⁵A time anchor is always a DATE (as defined in TimeML standard annotation) and its format follows the ISO-8601 standard: YYYY-MM-DD. The finest granularity admitted in the task for a time anchor is DAY. Other granularities admitted are MONTH (references as YYYY-MM) and YEAR (references as YYYY). A time anchor takes as value the point in time when the event occurred (in case of punctual events) or began (in case of durative events). Event ordering is based on temporal relations between events; more specifically on the before/after and includes/simultaneous relations as defined by ISO-TimeML. The system places the dates in the timeline from lowest to finest granularity.

344 b. At the same time, the teacher <EVENT eid="e2"> presents </EVENT>
 345 the ideas.
 346 <TLINK eventInstanceID="e1" relatedToTime="t1"
 347 relType="IS_INCLUDED"/>
 348 <TLINK eventInstanceID="e2" relatedToEventInstance="e1"
 349 relType="SIMULTANEOUS"/>

350 Two non-temporally compatible events are shown in Example 4.

351 (4) a. The <EVENT eid="e1"> meeting </EVENT> was
 352 <TIMEX3 tid="t1" value="2014-03-22T17:00"> yesterday at 17:00
 353 </TIMEX3>.
 354 b. After that, the teacher <EVENT eid="e2"> presents </EVENT> the ideas.
 355 <TLINK eventInstanceID="e1" relatedToTime="t1"
 356 relType="IS_INCLUDED"/>
 357 <TLINK eventInstanceID="e2" relatedToEventInstance="e1"
 358 relType="AFTER"/>

359 • *Cross-document temporal clustering*: Considering that in the previous
 360 step all the events of each document were assigned to a time anchor, in
 361 this step, this information is merged in a single timeline, in which all
 362 the events of the different documents are grouped together if they are
 363 happening at the same time.

364 (5) a. Document 1: The <EVENT eid="e1"> meeting </EVENT> was <TIMEX3
 365 tid="t1" value="2014-03-22"> yesterday </TIMEX3>.
 366 <TLINK eventInstanceID="e1" relatedToTime="t1"
 367 relType="IS_INCLUDED" />
 368 b. Document 2: The students <EVENT eid="e5"> met </EVENT> on <TIMEX3
 369 tid="t3" value="2014-03-22">Tuesday</TIMEX3>.
 370 <TLINK eventInstanceID="e5" relatedToTime="t3"
 371 relType="IS_INCLUDED" />

372 According to Example 5 and after performing the within-document
 373 temporal clustering, doc1-e1 is anchored to the date “2014-03-22”, and
 374 doc2-e5 is anchored to the same date. Therefore, in the cross-document
 375 temporal clustering step these two events will be considered part of the
 376 same group.

377 Finally, the temporal groups are chronologically ordered. For each line,
378 there is first a cardinal number indicating the position of an event in the
379 timeline, then the value of the anchor time attribute, and finally the list of
380 events anchored to this time attribute. Each event is represented as follows:
381 language (en/es), document identifier, sentence number and textual extent
382 of the event. For example, the event en-18315-7-leave is located in sentence
383 7 of document 18315 and it is in English. In this first clustering, if two
384 events have the same value for the anchor time attribute, they are placed in
385 the same group. In the next step, explained in the following section, these
386 temporal groups will be divided again according to their semantics.

387 3.1.4. *Semantic Clustering*

388 Two or more event mentions in the same time slot could refer to the same
389 real event. To detect these coreferential events, we have applied a clustering
390 process based on two kinds of semantic information: i) the event type; and,
391 ii) distributional semantic similarity between event mentions.

392 During the event extraction process, each event mention has been clas-
393 sified according to its type of event following TimeML standard [49]: oc-
394 currence, perception, reporting, aspectual, state, intentional state and in-
395 tentional action. All the event mentions with the same time slot have been
396 regrouped after also considering the type of event to which they have been
397 assigned.

398 Next, our approach clusters coreferential events (identifies all the events
399 that share the same time slot and the same type of event) according to the
400 compositional-distributional semantic similarity between them. The seman-
401 tics of the event structure is represented as a compositional-distributional
402 vector. Rather than creating a complex feature matrix to represent the se-
403 mantics of the argument, as described in [50], we propose a compact dis-
404 tributional semantic model. In this way, we consider the context of the
405 events as the main component that contributes to establishing the semantic
406 compatibility and, therefore, the event coreference. This relies on the fact
407 that distributional semantics are based on the contextual meaning of words
408 [51, 52]. Beyond trying to represent the meaning of words through lexicons
409 or ontologies, distributional semantics represent how words are used in real
410 context through vector spaces [53, 54]. These vectors are called contextual
411 vectors. Specifically, for each word of the event structure we have used the
412 English Word2Vec word embedding trained on the Google News corpus.

413 In our approach each event structure is formed, on the one hand by the

414 event head and, on the other hand by the nouns, verbs and adjectives of
 415 the main arguments. All this information is extracted by applying Freeling
 416 [55] as Part of Speech tagger and Semantic Role Labeling system. Following
 417 the additive model [56], these word vectors are added in a single composi-
 418 tional vector that represents the distributional meaning of the whole event
 419 structure.

420 An event structure (ES) with two arguments is formally represented as
 421 a tuple of three elements: two arguments ($A0$ and $A1$) and one event head
 422 (H):

$$ES = \langle A0, A1, H \rangle \quad (2)$$

423 Each argument is a compositional vector $\vec{V}(A)$ formed by the sum of the
 424 contextual vector $\vec{V}(w_n)$ of each word of the argument:

$$\vec{V}(A) = \sum^n \vec{V}(w_n) \quad (3)$$

425 where w_n represents each word of an argument and $\vec{V}(w_n)$ the contextual
 426 vector of each one of these words.

427 The event head H is the contextual vector of a single word. Finally, the
 428 compositional vector of the whole event structure $\vec{V}(ES)$ is:

$$\vec{V}(ES) = \vec{V}(A0) + \vec{V}(A1) + \vec{V}(H) \quad (4)$$

429 where $+$ means sum of vectors.

430 The similarity among all vectors two-to-two is represented by a square
 431 matrix. The final cluster is obtained applying a standard hierarchical cluster
 432 to this matrix. Specifically, we have applied an agglomerative clustering
 433 based on the average linkage criteria that uses the arithmetic mean of the
 434 distances between clusters to construct the dendrogram. We consider all
 435 event mentions grouped together at level one of this hierarchical cluster, that
 436 is, the second-most coarse-grained level under the root of the dendrogram.

437 3.1.5. Event cluster enrichment

438 The timeline consists of structured information in which all the event
 439 mentions related to the same event are grouped together according to the
 440 exact date when the event occurs. However, this information is not useful if
 441 the user that needs the information only has the event core (verb or nomi-
 442 nalization). The user will also need the arguments involved in the event to

443 obtain the accurate information about the event. Therefore, in this step, all
444 the arguments (semantic roles extracted in the previous step with Freeling)
445 of the events in each cluster are added to the timeline, enriching the infor-
446 mation provided for each event. In the Example 6, an enriched cluster of the
447 event mentions related to the same event is presented.

448 (6) 0 2008 en-82548-4-built:(A1,*The plane*),(A2,*with four Rolls–Royce_Trent 900 engines*)
449 (EN: In 2008, they built the plane with four Rolls-Royce_Trent 900
450 engines)
451 en-82548-2-made:(A1,*The first A380 superjumbo*),(A0,*by Airbus*)
452 (EN: In 2008, Airbus made the first A380 superjumbo)

453 In the example, for each event mention, all the arguments found in the
454 input document are added to the event mention with their corresponding
455 semantic role (A0, A1,...). Therefore, not only the event mention is used but
456 also the argument information.

457 3.2. *Abstractive summarization*

458 As previously mentioned, the aim of this module is to produce a narrative
459 abstractive summary with information given in an enriched timeline. This
460 summary is generated employing NLG techniques. In particular, we employ
461 a hybrid surface realization approach, based on over-generation and ranking
462 techniques. In these types of techniques, several possible outputs are gener-
463 ated and then ranked in order to select the best one, based on probability
464 models. For each of the enriched cluster of events from the enriched timeline,
465 the next steps are as follows:

- 466 • Argument selection: the arguments from the enriched timeline are se-
467 lected in the case that there is more than one argument for the same
468 semantic role. This selection is performed based on the probability of
469 the phrases contained in the arguments, which is calculated using a
470 language model.
- 471 • Obtaining verb frames: information about the frames corresponding to
472 the verbs of each event is obtained to generate a sentence without the
473 need to resort to grammar specifications.
- 474 • Sentence Generation: for each of the frames obtained a sentence is
475 generated, based on the frame structure.

- 476 • Sentence Ranking: a ranking is performed for selecting only one sen-
477 tence representing a specific event (cluster of event mentions) in the
478 timeline.

479 Before beginning the generation process, a language model is trained over
480 each of the input documents. This language model will be employed in some
481 of the steps of this module, and in particular, Factored Language Models
482 (FLM) are used to train it. FLM are an extension of the conventional lan-
483 guage models, proposed in [57], where a word is viewed as a vector of k factors
484 such that $w_t \equiv \{f_t^1, f_t^2, \dots, f_t^K\}$. The factors within this kind of model can
485 be anything, ranging from more basic elements, such as words or lemmas to
486 any other lexical, syntactic or semantic features needed for the task to be
487 addressed. The main objective of this type of model is to create a conditional
488 probability model over the selected factors: $P(f|f_1, \dots, f_N)$, being the pre-
489 diction of the factor f based on its N parents $\{f_1, \dots, f_N\}$. For the purpose of
490 this research, information about words, lemmas, Part-of-Speech (POS) tags
491 and synsets¹⁶ are used as the factors for training the FLMs. These factors
492 were selected due to the type of information they provide. In this regard,
493 syntactic and semantic information along with information about the words
494 themselves are needed in order to create a flexible abstractive summary in
495 relation to its vocabulary. To deal with these types of statistical models,
496 the SRILM [58] is used. This software is a toolkit for building and applying
497 statistical language model, which includes an implementation of FLM.

498 3.2.1. Argument selection

499 Taking as input the enriched timeline, for each of the events contained in
500 it, their arguments are checked to avoid duplicate semantic roles in the same
501 event.

502 In the case that two or more arguments for the same semantic role appear
503 within the event, the probability of the phrases contained in the arguments is
504 calculated employing the FLM previously trained. This probability is calcu-
505 lated employing only the words in the arguments either using the probability
506 given by the FLM when the phrase has 3 or less words, or otherwise, using
507 the chain rule (see Equation 5). In the chain rule, the probability of a phrase

¹⁶Set of cognitive synonyms related to a concept used in WordNet.

508 or a sentence is calculated as the product of the probability of all its words.

$$P(w_1, w_2 \dots w_n) = \prod_{i=1}^n P(w_i | w_1, w_2 \dots w_{i-1}) \quad (5)$$

509 When the probability of the different arguments for the same semantic
510 role is calculated, the argument with the highest probability is selected. In
511 Example 7 an event with several arguments for the same semantic role is
512 shown. In this example, the first argument for A1 (i.e. *Boeing*) will be
513 selected since its probability is higher than the one of the second argument
514 for A1 (i.e. *Civilian Deputy Undersecretary Darleen Druyun*).

515 (7) 0 2005 en-1173-35-hired:(A1,Boeing),(A1,CivilianDeputyUndersecretaryDarleenDruyun)
516 Probability of “Boeing”: 0.20
517 Probability of “Civilian Deputy Undersecretary Darleen Druyun”: 0.15

518 3.2.2. Verb frame extraction

519 After the different elements of the enriched timeline (i.e. their arguments)
520 are selected, the lexical resources VerbNet [59] and WordNet [60] are used to
521 obtain syntactic frames, from their event cores, which will be used during the
522 summary generation. VerbNet is one the largest verbs lexicons for English
523 including semantic and syntactic information about verbs. WordNet is a lex-
524 ical database composed by sets of synonym elements. Using both resources, a
525 set of frames containing the following information is extracted: i) the frames
526 from VerbNet comprise syntactic as well as semantic information about each
527 of the verbs of the lexicon; ii) WordNet provides a set of generic frames for all
528 the verbs. For every event, a set of frames from both, VerbNet and WordNet
529 are compiled. These frames are then analyzed to find out which elements of
530 the sentences need to be generated in the next step—the components of the
531 sentence, such as the subject or the object—. This avoids having to define a
532 grammar specification with the associated high cost.

533 When extracting the frames from VerbNet and WordNet, the “V“ in the
534 frames from Verbnnet represents the verb. WordNet, in this regard, is used
535 to extract the generic frames from a verb, which are consequently used to
536 produce a sentence for each of them.

537 Example 8 shows the frames which would be obtained from the event cores
538 of the Example 6 (i.e. *built* and *made*). Since the verbs *build* and *make*, for
539 the sense of constructing something combining materials and parts, belong to

540 the same VerbNet class and have the same synset in WordNet, the extracted
541 frames are the same for both.

542 (8) **VerbNet frames**

543 Agent V

544 Agent V Material

545 **WordNet frames**

546 Somebody - - -s something

547 3.2.3. Sentence generation

548 For each of the frames obtained in the previous step, a sentence is gen-
549 erated. If the specific event from which the verb frame was extracted has
550 arguments, the sentence is generated using these arguments along with the
551 information from the verb frame. The components of the frame may indicate
552 the need for some particular type of semantic role, such as an agent (i.e. A0,
553 A1) or an instrument (i.e. A2). Therefore, the sentence will be composed
554 using only the arguments needed and putting them in the order specified by
555 the frame. In certain cases, where the verb permits, if there is not an A0 but
556 an A1 argument, the A1 is treated as the Subject of the sentence, and this
557 sentence is generated in the passive voice.

558 In the case that the event does not have any arguments, a sentence is
559 generated following the structure given by the verb frame. For instance, if
560 the frame indicates the need for a Subject, it is generated based on the FLMs
561 trained, choosing the words with the highest probability appearing with the
562 corresponding verb of the event. The Object of the sentence is generated
563 using the same process, if needed.

564 In Example 9 the generated sentences for the frames shown in Example 8
565 can be seen. It is possible that, for the same verb, the frames obtained from
566 VerbNet and WordNet contain similar information to decide which arguments
567 of the event to select. In these cases, it is likely that the sentences generated
568 by both frames are the same, since they use the same arguments to generate
569 it.

570 (9) **build**

571 The plane was built.

572 The plane was built with four Rolls-Royce Trent 900 engines.

573 The plane was built with four Rolls-Royce Trent 900 engines.

574 **make**

575 by Airbus made.
576 by Airbus made the first A380 superjumbo, made by Airbus.
577 by Airbus made the first A380 superjumbo, made by Airbus.

578 3.2.4. Sentence ranking

579 Once a set of possible sentences containing the information of a specific
580 event is generated, a ranking is performed in order to select the sentence
581 which will form part of the chronological abstract summary. For selecting
582 the final sentences, the following process is applied: sentences are ranked
583 based on their probability which is computed by the chain rule (see Section
584 3.2.1).

585 The calculation of the probability of a word may differ depending of
586 the language model employed. Since, in this work, FLMs are used, the
587 probability of a word is calculated as the linear combination of FLMs as
588 suggested in [61] where a weight λ_i , is assigned to each of them (see Equation
589 6), being their total sum 1. In this Equation, f refers to a lemma, p refers to
590 a POS tag, and λ_i are set $\lambda_1 = 0.25$, $\lambda_2 = 0.25$ and $\lambda_3 = 0.5$. These values
591 were empirically determined by testing different values and comparing the
592 results obtained.

$$P(w_i) = \lambda_1 P(f_i | f_{i-2}, f_{i-1}) + \lambda_2 P(f_i | p_{i-2}, p_{i-1}) + \lambda_3 P(p_i | f_{i-2}, f_{i-1}) \quad (6)$$

593 The final selected sentence will be the one with the highest probability.
594 This sentence along with the date on which the event took place will be
595 considered as the sentence representing the information of the event.

596 Example 10 shows the final sentence selected from the ones in Example
597 9. The probabilities provided for each sentence are computed employing the
598 chain rule explained above (Equation (6)).

599 (10) Probability of “The plane was built.” : 0.16
600 Probability of “The plane was built with four Rolls-Royce Trent 900 engines.”: 0.25
601 Probability of “The plane was built with four Rolls-Royce Trent 900 engines.”: 0.25
602 Probability of “by Airbus made.”: 0.12
603 Probability of “by Airbus made The first A380 superjumbo, made by Airbus.”: 0.08
604 Probability of “by Airbus made The first A380 superjumbo, made by Airbus.”: 0.08
605 **Final Selected Sentence:** The plane was built with four Rolls-Royce Trent 900
606 engines.

607 Then, this sentence will be included in the final narrative abstractive
608 summary together with the remaining sentences generated by repeating this
609 process for each line in the enriched timeline.

610 4. Experimental Setup and Evaluation

611 NATSUM is focused on the transformation from a simple timeline to a
612 coherent narrative abstractive summary. For the evaluation of our system,
613 the test dataset provided for Task 4 at SemEval 2015 is used.¹⁷ This dataset
614 is composed of Wikinews articles about different topics: Airbus and Boeing;
615 General Motors, Chrysler and Ford; and the Stock Market. This evaluation
616 corpora consists of 90 documents (around 30,000 tokens and 915 events) and
617 they are very similar in terms of size. Each narrative abstractive summary
618 generated from the enriched timeline is entity-focused. This means that a
619 set of target entities is also provided within the corpus, and each timeline is
620 only composed of events related to this target entity. There is a total of 35
621 target entities in this dataset.

622 The following subsections provide information about the main experi-
623 ments carried out with the SemEval 2015 Task 4 dataset (Section 4.1), and
624 the evaluation methodology proposed (Section 4.2).

625 4.1. Main Experiments

626 Regarding the experiments conducted, for each target entity in the Se-
627 mEval 2015 Task 4 dataset, a narrative abstractive summary was generated
628 considering two configurations: (i) gold-standard experiment and (ii) over-
629 all system experiment. In total, 70 narrative summaries were generated (35
630 summaries for each experiment). For the gold-standard experiment, gold-
631 standard timelines provided in SemEval 2015 Task 4 are used. Using these
632 gold-standard timelines it is possible to measure the abstractive summariza-
633 tion module, avoiding the errors derived from the enriched timeline genera-
634 tion task. For the overall experiment, unannotated data is used to evaluate
635 the system in a real scenario in which our narrative abstractive summaries
636 could be applied. In this manner, the raw data of the Semeval corpus was
637 used as input, and then, the Enriched Timeline Extraction module provided
638 an intermediate scheme. The scheme contains the events and temporal in-
639 formation to be used by the Abstractive Summarization module to generate

¹⁷<http://alt.qcri.org/semeval2015/task4/index.php?id=data>

640 the sentences that will compose the final narrative summary. Furthermore,
641 the Timeline Extraction module was evaluated in isolation obtaining the
642 following results for English: F1-measure 27.63%, Precision 25.28%, Recall
643 30.47%. These results surpass the evaluation presented in [45], but evalu-
644 ating the Enriched Timeline Extraction module is beyond the scope of this
645 work. In addition, several state-of-the-art extractive summarization systems
646 were also used for the experiments for comparison purposes. In particular,
647 we selected the following systems: COMPENDIUM [62], GRAFENO [63]
648 and Open Text Summarizer (OTS) [64], since they provide either a visual
649 interface or the program to generate the summaries. In order to generate
650 multi-document and entity-focused extractive summaries that contain the
651 relevant information about a given entity, the input documents were prepro-
652 cessed following a two-step strategy. Firstly, all the documents belonging to
653 the same corpus were merged into a single macro-document; and secondly,
654 noisy sentences were removed from the input macro-document, i.e., the sen-
655 tences not talking about the focused entity or referring to them. By this
656 means, the job of summarization systems was only focused on determining
657 the relevant information to generate the final extractive summary, so the
658 techniques they implemented remained the same. In the end, 35 summaries
659 were produced by each system.

660 Finally, two baselines for narrative abstractive summarization were also
661 proposed (*FirstEvent* and *LongestEvent*). These baselines generate the nar-
662 rative summary using either only the first event (*FirstEvent*), or the event
663 with the highest number of arguments (*LongestEvent*) of each cluster pro-
664 vided by the gold-standard timelines—for experiment (i)—, or by the en-
665 riched timeline—for experiment (ii)—.

666 4.2. Evaluation Methodology

667 To assess the appropriateness of the resulting summary in terms of its
668 content and fluency, two types of quantitative evaluation were performed,
669 together with an additional human linguistic evaluation.

670 The first quantitative evaluation involved the analysis of extractive sum-
671 maries generated by state-of-the-art summarization systems. The goal of
672 this evaluation was to determine to what extent extractive summarization
673 systems were able to capture the relevant events and temporal information
674 contained in the input documents, and whether these systems were appropri-
675 ate for conducting narrative summarization or not. For this, we computed

676 the number of events and temporal information, comparing them to the gold-
677 standard annotations of the corpus employed. In order to avoid the errors
678 that may be obtained by just computing whether an event is present or not
679 in the summary we also took into account the location of the event, i.e., the
680 sentence in which it appears. For instance, the summary may contain a verb
681 but this does not necessarily refer to the same event of the gold-standard,
682 underscoring the importance of identifying the context in which the event
683 occurred so as to verify the accuracy of the generated summary.

684 The second type of quantitative evaluation is based on the hypothesis that
685 our abstractive summarization proposal enhances the quality of the narrative
686 summaries, relying on NLG techniques and using temporal information. For
687 this purpose, ROUGE tool [65] was used. ROUGE evaluates how informative
688 an automatic summary is by comparing its content to one or more reference
689 summaries. Such comparison is made in terms of n-gram co-occurrence (e.g.,
690 unigrams, bigrams, or word sequences). Moreover, ROUGE implements dif-
691 ferent metrics, such as unigram similarity (ROUGE-1); bigram similarity
692 (ROUGE-2); longest common subsequence (ROUGE-L) and bigram similar-
693 ity skipping unigrams (ROUGE-SU4). For each of these metrics, it provides
694 the commonly used HLT measures (precision, recall and F1-measure):

$$695 \textit{Precision} = \frac{\#CorrectPhrasesExtracted}{\#TotalPhrasesExtracted}, \quad (7)$$

$$\textit{Recall} = \frac{\#CorrectPhrasesExtracted}{\#CorrectPhrasesTest}, \quad (8)$$

$$696 \textit{F1 - measure} = \frac{2 * \textit{Precision} * \textit{Recall}}{\textit{Precision} + \textit{Recall}}, \quad (9)$$

697 where $\#CorrectPhrasesExtracted$ is the number of correct sentences that the
698 evaluated system extracts, $\#TotalPhrasesExtracted$ the total number of sen-
699 tences that the evaluated system extracts and $\#CorrectPhrasesTest$ the total
700 number of sentences included in the reference summaries.

701 ROUGE requires reference summaries and the creation of them is a time-
702 consuming and costly task. Therefore, a semi-automatic process was imple-
703 mented in order to generate a reference summary directly created from the
704 gold-standard timelines that were available within the corpus used for the
705 experiments. This process is further described in Section 4.2.1.

706 After having created the set of reference summaries, we directly compared
707 the content of the generated summaries to the reference ones. For this evalua-

708 tion, apart from our proposed narrative abstractive summarization approach
709 (NATSUM) , we also considered the extractive systems previously analyzed
710 (COMPENDIUM, GRAFENO and OTS), as well as the two proposed base-
711 lines (FirstEvent and LongestEvent). This enabled a comparison of this
712 paper’s proposal with other approaches, as well as verifying whether extrac-
713 tive summarization systems present limitations when it comes to performing
714 this task.

715 Using ROUGE for conducting this evaluation is appropriate as the events
716 are represented with words (generally verbs). Therefore, if the automatic
717 summary correctly captures the relevant events together with the right ar-
718 guments, the result for the ROUGE metrics will increase because the gen-
719 erated summary and the reference summary (gold-standard) are similar. In
720 this context, the summaries contain the key information of the documents.
721 However, using ROUGE exclusively for the evaluation is limited, since it is
722 not useful for determining the linguistic quality of the generated summaries
723 and is incapable of deciding the degree of grammatical correctness and mean-
724 ingfulness of the summaries. In this manner, a human evaluation was also
725 carried out involving several assessors that evaluated the linguistic quality of
726 the generated summaries. Hence, quantitative as well as qualitative results
727 were obtained (reported and explained in Section 5). The linguistic quality
728 of the generated abstractive summaries was assessed taking the readability
729 and linguistic criteria of the well-known summarization tracks for DUC¹⁸
730 and TAC¹⁹ conferences as a benchmark. Specifically, we evaluated the read-
731 ability/fluency of the summaries, including different criteria, such as the
732 summary’s grammaticality, non-redundancy, referential clarity, focus, as well
733 as structure and coherence. Moreover, the summary’s overall responsiveness
734 was also evaluated to determine the extent to which the amount of informa-
735 tion in the summary actually helped satisfy the information requirement.

736 For this, 12 humans with an advanced level of English participated in
737 this evaluation. The task consisted of completing a questionnaire²⁰ that
738 tackled the previously mentioned linguistic issues. Finally, also as part of
739 the manual evaluation, a human relevance judgement evaluation was carried
740 out. In this manner, we could check from a human perspective, which system

¹⁸<https://www-nlpir.nist.gov/projects/duc/index.html>

¹⁹<https://tac.nist.gov/>

²⁰<https://goo.gl/buC68B>

741 generated the summaries that were most preferred by users. To conduct this
742 task, assessors had to assign a preference ranking for a set of summaries,
743 indicating their most preferred, second most preferred and least preferred
744 summary. A second questionnaire was designed for this purpose²¹.

745 4.2.1. Generation of reference summaries

746 In this section, we explain the process for creating the reference sum-
747 maries that will be used in the quantitative evaluation. To create reference
748 summaries that allow us to evaluate the proposal, a set of patterns are applied
749 over the gold enriched timelines.

750 The following steps are performed in order to generate each sentence that
751 will compose the reference summary:

- 752 • *Verb selection*: Since the cluster contains different event mentions for
753 the same event, in the reference summary the first verb in the cluster
754 is used as representative of all the events in the cluster.
- 755 • *Arguments selection*: In order to create the sentence, only one of each
756 type of argument is necessary. In case there is more than one, the
757 longest one is chosen, since it is the most complete one, and it would
758 contain more information about the argument, thus leading to a more
759 informative sentence.
- 760 • *Sentence generation*: For each cluster, a sentence following this pattern
761 is generated:

762 (11) **Pattern:** *Time* A0 *event* A1 A2 A3 A4

763 Only the arguments available are used. A2, A3 and A4 are optional,
764 but in case there is no A0, or A1, the target entity is used.

765 In case of nominalizations, since they are not verbs, it is not possible to
766 obtain any semantic role. For these cases, we create a sentence using
767 the pattern:

768 (12) **Pattern:** *Time TargetEntity* had a *NominalizationEvent*

769 **Example:** On February (*Time*) Airbus (*TargetEntity*) had a
770 crush (*Nominalization*)

²¹<https://goo.gl/Mrj8yY>

771 **5. Results and Discussion**

772 In this section, we show the results obtained through the different eval-
 773 uations described in the previous section, as well as the analysis of these
 774 results.

775 *5.1. Limitations of Extractive Summarization*

776 Table 1 shows the results obtained after analyzing both the number of
 777 relevant events and the presence of temporal information that were contained
 778 in the extractive summaries generated by COMPENDIUM, GRAFENO and
 779 OTS. As observed, although the extractive summarization systems were
 780 adapted to be multi-document and entity-focused, they are only able to
 781 capture a small percentage of the relevant events and temporal information
 782 that should be included in the narrative summary. Concerning the number
 783 of events reflected in the summary, the highest result was obtained by the
 784 GRAFENO system (38.49%), but this result still represents less than half
 785 of the relevant events identified in the gold-standard. As for the temporal
 786 information, we noted that GRAFENO is the extractive system that obtains
 787 the poorest results, reflecting 7% of the temporal information, which may
 788 render difficult the comprehension of the summary with respect to the dates
 789 of the different events. COMPENDIUM and OTS, the other systems used,
 790 both exhibit similar performance.

791 Given that several relevant events were not captured and temporal infor-
 792 mation was omitted— hence, these items were not extracted as part of the
 793 output summary— we can conclude that traditional extractive summariza-
 794 tion systems are not effective in terms of generating narrative summaries.

System	Events	Temporal information
COMPENDIUM	26.86%	18.90%
GRAFENO	38.49%	7.10%
OTS	22.04%	18.04%

Table 1: Average percentage of events and temporal information reflected in extractive summaries.

795

796 *5.2. Summarization Results*

797 This section describes the automatic and manual evaluation for NAT-
798 SUM within the two experiments conducted: i) gold-standard experiment,
799 and ii) overall system experiment. Section 5.2.1 specifically reports the re-
800 sults obtained after automatically evaluating the content of summaries using
801 ROUGE tool, whereas Section 5.2.2 provides the results for the manually
802 conducted linguistic and readability evaluation. For both subsections, we
803 also compare NATSUM with respect to other summarization systems and
804 baselines.

805 *5.2.1. Automatic Evaluation*

806 The results shown in this section refer to the content assessment of the
807 narrative summaries generated by NATSUM compared to reference sum-
808 maries. As previously stated in Section 4, ROUGE was selected as the tool
809 for automatically evaluating our summaries, since it is a widespread summa-
810 rization evaluation tool that has been shown to correlate well with human
811 evaluations [66]. The most recent version of ROUGE (ROUGE-1.5.5) was
812 used.

813 Table 2 and Table 3 report the average ROUGE recall (R), precision (P)
814 and F1-measure (F) for the following metrics: ROUGE-1 and ROUGE-2—
815 compute the number of overlapping unigrams and bigrams, respectively—;
816 ROUGE-L—calculates the longest common subsequence between an auto-
817 matic and a reference summary—; and, ROUGE-SU4—measures the overlap
818 of skip-bigrams an automatic summary contains with respect to a model
819 one, with a maximum distance of four words between them—. The higher
820 the recall, precision and F1-measure values, the better.

821 The two tables differ in the input given for the Abstractive Summarization
822 module corresponding to the experimental scenarios described in Section 4.1:
823 i) the gold-standard, and ii) the overall experiment, respectively. Whereas in
824 Table 2, the input for this module is derived from the gold-standard timelines
825 available in the corpus, Table 3 reports the results of the system in a real
826 scenario, thus allowing us to also analyze how the overall system performs.

827 Furthermore, the “FirstEvent” refers to the narrative summary approach
828 generated, only taking into account the first event provided by the enriched
829 timeline, which is considered as a baseline. The “LongestEvent” refers to
830 an additional narrative summarization approach that takes into account,
831 for each line of the given timeline, the event with the higher number of

832 arguments, to generate a sentence from it. We also computed the per-
 833 formance of the extractive summarization approaches previously analyzed
 834 (COMPENDIUM, GRAFENO, OTS).

	ROUGE-1			ROUGE-2			ROUGE-L			ROUGE-SU4		
	R	P	F	R	P	F	R	P	F	R	P	F
COMPENDIUM	0.317	0.370	0.312	0.114	0.154	0.121	0.296	0.348	0.293	0.142	0.180	0.145
GRAFENO	0.285	0.415	0.295	0.102	0.199	0.118	0.261	0.384	0.272	0.127	0.140	0.139
OTS	0.305	0.362	0.303	0.106	0.148	0.114	0.280	0.335	0.280	0.133	0.173	0.138
FirstEvent	0.323	0.583	0.402	0.141	0.270	0.179	0.316	0.570	0.392	0.140	0.264	0.176
LongestEvent	0.351	0.688	0.445	0.166	0.335	0.215	0.340	0.665	0.431	0.165	0.339	0.214
NATSUM	0.576	0.735	0.637	0.420	0.544	0.467	0.559	0.714	0.619	0.400	0.518	0.445

Table 2: Average values for recall, precision and F1-measure for the gold-standard annotations ((i) gold-standard experiment). Comparison between different summarization and baseline approaches.

	ROUGE-1			ROUGE-2			ROUGE-L			ROUGE-SU4		
	R	P	F	R	P	F	R	P	F	R	P	F
COMPENDIUM	0.317	0.370	0.312	0.114	0.154	0.121	0.296	0.348	0.293	0.142	0.180	0.145
GRAFENO	0.285	0.415	0.295	0.102	0.199	0.118	0.261	0.384	0.272	0.127	0.140	0.139
OTS	0.305	0.362	0.303	0.106	0.148	0.114	0.280	0.335	0.280	0.133	0.173	0.138
FirstEvent	0.258	0.463	0.302	0.083	0.164	0.101	0.250	0.444	0.293	0.100	0.194	0.119
LongestEvent	0.251	0.524	0.312	0.088	0.196	0.114	0.245	0.510	0.305	0.099	0.225	0.125
NATSUM	0.433	0.595	0.470	0.263	0.363	0.284	0.422	0.579	0.457	0.260	0.360	0.282

Table 3: Average values for recall, precision and F1-measure when using raw data without any type of annotation as input ((ii) overall system experiment). Comparison between different summarization and baseline approaches in a real scenario.

835 For each table, rows 3-5 refer to the extractive summarization approaches,
 836 whereas rows 6-8 refer to abstractive summarization. The results indicate
 837 that regardless of the input type used for the Abstractive Summarization
 838 module (either the gold-standard timelines for event identification available
 839 in the corpus, or the ones produced by the Enriched Timeline Extraction
 840 module), our system outperforms the remaining ones. This means that in-
 841 tegrating the module for identifying events, as well as extracting temporal
 842 information enhances narrative summarization. When the complete system
 843 is evaluated, the results for the last two rows in Table 3 are lower than
 844 the corresponding ones in Table 2. This is explained by the errors that the
 845 Enriched Timeline Extraction module may introduce in the overall system.
 846 However, despite this issue, in both evaluations, NATSUM obtains better
 847 results than the others.

848 Table 4 and Table 5 provide the percentage of improvement obtained by
 849 NATSUM compared to the remaining summarization systems and baselines,
 850 taking only into account the F1-measure values.

	COMPENDIUM				GRAFENO				OTS				FirstEvent				LongestEvent			
	R1	R2	RL	RSU4	R1	R2	RL	RSU4	R1	R2	RL	RSU4	R1	R2	RL	RSU4	R1	R2	RL	RSU4
NATSUM	105	286	111	207	116	295	128	220	110	309	121	223	59	160	58	153	43	117	43	108

Table 4: Percentage of improvement for the F1-measure metric when comparing NATSUM with respect to the extractive summarization approaches and abstractive baselines for the gold-standard annotations ((i) gold-standard experiment). R1, R2, RL, and RSU4 refer to ROUGE-1, ROUGE-2, ROUGE-3 and ROUGE-SU4, respectively.

	COMPENDIUM				GRAFENO				OTS				FirstEvent				LongestEvent			
	R1	R2	RL	RSU4	R1	R2	RL	RSU4	R1	R2	RL	RSU4	R1	R2	RL	RSU4	R1	R2	RL	RSU4
NATSUM	51	135	56	95	59	140	68	103	55	149	65	105	56	182	56	137	51	153	50	125

Table 5: Percentage of improvement for the F1-measure metric when comparing NATSUM with respect to the extractive summarization approaches and abstractive baselines when using raw data without any type of annotation as input ((ii) overall system experiment). R1, R2, RL, and RSU4 refer to ROUGE-1, ROUGE-2, ROUGE-3 and ROUGE-SU4, respectively.

851 The results indicate that NATSUM performs better than other summa-
 852 rization approaches. This improvement is even greater when compared to the
 853 extractive summarization approaches. Moreover, despite the LongestEvent
 854 baseline being more competitive than the FirstEvent baseline, NATSUM is
 855 still capable of delivering a better performance. On the one hand, when
 856 considering the gold-standard timelines (i.e., only the Abstractive Summa-
 857 rization module without using the Enriched Timeline Extraction module),
 858 NATSUM’s performance increases for the F1-measure by 59% for ROUGE-
 859 1; 160% for ROUGE-2; 58% for ROUGE-L; and 153% for ROUGE-SU4
 860 compared to the FirstVerb baseline; and by 43% for ROUGE-1; 117% for
 861 ROUGE-2; 43% for ROUGE-L; and 108% for ROUGE-SU4 compared to the
 862 LongestEvent baseline. On the other hand, when considering the raw data
 863 without any kind of annotation as input—i.e. our complete approach, inte-
 864 grating both modules explained in Section 3—, NATSUM’s performance is
 865 also increased compared to the baselines as can be seen in Tables 3 and 5.

866 NATSUM also performs better than the multi-document entity-focused
 867 extractive summarization tested. The extractive summarization system with

868 the best F1-measure results for all ROUGE metrics —COMPENDIUM— is
869 improved by 51% for ROUGE-1, when our narrative abstractive approach is
870 compared to the best extractive summarization system in the real scenario—
871 i.e., with raw text as input data for the approach without any type of an-
872 notation on events—. When gold-standard timelines are considered, this
873 improvement increases by 105% for ROUGE-1.

874 Additionally, the use of NLG techniques does not decrease the perfor-
875 mance of the resulting summaries, as demonstrated by the results of Table 2,
876 when the input for the Abstractive Summarization module comes from gold
877 standard event and temporal annotations, thus indicating that NLG can
878 benefit abstractive summarization. This reconfirms our initial claim that
879 extractive summarization is not sufficient for generating effective narrative
880 summaries.

881 Finally, the main conclusion of this quantitative evaluation using ROUGE
882 is that NATSUM’s approach of integrating the Enriched Timeline Extraction
883 module for identifying co-referent events and temporal information in differ-
884 ent related documents, together with an Abstractive Summarization module
885 using NLG techniques is highly effective for producing narrative summaries.

886 In Example 13, a fragment of a generated narrative abstractive summary
887 about “Boeing” using our NATSUM system is shown.

- 888 (13) 2006-01: The first of the new airliner delivered to Pakistan International Airlines.
889 2007-06-10: The aircraft have a pre-modification catalogue value of US \$ 3.5 billion.
890 2007-07-07: Announced 35 new orders from German airline Air Berlin and ALAFCO
891 Aviation Lease & Finance of Kuwait.
892 2007-07-08: Boeing received a congratulatory letter from Airbus.
893 2007-07-08: The plane promises as it is the first model to be built out of plastic
894 and carbon composites, more lightweight than conventional materials.

895 5.2.2. Readability Evaluation

896 This section reports the results obtained for the manual readability eval-
897 uation. As previously explained in Section 4, a linguistic evaluation with
898 human assessors was also conducted to determine whether the abstractive
899 summaries were appropriate from a readability perspective.

900 For this evaluation, we only compared the abstractive summaries, NAT-
901 SUM and the two baselines — FirstEvent and LongestEvent— since they
902 used NLG techniques to create the summaries. Therefore, to verify the lin-
903 guistic quality of the generated content was more critical in this case, whereas

904 extractive summaries just copy and paste the same content available from
 905 the original documents.

906 Table 6 and Table 7 report the average results obtained for i) the gold-
 907 standard, and ii) the overall experiment, respectively.

	Readability/Fluency						Overall Responsiveness
	Grammaticality	Non-redundancy	Referential clarity	Focus	Structure and Coherence	Average	
FistEvent	2.47	2.70	2.73	2.42	1.97	2.46	2.16
LongestEvent	2.08	2.77	2.80	2.30	1.85	2.36	2.03
NATSUM	2.78	3.18	3.36	3.25	2.83	3.08	2.89

Table 6: Average values for readability/fluency (including the average values for summary’s grammaticality, non-redundancy, referential clarity, focus and structure and coherence) and for the summary’s overall responsiveness for the (i) gold-standard experiment.

	Readability/Fluency						Overall Responsiveness
	Grammaticality	Non-redundancy	Referential clarity	Focus	Structure and Coherence	Average	
FistEvent	2.52	2.81	2.84	3.00	2.33	2.70	2.74
LongestEvent	2.45	2.76	3.05	2.90	2.21	2.67	2.66
NATSUM	2.69	3.41	3.53	3.79	3.07	3.30	3.60

Table 7: Average values for readability/fluency (including the average values for summary’s grammaticality, non-redundancy, referential clarity, focus and structure and coherence) and for the summary’s overall responsiveness for the (ii) the overall system experiment.

908 As can be seen in the tables, in both experiments NATSUM obtains better
 909 results than the ones obtained by the two baselines. These results indicate
 910 that NATSUM improves the linguistic quality of the generated summaries
 911 in comparison to the baselines, thus corroborating the results achieved in
 912 the automatic evaluation. In terms of readability/fluency results, the sum-
 913 maries generated by NATSUM have a higher structure and coherence than
 914 the baselines summaries. In addition to this, they present less redundancy
 915 and more referential clarity as well as more grammaticality than the ones
 916 from the baselines, maintaining a better focused summary. Moreover, in
 917 terms of overall responsiveness, NATSUM summaries have scored higher for
 918 both experiments.

919 Furthermore, as mentioned, a human relevance judgement evaluation was
 920 carried out. In this case, the assessors preferred the summaries generated by
 921 NATSUM for both experiments –79.45% and 79.66% for the gold-standard
 922 and overall experiments, respectively–.

923 6. Assessing NATSUM in the context of Timeline Summarization

924 To the best of our knowledge, there is no specific dataset with reference
925 summaries that could be appropriate for the specific features of NATSUM
926 (i.e., narrative chronological abstractive summarization). However, having
927 obtaining good results in the evaluation conducted in Section 4.2, it would
928 be also important to validate these results and findings by benchmarking
929 NATSUM against additional existing datasets developed for a similar task
930 (i.e., timeline summarization). Besides the comparison with the extractive
931 systems already used throughout this research work (i.e., COMPENDIUM
932 [62], GRAFENO [63] and Open Text Summarizer (OTS) [64]), this would
933 allow us to compare NATSUM with more task-oriented and focused state-
934 of-the-art systems.

935 Summaries generated for the task of timeline summarization mainly dif-
936 fer from those generated by NATSUM in that the latter aims to generate
937 narrative summaries and not timelines. In the case of NATSUM, timelines
938 constitute the means to generate the final narrative structure. In this sense,
939 the input of the abstractive summarization module is not a set of documents,
940 but a target oriented timeline. In contrast, in the case of timeline summa-
941 rization, the final aim is to generate a timeline that serves as the summary
942 of one or more input documents.

943 Regardless of these differences, and considering that the final timelines in
944 timeline summarization contain short summaries temporally ordered by the
945 document creation time, NATSUM is evaluated using an specific available
946 dataset for the task of timeline summarization. The dataset finally chosen
947 for the evaluation and comparison is Timeline17 dataset, which is the one
948 used in [43] and [44]. The reasons for using this dataset were twofold. On the
949 one hand, it was selected because it is available online²² and, on the other
950 hand, a comparison with other timeline summarization systems is presented
951 as well. Therefore, using the same dataset, the ultimate goal of this evalu-
952 ation is to compare NATSUM with all the timeline summarization systems
953 presented in [43] and [44], as well as compared it with the extractive multi-
954 document summarization systems presented throughout this research work
955 (COMPENDIUM, OTS and GRAFENO) to confirm and validate whether
956 the summaries generated by NATSUM offer an added value with respect to
957 a standard timeline extractive summary.

²²<http://www.l3s.de/~gtran/timeline/>

958 In the next subsections, we describe the dataset in more detail (Section
959 6.1) together with the results obtained (Section 6.2).

960 *6.1. Timeline17 Dataset Description*

961 This dataset is composed of news articles from different media outlets
962 about 9 different topics: BP Oil, Michael Jackson Death, H1N1, Haiti Earth-
963 quake, Financial Crisis, Libyan War, Iraq War, Egyptian Protest, and Syrian
964 Crisis. The dataset, created by the authors of [43] and [44], was gathered in
965 two steps:

- 966 • Collecting human timelines (ground truth): They collected available
967 timelines published by popular news agencies such as CNN, BBC, NBC-
968 news, etc. that discuss the previous 9 topics. From these topics, 17
969 timelines were manually built. This human timelines are the gold stan-
970 dard (i.e., reference summaries) for the evaluation performed in the
971 next section.
- 972 • Retrieving news articles: For each timeline, they used Google Web
973 Search Engine²³ to retrieve news articles from the same news agency
974 of the timeline (i.e. BBC news articles for BBC-published timeline,...)
975 using topics as query. In the end, they obtained 4,650 news articles
976 after removing duplicate news. All these news articles are the input to
977 NATSUM system.

978 *6.2. Results and Comparison with Timeline Summarization Systems*

979 In order to apply NATSUM to the timeline summarization dataset de-
980 scribed in the previous section, the system needs to use the different top-
981 ics as target entities for each timeline generated (BP Oil, Michael Jackson
982 Death, H1N1, Haiti Earthquake, Financial Crisis, Libyan War, Iraq War,
983 Egyptian Protest and Syrian Crisis). Then, the two modules of the proposal
984 are applied to the input documents to create the different narrative abstrac-
985 tive summaries. Once the summaries were generated, they were evaluated
986 with ROUGE with respect to the reference timeline summaries available in
987 the dataset. In order to evaluate the summaries under the same conditions,
988 ROUGE was set to truncate the length of the generated summaries to the
989 same length as the reference timelines had.

²³<https://www.google.com/>

990 Table 8 reports the average F1-measure (F) results for ROUGE-1, ROUGE-
 991 2 and ROUGE-SU4 results. Rows 3-5 refer to the performance of the ex-
 992 tractive summarization approaches previously analyzed (COMPENDIUM,
 993 GRAFENO, OTS), whereas rows 6-10 refers to the timeline summarization
 994 systems presented in [43] and [44]. Finally, the last row provides NATSUM
 995 performance²⁴.

	ROUGE-1	ROUGE-2	ROUGE-SU4
	F	F	F
COMPENDIUM [62]	0.340	0.085	0.133
GRAFENO [63]	0.267	0.069	0.102
OTS [64]	0.337	0.076	0.127
Chieu et al.[39]	0.202	0.037	0.041
MEAD[67]	0.208	0.049	0.039
ETS[40]	0.207	0.047	0.042
Tran Linear Regression[44]	0.218	0.050	0.046
Tran LTR[43]	0.230	0.053	0.050
NATSUM	0.413	0.121	0.176

Table 8: Average F1-measure values when using Timeline17 dataset as input. Comparison between different multi-document and timeline summarization approaches.

996 As shown in Table 8, NATSUM greatly overperforms timeline summa-
 997 rization systems for all ROUGE metrics, being the main reason that the
 998 summarization module is using an enriched timeline as input. The approach
 999 exploits not only the temporal information about the document creation
 1000 time (as timeline summarization does) but also all the temporal links and
 1001 expressions related to the events referring to the target entity across different
 1002 documents. This implies a temporal information processing that goes further
 1003 in terms of exploiting temporal information than merely using the document
 1004 creation time. Furthermore, NATSUM approach is using the events in the
 1005 timeline, and their arguments, to generate a sentence that covers all the argu-
 1006 ments of the event. Since NATSUM is dealing with the coreference of events,
 1007 for the same event, named in different ways in different documents, our final

²⁴Only F1-measure for ROUGE-1, ROUGE-2, and ROUGE-SU4 is presented since this is the measure reported in referenced papers.

1008 summary is generating a single sentence which condenses all the information
1009 related to the event in question, which results in avoiding redundancy in the
1010 resulting summary. Furthermore, the results obtained corroborate the pre-
1011 vious evaluation of NATSUM in comparison with extractive multi-document
1012 summarization systems. Despite using a different input corpora, NATSUM
1013 performs better than COMPENDIUM, GRAFENO and OTS. It is also worth
1014 noting that extractive summaries obtain higher ROUGE results than timeline
1015 summaries. This could be explained by the fact that those systems are very
1016 competitive as far as detecting relevant information from input documents
1017 is concerned.

1018 Finally, the results also indicate that providing a narrative abstractive
1019 summary instead of just a timeline summary is better, since besides includ-
1020 ing dates, they also provide relevant information that is generate from the
1021 information found in different sources about the same event. This validates
1022 the appropriateness of the NLG techniques used within the NATSUM system
1023 for generating abstractive summaries.

1024 7. Conclusions

1025 This work presents NATSUM, a narrative abstractive summarization ap-
1026 proach that integrates structured timeline knowledge together with natural
1027 language generation techniques to enhance the creation of such type of sum-
1028 maries. Our integrated approach was motivated by two aspects: First, it is
1029 based on the fact that humans tend to apply chronological ordering of events
1030 in the summarizing process, which implies the need for timelines. Second,
1031 when using an abstractive summarization approach, rather than an extrac-
1032 tive one, the relevant information (e.g., *who? what?, when?, where?,...*)
1033 can be fused together, leading to the generation of more complete sentences,
1034 and thus, more comprehensible and effective summaries. Hence, NATSUM's
1035 architecture comprises two main modules: i) Enriched Timeline Extraction
1036 module, and ii) Abstractive Summarization module. The former module uses
1037 a set of plain news documents and a target entity as input, and obtains a
1038 structured timeline document plan that is enriched with all the arguments of
1039 each event involved in the timeline for the particular target entity. Specifi-
1040 cally, for each line of the timeline, there is a cluster with the exact date of the
1041 event and a set of event mentions together with their arguments, extracted
1042 from different documents, that refer to the same event. The latter module
1043 generates a narrative abstractive summary using the enriched timeline. For

1044 this, a hybrid surface realization approach, based on over-generation and
1045 ranking techniques is used.

1046 The evaluation conducted and the results obtained show that extractive
1047 summaries lose between 22% (OTS) and 38%(GRAFENO) of the *events* re-
1048 lated with the target entity; and between 7% (GRAFENO) and 19% (COM-
1049 PENDIUM) of the *temporal information*. Moreover, regarding the content
1050 evaluation of the narrative abstractive summaries, the F1-measure for all
1051 ROUGE metrics improves by at least 50% in the worst case, when our nar-
1052 rative abstractive system (NATSUM) is compared to the extractive summa-
1053 rization systems, as well as to the baselines in the real scenario—i.e., with
1054 raw text as input data for the approach without any type of annotation about
1055 events—. Remarkable improvements are also obtained for the gold-standard
1056 experiment.

1057 In addition, a manual evaluation was carried out between the summaries
1058 generated by the two baselines and NATSUM to measure the readability/fluency
1059 and overall responsiveness of the summaries. The results obtained corrob-
1060 orate the ones from the automatic evaluation, with the summaries from
1061 NATSUM being better than both of the baseline ones for both experiments
1062 ((i) gold-standard and (ii) overall experiments). Besides, a human relevance
1063 judgement evaluation was performed, where the NATSUM summaries were
1064 preferred in almost 80% of the cases for both experiments. Finally, in order
1065 to compare NATSUM with other systems, a timeline summarization dataset
1066 is used as input, since it is the most similar task to our proposal, conclud-
1067 ing that NATSUM greatly improves the results obtained by state-of-the-art
1068 timeline summarization and extractive systems.

1069 Although NATSUM has shown very good and promising results, also im-
1070 proving the performance of extractive summarization approaches, there are
1071 several aspects to consider for future development concerning the individual
1072 modules that are integrated into NATSUM. First, the Enriched Timeline
1073 Extraction module should be improved to better identify co-referent events
1074 and temporal relationships between events, especially when these relation-
1075 ships are implicit. This would narrow the gap between the results obtained
1076 when using gold-standard timelines. Second, the Abstractive Summarization
1077 module should be improved so that it would include appropriate discourse
1078 markers for connecting individual sentences to increase the coherence of the
1079 produced narrative summaries, rather than listing a set of relevant newly gen-
1080 erated sentences. This would enhance the quality of the resulting narrative
1081 summaries generated by NATSUM.

1082 **Acknowledgements**

1083 This research work has been partially funded by the Spanish Government
1084 through projects TIN2015-65100-R, TIN2015-65136-C2-2-R, as well as by the
1085 project “Análisis de Sentimientos Aplicado a la Prevención del Suicidio en
1086 las Redes Sociales (ASAP)” funded by Ayudas Fundación BBVA a equipos
1087 de investigación científica. Moreover, it has been also funded by Generalitat
1088 Valenciana through project “SIIA: Tecnologías del lenguaje humano para
1089 una sociedad inclusiva, igualitaria, y accesible” with grant reference PROM-
1090 ETEO/2018/089

1091 **References**

- 1092 [1] E. Lloret, M. Palomar, Text summarisation in progress: A literature
1093 review, *Artif. Intell. Rev.* 37 (2012) 1–41.
- 1094 [2] I. Mani, *Advances in Automatic Text Summarization*, MIT Press, Cam-
1095 bridge, MA, USA, 1999.
- 1096 [3] J. Gottschall, *The Storytelling Animal*, Houghton Mifflin Harcourt,
1097 2012.
- 1098 [4] M. R. Hovav, E. Doron, I. Sichel, *Lexical Semantics, Syntax, and Event*
1099 *Structure*, Oxford University Press, Oxford, 2010.
- 1100 [5] I. Mani, J. Pustejovsky, R. Gaizauskas, *The Language of Time*, Oxford
1101 University Press, Oxford, 2005.
- 1102 [6] H. Ji, R. Grishman, Z. Chen, P. Gupta, Cross-document event ex-
1103 traction and tracking: Task, evaluation, techniques and challenges, in:
1104 *Proceedings of the International Conference RANLP-2009*, Association
1105 for Computational Linguistics, 2009, pp. 166–172.
- 1106 [7] M. Palmer, D. Gildea, P. Kingsbury, The Proposition Bank: An Anno-
1107 tated Corpus of Semantic Roles, *Computational Linguistics* 31 (2005).
- 1108 [8] A.-L. Minard, M. Speranza, E. Agirre, I. Aldabe, M. van Erp,
1109 B. Magnini, G. Rigau, R. Urizar, Semeval-2015 task 4: Timeline:
1110 Cross-document event ordering, in: *Proceedings of the 9th Interna-*
1111 *tional Workshop on Semantic Evaluation, SemEval ’15*, Association for
1112 Computational Linguistics, 2015, pp. 778–786.

- 1113 [9] Semeval-2015, International Workshop on Semantic Evaluation, 2015.
- 1114 [10] C. D. Manning, M. Surdeanu, J. Bauer, J. Finkel, S. J. Bethard, D. Mc-
1115 Closky, The Stanford CoreNLP natural language processing toolkit, in:
1116 Proceedings of 52nd Annual Meeting of the Association for Computa-
1117 tional Linguistics: System Demonstrations, pp. 55–60.
- 1118 [11] T. Caselli, A. Fokkens, R. Morante, P. Vossen, SPINOZA_VU: An NLP
1119 Pipeline for Cross Document TimeLines, in: Proceedings of the 9th
1120 International Workshop on Semantic Evaluation (SemEval 2015), As-
1121 sociation for Computational Linguistics, Denver, Colorado, 2015, pp.
1122 787–791.
- 1123 [12] B. Moulahi, J. Strötgen, M. Gertz, L. Tamine, HeidelToul: A Baseline
1124 Approach for Cross-document Event Ordering, in: Proceedings of the
1125 9th International Workshop on Semantic Evaluation (SemEval 2015),
1126 Association for Computational Linguistics, Denver, Colorado, 2015, pp.
1127 825–829.
- 1128 [13] B. Navarro, E. Saquete, Gplsiua: Combining temporal information and
1129 topic modeling for cross-document event ordering, in: Proceedings of the
1130 9th International Workshop on Semantic Evaluation (SemEval 2015),
1131 Association for Computational Linguistics, Denver, Colorado, 2015, pp.
1132 820–824.
- 1133 [14] H. Llorens, E. Saquete, B. Navarro-Colorado, Applying Semantic
1134 Knowledge to the Automatic Processing of Temporal Expressions and
1135 Events in Natural Language, *Information Processing & Management* 49
1136 (2013) 179–197.
- 1137 [15] H. Llorens, E. Saquete, B. Navarro-Colorado, Automatic System for
1138 Identifying and Categorizing Temporal Relations in Natural Language,
1139 *International Journal of Intelligent Systems* 27 (2012) 680–703.
- 1140 [16] E. Laparra, R. Agerri, I. Aldabe, G. Rigau, Multilingual and cross-
1141 lingual timeline extraction, *CoRR* abs/1702.00700 (2017).
- 1142 [17] E. Reiter, R. Dale, *Building Natural Language Generation Systems*,
1143 Cambridge University Press, New York, NY, USA, 2000.

- 1144 [18] A. Gatt, E. Reiter, Simplenlg: A realisation engine for practical appli-
1145 cations, in: Proceedings of the 12th European Workshop on Natural
1146 Language Generation, ENLG '09, Association for Computational Lin-
1147 guistics, Stroudsburg, PA, USA, 2009, pp. 90–93.
- 1148 [19] D. A. Smith, H. Lieberman, Generating and interpreting referring ex-
1149 pressions as belief state planning and plan recognition, in: A. Gatt,
1150 H. Saggion (Eds.), ENLG 2013 - Proceedings of the 14th European
1151 Workshop on Natural Language Generation, August 8-9, 2013, Sofia,
1152 Bulgaria, The Association for Computer Linguistics, 2013, pp. 61–71.
- 1153 [20] D. Duma, E. Klein, Generating natural language from linked data:
1154 Unsupervised template extraction, in: Proceedings of the 10th Inter-
1155 national Conference on Computational Semantics (IWCS 2013) – Long
1156 Papers, Association for Computational Linguistics, Potsdam, Germany,
1157 2013, pp. 83–94.
- 1158 [21] I. Konstas, M. Lapata, Inducing document plans for concept-to-text gen-
1159 eration, in: Proceedings of the 2013 Conference on Empirical Methods
1160 in Natural Language Processing, Association for Computational Lin-
1161 guistics, Seattle, Washington, USA, 2013, pp. 1503–1514.
- 1162 [22] R. Kondadadi, B. Howald, F. Schilder, A statistical nlg framework for
1163 aggregated planning and realization, in: Proceedings of the 51st An-
1164 nual Meeting of the Association for Computational Linguistics (Volume
1165 1: Long Papers), Association for Computational Linguistics, Sofia, Bul-
1166 garia, 2013, pp. 1406–1415.
- 1167 [23] M. E. Vicente, C. Barros, E. Lloret, Statistical language modelling for
1168 automatic story generation, *Journal of Intelligent and Fuzzy Systems*
1169 34 (2018) 3069–3079.
- 1170 [24] D. Dannélls, Multilingual text generation from structured formal repre-
1171 sentations., University of Gothenburg, Göteborg, 2012.
- 1172 [25] C. D. Manning, H. Schütze, Foundations of Statistical Natural Language
1173 Processing, MIT Press, Cambridge, MA, USA, 1999.
- 1174 [26] A. Mnih, Y. W. Teh, A fast and simple algorithm for training neural
1175 probabilistic language models, in: Proceedings of the 29th Interna-

- 1176 tional Conference on Machine Learning, ICML 2012, Edinburgh, Scot-
1177 land, UK, June 26 - July 1, 2012.
- 1178 [27] M. Ballesteros, B. Bohnet, S. Mille, L. Wanner, Data-driven sentence
1179 generation with non-isomorphic trees, in: Proceedings of the 2015 Con-
1180 ference of the North American Chapter of the Association for Com-
1181 putational Linguistics: Human Language Technologies, Association for
1182 Computational Linguistics, Denver, Colorado, 2015, pp. 387–397.
- 1183 [28] A. See, P. J. Liu, C. D. Manning, Get to the point: Summarization with
1184 pointer-generator networks, in: Proceedings of the 55th Annual Meeting
1185 of the Association for Computational Linguistics (Volume 1: Long Pa-
1186 pers), Association for Computational Linguistics, 2017, pp. 1073–1083.
- 1187 [29] Y.-C. Chen, M. Bansal, Fast abstractive summarization with reinforce-
1188 selected sentence rewriting, in: Proceedings of the 56th Annual Meeting
1189 of the Association for Computational Linguistics (Volume 1: Long Pa-
1190 pers), Association for Computational Linguistics, 2018, pp. 675–686.
- 1191 [30] J. Cordeiro, G. Dias, P. Brazdil, Rule induction for sentence reduction,
1192 in: L. Correia, L. P. Reis, J. Cascalho (Eds.), Progress in Artificial
1193 Intelligence, Springer Berlin Heidelberg, Berlin, Heidelberg, 2013, pp.
1194 528–539.
- 1195 [31] M. Valizadeh, P. Brazdil, Exploring actor–object relationships for query-
1196 focused multi-document summarization, *Soft Computing* 19 (2015)
1197 3109–3121.
- 1198 [32] N. Chambers, D. Jurafsky, Unsupervised learning of narrative event
1199 chains, in: K. R. McKeown, J. D. Moore, S. Teufel, J. Allan, S. Fu-
1200 rui (Eds.), ACL 2008, Proceedings of the 46th Annual Meeting of the
1201 Association for Computational Linguistics, June 15-20, 2008, Colum-
1202 bus, Ohio, USA, The Association for Computer Linguistics, 2008, pp.
1203 789–797.
- 1204 [33] N. Chambers, D. Jurafsky, Unsupervised learning of narrative schemas
1205 and their participants, in: K. Su, J. Su, J. Wiebe (Eds.), ACL 2009,
1206 Proceedings of the 47th Annual Meeting of the Association for Computa-
1207 tional Linguistics and the 4th International Joint Conference on Natural

- 1208 Language Processing of the AFNLP, 2-7 August 2009, Singapore, The
1209 Association for Computer Linguistics, 2009, pp. 602–610.
- 1210 [34] J. C. K. Cheung, H. Poon, L. Vanderwende, Probabilistic Frame Induc-
1211 tion (2013).
- 1212 [35] N. Chambers, Event Schema Induction with a Probabilistic Entity-
1213 Driven Model, Proceedings of the 2013 Conference on Empirical Meth-
1214 ods in Natural Language Processing (EMNLP 2013) (2013) 1797–1807.
- 1215 [36] N. Mostafazadeh, From Event to Story Understanding, Ph.D. thesis,
1216 University of Rochester, 2017.
- 1217 [37] K. Markert, S. Martschat, Improving ROUGE for timeline summariza-
1218 tion, in: M. Lapata, P. Blunsom, A. Koller (Eds.), Proceedings of the
1219 15th Conference of the European Chapter of the Association for Com-
1220 putational Linguistics, EACL 2017, Valencia, Spain, April 3-7, 2017,
1221 Volume 2: Short Papers, Association for Computational Linguistics,
1222 2017, pp. 285–290.
- 1223 [38] J. Allan, R. Gupta, V. Khandelwal, Temporal summaries of news topics,
1224 in: W. B. Croft, D. J. Harper, D. H. Kraft, J. Zobel (Eds.), SIGIR 2001:
1225 Proceedings of the 24th Annual International ACM SIGIR Conference
1226 on Research and Development in Information Retrieval, September 9-13,
1227 2001, New Orleans, Louisiana, USA, ACM, 2001, pp. 10–18.
- 1228 [39] H. L. Chieu, Y. K. Lee, Query based event extraction along a timeline,
1229 in: Proceedings of the 27th Annual International ACM SIGIR Confer-
1230 ence on Research and Development in Information Retrieval, SIGIR '04,
1231 ACM, New York, NY, USA, 2004, pp. 425–432.
- 1232 [40] R. Yan, X. Wan, J. Otterbacher, L. Kong, X. Li, Y. Zhang, Evolu-
1233 tionary timeline summarization: A balanced optimization framework
1234 via iterative substitution, in: In Proceedings of the 34th International
1235 ACM SIGIR Conference on Research and Development in Information
1236 Retrieval, SIGIR '11, ACM, 2011, pp. 745–754.
- 1237 [41] G. Tran, M. Alrifai, E. Herder, Timeline summarization from rele-
1238 vant headlines, in: A. Hanbury, G. Kazai, A. Rauber, N. Fuhr (Eds.),
1239 Advances in Information Retrieval, Springer International Publishing,
1240 Cham, 2015, pp. 245–256.

- 1241 [42] W. Y. Wang, Y. Mehdad, D. R. Radev, A. Stent, A low-rank approxima-
1242 tion approach to learning joint embeddings of news stories and images
1243 for timeline summarization, in: Proceedings of the 2016 Conference of
1244 the North American Chapter of the Association for Computational Lin-
1245 guistics: Human Language Technologies, Association for Computational
1246 Linguistics, 2016, pp. 58–68.
- 1247 [43] G. B. Tran, A. T. Tran, N.-K. Tran, M. Alrifai, N. Kanhabua, Lever-
1248 age learning to rank in an optimization framework for timeline sum-
1249 marization, SIGIR 2013 Workshop on Time-aware Information Access
1250 (TAIA’2013) (2013).
- 1251 [44] G. Binh Tran, M. Alrifai, D. Quoc Nguyen, Predicting relevant news
1252 events for timeline summaries, in: Proceedings of the 22Nd International
1253 Conference on World Wide Web, WWW ’13 Companion, ACM, New
1254 York, NY, USA, 2013, pp. 91–92.
- 1255 [45] B. Navarro-Colorado, E. Saquete, Cross-document event ordering
1256 through temporal, lexical and distributional knowledge, Knowl.-Based
1257 Syst. 110 (2016) 244–254.
- 1258 [46] A.-L. Minard, M. Speranza, R. Urizar, B. Altuna, M. van Erp,
1259 A. Schoen, C. van Son, Meantime, the newsreader multilingual event
1260 and time corpus, in: Proceedings of the Tenth International Conference
1261 on Language Resources and Evaluation (LREC 2016).
- 1262 [47] R. Saurí, J. Littman, R. Knippen, R. Gaizauskas, A. Set-
1263 zer, J. Pustejovsky, TimeML Annotation Guidelines 1.2.1
1264 (<http://www.timeml.org/>), 2006.
- 1265 [48] J. Pustejovsky, J. M. Castaño, R. Ingria, R. Saurí, R. J. Gaizauskas,
1266 A. Setzer, G. Katz, D. R. Radev, Timeml: Robust specification of
1267 event and temporal expressions in text, in: M. T. Maybury (Ed.), New
1268 Directions in Question Answering, Papers from 2003 AAAI Spring Sym-
1269 posium, Stanford University, Stanford, CA, USA, AAAI Press, 2003, pp.
1270 28–34.
- 1271 [49] I. T. W. Group, ISO TimeML TC37 draft international standard DIS
1272 24617-1, 2008.

- 1273 [50] C. A. Bejan, S. Harabagiu, Unsupervised Event Coreference Resolution,
1274 Computational Linguistics 40 (2014) 311–347.
- 1275 [51] J. R. Firth, Papers in Linguistics (1934-1951), Oxford University Press,
1276 Oxford, 1957.
- 1277 [52] Z. Harris, Mathematical structures of language, Wiley, New York, 1968.
- 1278 [53] P. Turney, P. Pantel, From frequency to meaning: Vector space models
1279 of semantics, Journal of Artificial Intelligence Research 37 (2010) 141–
1280 188.
- 1281 [54] P. Gärdenfors, The geometry of meaning. Semantics based on conceptual
1282 spaces., MIT Press, Cambridge, Mass., 2014.
- 1283 [55] L. Padró, E. Stanilovsky, FreeLing 3.0: Towards Wider Multilinguality,
1284 in: Proceedings of the Language Resources and Evaluation Conference
1285 (LREC 2012), ELRA, Istanbul, Turkey.
- 1286 [56] J. Mitchell, M. Lapata, Composition in Distributional Models of Se-
1287 mantics, Cognitive Science 34 (2010) 1388–1429.
- 1288 [57] J. A. Bilmes, K. Kirchhoff, Factored language models and generalized
1289 parallel backoff, in: Proceedings of the 2003 Conference of the North
1290 American Chapter of the Association for Computational Linguistics on
1291 Human Language Technology: Companion Volume of the Proceedings
1292 of HLT-NAACL 2003–short Papers - Volume 2, pp. 4–6.
- 1293 [58] A. Stolcke, Srlm – an extensible language modeling toolkit, in: IN
1294 PROCEEDINGS OF THE 7TH INTERNATIONAL CONFERENCE
1295 ON SPOKEN LANGUAGE PROCESSING (ICSLP 2002, pp. 901–904.
- 1296 [59] K. K. Schuler, Verbnet: A Broad-coverage, Comprehensive Verb Lexi-
1297 con, Ph.D. thesis, 2005.
- 1298 [60] C. Fellbaum, WordNet: An Electronic Lexical Database., MIT Press,
1299 1998.
- 1300 [61] A. Isard, C. Brockmann, J. Oberlander, Individuality and alignment
1301 in generated dialogues, in: Proceedings of the INLG, Association for
1302 Computational Linguistics, 2006, pp. 25–32.

- 1303 [62] E. Lloret, M. Palomar, COMPENDIUM: a text summarisation tool
1304 for generating summaries of multiple purposes, domains, and genres,
1305 Natural Language Engineering 19 (2013) 147–186.
- 1306 [63] A. F. Sevilla, A. Fernandez-Isabel, A. Díaz, Enriched semantic graphs
1307 for extractive text summarization, in: Conference of the Spanish As-
1308 sociation for Artificial Intelligence, Springer International Publishing,
1309 Springer International Publishing, 2016.
- 1310 [64] F. Andonov, V. Slavova, G. Petrov, On the open text summarizer,
1311 International Journal “Information Content and Processing” 3 (2016)
1312 278–287.
- 1313 [65] C.-Y. Lin, ROUGE: A Package for Automatic Evaluation of Summaries,
1314 in: Text Summarization Branches Out: Proceedings of the Associa-
1315 tion for Computational Linguistics Workshop, Association for Compu-
1316 tational Linguistics, 2004, pp. 74–81.
- 1317 [66] C.-Y. Lin, E. Hovy, Automatic evaluation of summaries using n-gram
1318 co-occurrence statistics, in: Proceedings of the 2003 Conference of the
1319 North American Chapter of the Association for Computational Linguis-
1320 tics on Human Language Technology - Volume 1, NAACL ’03, Associ-
1321 ation for Computational Linguistics, Stroudsburg, PA, USA, 2003, pp.
1322 71–78.
- 1323 [67] D. Radev, T. Allison, S. Blair-Goldensohn, J. Blitzer, A. Çelebi, S. Dim-
1324 itrov, E. Drabek, A. Hakim, W. Lam, D. Liu, J. Otterbacher, H. Qi,
1325 H. Saggion, S. Teufel, M. Topper, A. Winkel, Z. Zhang, Mead - a
1326 platform for multidocument multilingual text summarization, in: Pro-
1327 ceedings of the Fourth International Conference on Language Resources
1328 and Evaluation (LREC’04), European Language Resources Association
1329 (ELRA), 2004.