

Information Extraction using Document Spanners

Information Extraction

Information extraction is the concept that involves structuring and combining data, whether it is explicitly stated or inferred, in texts. For a rule language, its value and tradeoffs can be analyzed by two concepts: expressivity and simplicity. Whether expressivity allows greater quality on the results and simplicity allows better human understandability, higher throughput and better performance to train models. Nevertheless, the lack of formalisms to abstract and define the expressivity, has made it difficult to analyze and compare the rule languages.

Document Spanners

In order to solve the lack of formalisms, a framework called document spanners has been defined by leveraging the principles of database management and setting up mathematical formulas to abstract a rule language.

A document spanner models a program for Information Extraction (IE) as a function that takes as input a text document (string over a finite alphabet) and produces a relation of spans (intervals in the document) over a predefined schema.

> String: Jon Doe is an artist. $\Upsilon st = \Sigma^* \cdot x \{ Jon \lor Doe \} \cdot \Sigma^*$ Spans: $[1,3\rangle [5,7\rangle$

State-of-the-Art

Document spanners on SpLog	Decis
Freydenberger proposes Subword Property Logic (SpLog), a subset of EC ^{REG} that has the same expressive power as any of the originally proposed core spanners and define methods to transform from core spanner logic to SpLog.	Freydenberg performance patterns (w terminals), w = y, where expressions They found operators ar techniques their use. Th
Coupling ontologies with document spanner	Recursive
Lembo and Scafoglieri introduce the concept of inking document spanners with existing ontologies following the Ontology-based Data Access (OBDA) framework. Such framework provides users with access to the information in their data sources through a three-level architecture, constituted by the ontology, the sources, and the mappings between the two. The Ontology-based Document Spanning (OBDS) system creates the mappings for the OBDA, associating the document spanner to the conjunctive query (CQ). Furthermore, they do not only introduce the concept, but they provided an algorithm that rewrites every CQ issued over a DL- Lite _R OBDS System into a spanner through means of the GLAV schema .	Peterfreund (spanners burregular exp language RC Datalog over four main to than the original can tell whe length), 2) determine wo string (avoid equality pre- can be evan RGXlog capa document so need to exp



sion problems performance

Holldack compared and the er document spanners against of consisting of variables and word word equations (equations of the form x x and y are patterns), and regular (which allow repetition operators). out that string equality and join re really expensive and thus that some will need to be developed to restrict is will lead to the creation of subclasses of document spanners.

programs for document spanners

et al. go beyond regular spanners uilt on top of non-recursive Datalog and ressions) and develop the spanner GXlog which is made utilizing recursive r regular expressions. They found out things: 1) RGXlog is more expressible ginally proposed core spanners (since it ether or not two spans have the same it is possible to write a program to whether or not two spans are the same ling the necessity to include the stringedicate), 3) that every RGXlog program aluated in polynomial time (making able of expressing all polynomial time panners), and finally 4) that they only press their regex formulas utilizing two variables in order to make it polynomial time.