

Extracting Biological Processes with Global Constraints

Author 1

XYZ Company
111 Anywhere Street
Mytown, NY 10000, USA
author1@xyz.org

Author 2

ABC University
900 Main Street
Ourcity, PQ, Canada A1A 1T2
author2@abc.ca

Abstract

Reasoning over processes is fundamental for language understanding applications such as Question Answering. In this paper we propose a method for extracting relations between events in a process. We annotate 150 paragraphs describing biological processes and show that by taking advantage of the global structure of a process we can substantially improve performance. In addition, we release our data set.

1 Introduction

Motivation: Being able to reason over processes is crucial for language understanding applications. Consider for example the paragraph in Figure 1, which describes the process of ATP synthesis. A human reading this paragraph will be able to answer complex questions that require understanding of the process structure such as

1. *How do H^+ ions contribute to the production of ATP?*
2. *What causes the rotor to spin?*
3. *In ATP synthesis, what happens if the rotor fails to spin?*

All these questions require understanding and reasoning over processes, and thus systems that have only bag-of-words representations will fail. In this paper we suggest a method for extracting a process structure that will facilitate answering of complex questions.

Relation to previous work: Extracting processes is related to two lines of works in Information Extraction - event extraction and timeline construction. Recent work in event extraction (Riedel and McCallum, 2011; McClosky et al., 2011) is based on BioNLP challenges and focuses on extraction of a closed set of events such as *regulation* and *phosphorylation* from a single sentence and their relations to proteins. However, a process is typically described over multiple sentences and involves a large number of possible events. Work on timeline construction (Do et al., 2012; McClosky and Manning, 2012) requires partially ordering a set of events that is described in a sequence of sentence. However, fully capturing process structure requires a rich set of relations (*cause*, *super*) that is missing from this line of work.

Emphasizing this work: In this paper, we find the structure of a biological process by extracting the relations between the process events. Properties of our task (a) spans multiple sentences (b) open-ended set of events (c) rich set of relations comparing to timeline construction (d) the nature of the text - it is a textbook rather than abstracts. (e) We do not use domain-specific knowledge. Some sentence that says that by doing this we will be able to answer the complex question from the beginning - linking to language understanding.

Technical contribution Processes have a global structure and we want to take advantage of that when extracting the relations. For example, all events a process description are connected to one another and there are various constraints such as if two event mentions refer to the same event then they must be

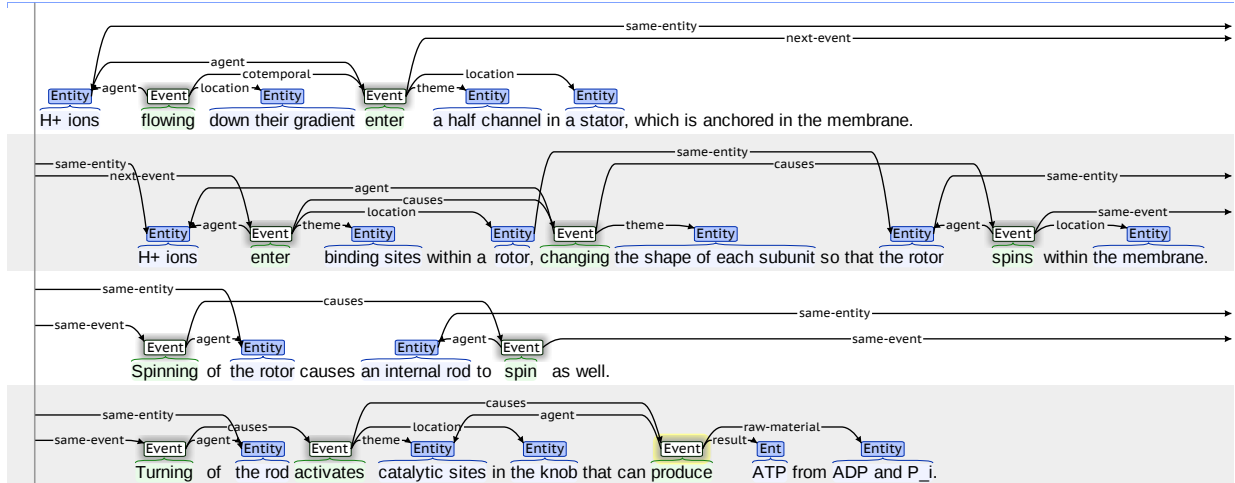


Figure 1: An annotation of the ATP synthesis process

related in a similar way to a third event. Similar to many recent works in NLP () we model global constraints using ILP, however since many of the constraints can be violated we use soft constraints. We show that by encoding global constraints we can substantially improve performance.

Contributions

- We define the task of process extraction - what is a process - we also identify properties of the process structure
- We propose a method for process extraction that uses global constraints and show that it improves performance
- We release a set of 150 biological processes, annotated by biologists.

structure Background, Definition and properties of processes, Local classifier (old features and new features), Global model, Experiments and maybe analysis

2 Related Work

BioNLP work

Timeline construction work.

Scripts work - Chambers, Poon 2013.

Work that uses global constraints with ILP or dual decomposition or whatever.

3 Process Definition and Data Set

A process is defined by a series of events, where each event involves some participants, and also by relations between the events (temporal, causal, etc.).

Notation and definition of a process. [think whether all of the notation is needed]

A process $\mathcal{P} = (V, E)$ is a graph with labeled edges (maybe the sentence \mathbf{x} is also part of it?), where the nodes $V = \{1, \dots, |V|\}$ are event mentions and edges represent event-event relations. Given a paragraph $\mathbf{x} = \{x_1, \dots, x_{|\mathbf{x}|}\}$, we define $x_{i:j}$ to be a span of words $\{x_i, x_{i+1}, \dots, x_j\}$. An event mention v is defined by an event trigger t_v , which is some span of words $x_{i:j}$ and by a set of arguments A_v , where each argument $a_v \in A_v$ is again a span of words and a semantic role label a_l taken from a set \mathcal{L} . an event-event relation is a labeled edge (u, v, r) where $r \in \mathcal{R}$ is a closed set of relations¹.

For example, the first sentence of Figure 1 contains two event mentions, where $t_1 = \text{flowing}$ and $t_2 = \dots$. An example for the notation. [example that will make the notation complete and clear]

In general, one can think of process extraction as comprising of two steps - the first is standard semantic role labeling which comprises trigger identification, argument identification and role labeling, and the second is extracting a rich set of event-event relations. We focus on the second task - basically we

¹Our process annotation also contains coreference relations between arguments but we omit this since it is not very relevant.

get as input the set of event mention triggers \mathcal{T} and extract the E .

Next we will describe the full set of semantic role labels \mathcal{L} , and more importantly the set of event-event relations \mathcal{R} .

The set of semantic role labels \mathcal{L} contains standard labels such as *agent*, *theme*, *origin*, *destination* and *location*. In addition we have two semantic role labels that are relevant for the biological text domain - *result* and *raw-material*, which correspond to arguments that are the result of the event and materials used during the event. The last sentence in Figure 1 gives an example for these two labels.

Describe the set of relations \mathcal{R} : (a) NextEvent - directed relation (b) Causes - directed relation (c) Enables - directed relation (d) SuperEvent - directed relation (e) Cotemporal - undirected relation (f) SameEvent - event coreference. (g) None. Differences from previous work: (1) we have not only the temporal relation but also "cause" and "enable" that are important in our domain. (2) We have SuperEvent that most work on temporal ordering did not deal with (we should write who did deal with it) (3) We add coreference as just one of the event-event relations. This allows us to add constraints between cored and other relations in our global formulation. Write something about that we believe this is a succinct and good set of relation for process representation.

Properties of processes Naturally there are various properties to coherent processes (1) in the semantic role labeling part - different event triggers don't overlap, different argument of the same event don't overlap (this was used by Toutanova and Haghighi, 2006) [we have to think what to mention depending on whether we show results of the full pipeline] (2) in the event-event relations - all events are somehow linked to one another. Also in general the events are "chain-like" that is most events are related to one or two other events but not more than that (see Table that shows distribution of degrees of event mentions). In Section ... we will show that by using these global properties of processes and more we can improve performance. (3) Some triads are not possible etc. (4) more? [In general I think it is interesting and relevant to discuss here the global properties of processes even if we eventually don't use all of them in our final model - this is simply

since the local model does a good enough job but still identifying these properties is of interest in this paper].

Data set. We annotated 150 processes using the described scheme. Table ... describes the statistics of this data set. Stats in the data set: average number of tokens, average number of events mentions, average number of relations, more?

Next we will describe our method that given an annotation of event triggers \mathcal{T} extracts all the event-event relations. We first describe a local pairwise model that classifies each pair of events independently of others and then show our full model that uses global properties of structure.

References

- Quang Do, Wei Lu, and Dan Roth. 2012. Joint inference for event timeline construction. In *EMNLP-CoNLL*, pages 677–687.
- David McClosky and Christopher D. Manning. 2012. Learning constraints for consistent timeline extraction. In *EMNLP-CoNLL*, pages 873–882.
- David McClosky, Mihai Surdeanu, and Christopher D. Manning. 2011. Event extraction as dependency parsing. In *ACL*, pages 1626–1635.
- Sebastian Riedel and Andrew McCallum. 2011. Fast and robust joint models for biomedical event extraction. In *Proceedings of the Conference on Empirical methods in natural language processing (EMNLP '11)*.