

Discovering Content through Text Mining for a Synthetic Biology Knowledge System



Mai H. Nguyen
Gaurav Nakum
Jaiwei Tang
Xuanyu Wu



Bridget T. McInnes
Nicholas E. Rodriguez



Eric Young
Kevin Keating

- **Motivation**

- Scientific articles contain wealth of information about experimental methods and design results
- The number of scientific articles published is growing exponentially
- Automating process of extracting information
 - Extracting information from text is difficult due to ambiguity, variability, and volume

- **Approach**

- Use Named Entity Recognition (NER) to mine existing literature
- Goal of NER: Locate and classify entities in text into categories
 - e.g., genes, chemical, species

Species

Gene or Protein

Gene or Protein

Chemical

endogenous *S. cerevisiae* enzymes such as the reductase Oye2 and acetyltransferase Aft1 are known to degrade geraniol.

Synthetic Biology Knowledge System: Text Mining Pipeline

- Text mining Pipeline:
 - Parse: XML formatted full articles
 - Extract: Entities using our named entity recognition (NER) system
 - General entities (eg. gene and proteins, chemicals, cell line, and species)
 - Synthetic biology specific entities (eg. Promoters)
 - Validate: NER-discovered annotations are validated by domain experts
 - Refine: Fine tune NER system with validated annotations
 - Repeat: Continue until the system obtains a sufficient precision and recall

Entities identified are linked to the SBKS open knowledge system

```

graph TD
    Input1["'S.'"] --> BioBERT["BioBERT Contextualized Language Model"]
    Input2["'cerevisiae'"] --> BioBERT
    Input3["'enzymes'"] --> BioBERT
    BioBERT --> E1["'S' embedding"]
    BioBERT --> E2["'.' embedding"]
    BioBERT --> E3["'cere' embedding"]
    BioBERT --> E4["'visiae' embedding"]
    BioBERT --> E5["'enzymes' embedding"]
    E1 --> T1["TopModel"]
    E2 --> T2["TopModel"]
    E3 --> T3["TopModel"]
    E4 --> T4["TopModel"]
    E5 --> T5["TopModel"]
    T1 --> P1["'Species'"]
    T2 --> P2["'Species'"]
    T3 --> P3["'Species'"]
    T4 --> P4["'O'"]
    T5 --> P5["'O'"]
    P1 --> Output1["'Species'"]
    P2 --> Output2["'Species'"]
    P3 --> Output3["'Species'"]
    P4 --> Output4["'O'"]
    P5 --> Output5["'O'"]
  
```

[illegible]

Word cloud of Chemical entity mentions in ACS dataset

Annotations discovered by NER model in an ACS article