

# High-precision biological event extraction with a concept recognizer

K. Bretonnel Cohen\*, Karin Verspoor\*, Helen L. Johnson,  
Chris Roeder, Philip V. Ogren, William A. Baumgartner Jr.,  
Elizabeth White, Hannah Tipney, and Lawrence Hunter

**U. Colorado Denver School of Medicine**

\*These two authors contributed equally to the work reported here

# Introduction

- All three tasks approached as concept recognition and analysis using the OpenDMAP semantic parser
- Manually-written rules
- Achieved highest precision for two of the three tasks; recall was low

# Methods

- Named entity recognition
  - Proteins: provided by organizers
  - Other semantic classes:
    - LingPipe with GENIA model
    - ConceptMapper (Gene Ontology cellular components, Cell Type Ontology, BRENDA Tissue Ontology, and Sequence Ontology)
      - Only modifications: remove *cell* from CTO and add synonym *nuclear* to GOCC

# Methods

- Coördination
  - Retrained OpenNLP constituent parser with...
    - 500 abstracts from beta version of GENIA treebank
    - 10 full-text articles from the CRAFT corpus
  - Distributed meaning assumed
  - Used for proteins only

# Methods

- OpenDMAP (Hunter et al. 2008)
  - Uses ontology as central organizing structure
    - Commitment to using community-consensus ontologies (mostly OBO)
  - Semantic grammars allow mixing of terminals and semantically typed non-terminals
  - Context-free power with variable ordering
  - Slot fillers constrained by ontology

# Methods

Protein\_transport :=

[TRANSPORTED-ENTITY] translocation

@(from {DET}? [TRANSPORT-ORIGIN])

@(to {DET}? [TRANSPORT-DESTINATION])

*Bax translocation to mitochondria from the cytosol*

*Bax translocation from the cytosol to the mitochondria*

# Methods

Protein\_transport :=

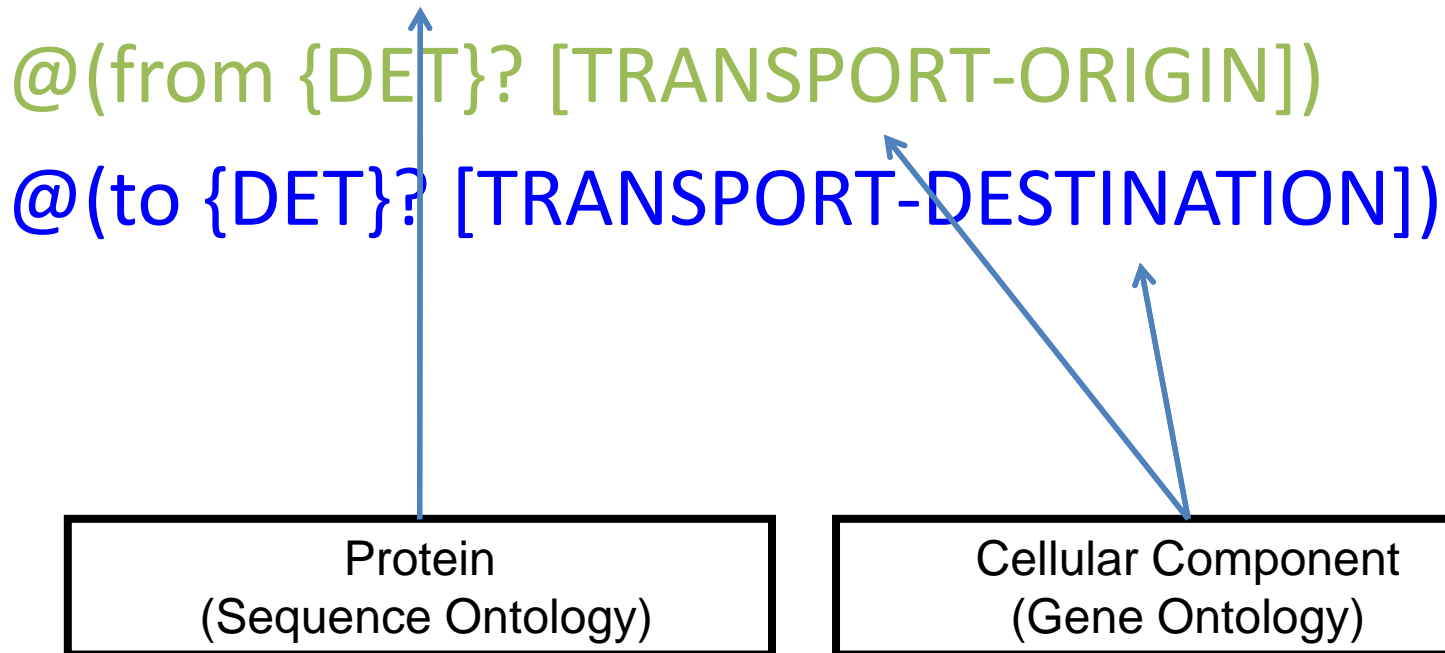
[TRANSPORTED-ENTITY] translocation

@(from {DET}? [TRANSPORT-ORIGIN])

@(to {DET}? [TRANSPORT-DESTINATION])

Protein  
(Sequence Ontology)

Cellular Component  
(Gene Ontology)



# Methods

- All event types represented as frames
  - Elements from ontology constrain every slot

EVENT TYPE: REGULATION

AtLoc: instance of biological\_entity

Cause: instance of protein

CSite: instance of biological\_concept or polypeptide\_region

Event\_action: instance of trigger\_word or detection\_method

Site: instance of biological\_concept or polypeptide\_region

Theme: instance of protein or biological\_process

ToLoc: instance of biological\_entity



# Methods

The screenshot displays a software interface with two main panels: 'CLASS BROWSER' and 'CLASS EDITOR'. The 'CLASS BROWSER' panel shows a class hierarchy for the project 'bionlp09-ontology', with 'regulation' selected under 'biological process'. The 'CLASS EDITOR' panel shows details for the 'regulation' class, including its name, role (Concrete), and a table of template slots.

**CLASS BROWSER**  
For Project: ● bionlp09-ontology

Class Hierarchy

- :THING
  - ▶ ● :SYSTEM-CLASS
  - ▶ ● entrez\_record
  - ▼ ● biological concept
    - ▶ ● biological\_entity
    - ▼ ● biological process
      - ▶ ● regulation
      - gene\_expression
      - transcription
      - protein\_catabolism
      - localization
      - ▶ ● binding
        - phosphorylation
      - detection\_method
    - ▶ ● linguistic concept

**CLASS EDITOR**  
For Class: ● regulation (instance of :STANDARD-CLASS)

Name: regulation

Role: Concrete ●

Documentation: Any process that modu process. Biological proc include the control of g with a protein or substr

Template Slots

| Name         | Cardinality     | Type   |
|--------------|-----------------|--|
| AtLoc        | single          | Instance of biological_entity                                    |
| Cause        | single          | Instance of protein or protein_conjunction                       |
| CSite        | single          | Instance of biological_concept or polypeptide_region             |
| event_action | single          | Instance of trigger_words or detection_method                    |
| Site         | single          | Instance of biological_entity or polypeptide_region              |
| Theme        | required single | Instance of protein_conjunction or biological process or protein |
| ToLoc        | single          | Instance of biological_entity                                    |

# Methods

- All event types represented as frames
  - Elements from ontology constrain every slot

EVENT TYPE: REGULATION

AtLoc: instance of biological\_entity

Cause: instance of protein

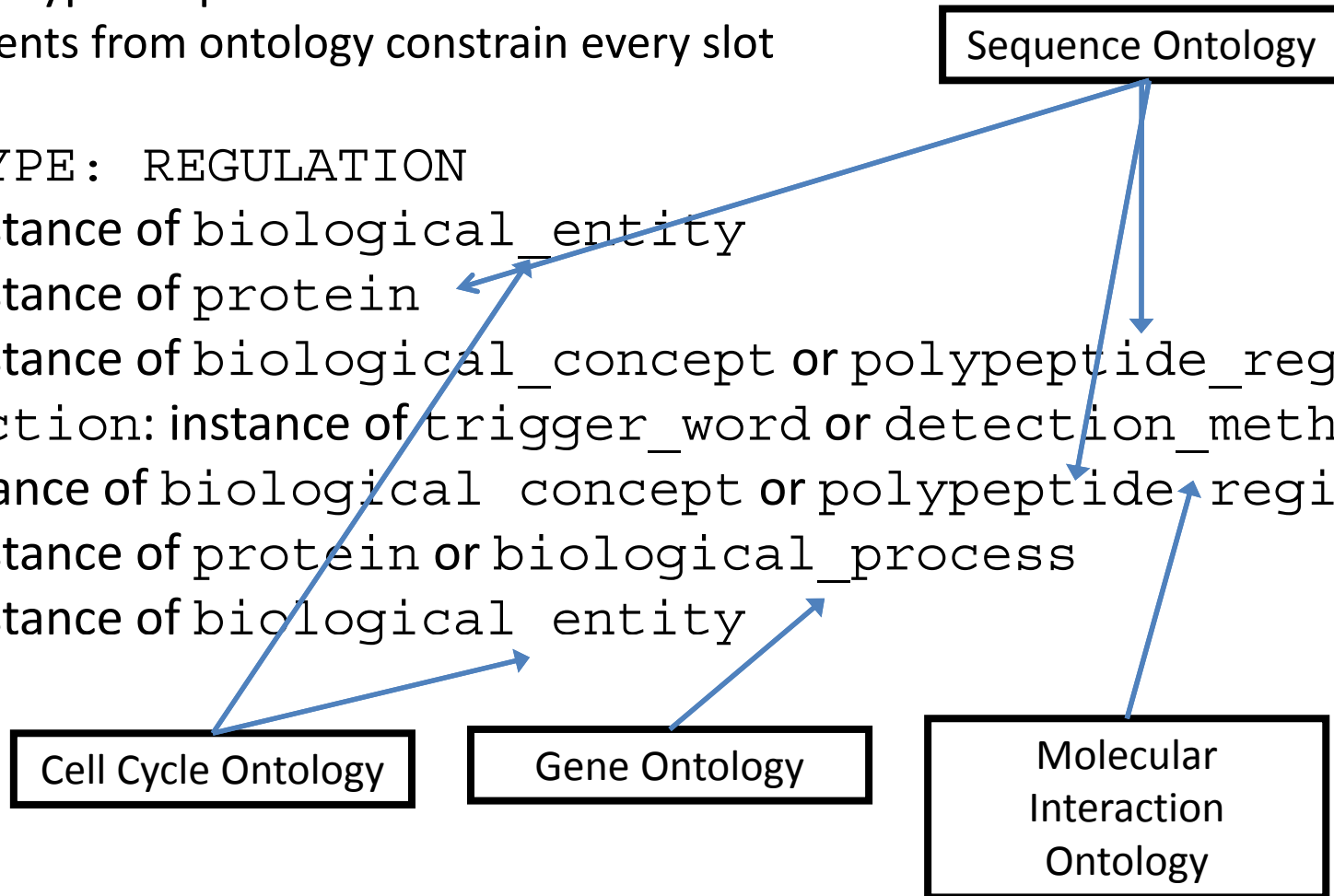
CSite: instance of biological\_concept or polypeptide\_region

Event\_action: instance of trigger\_word or detection\_method

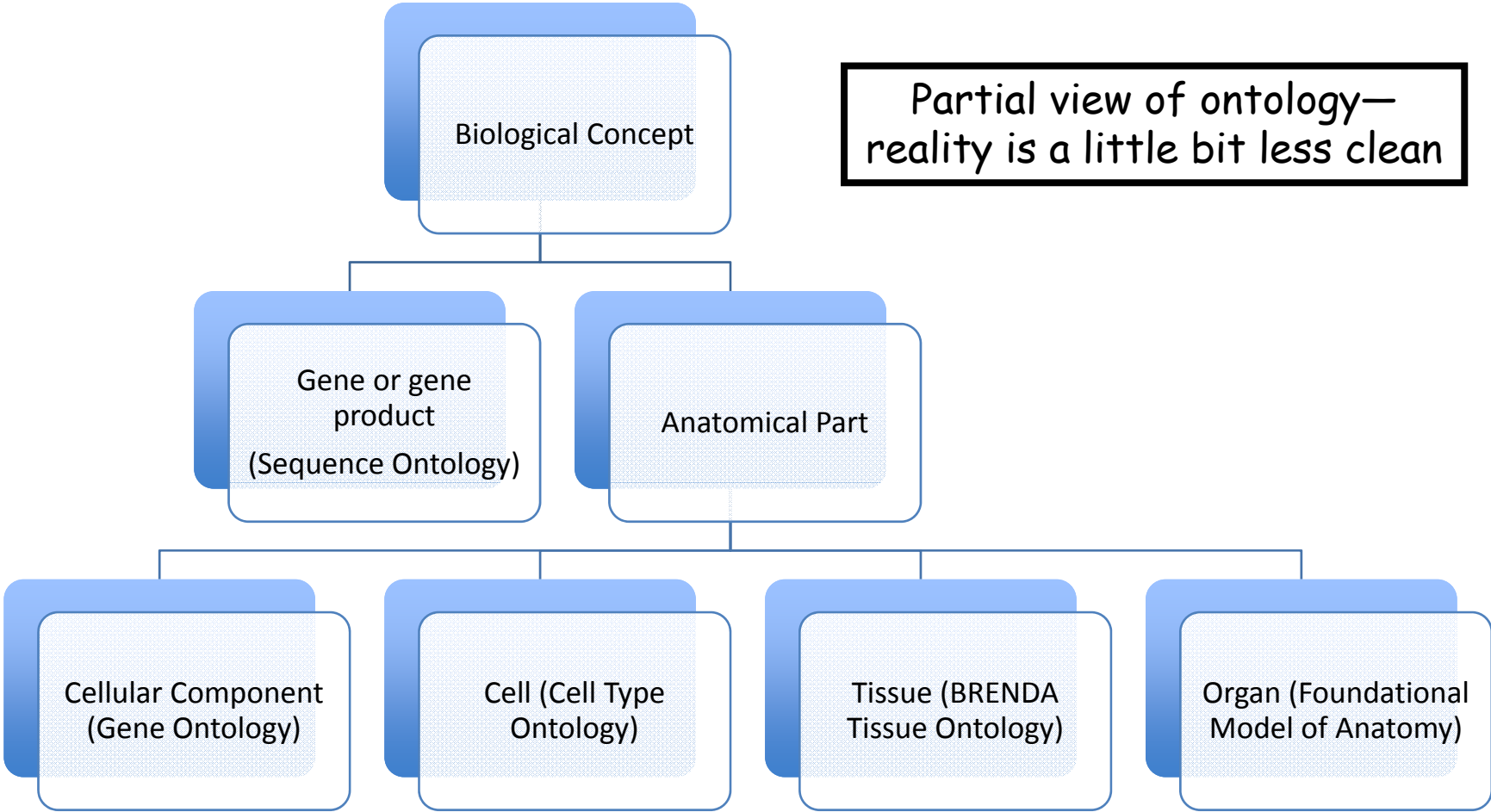
Site: instance of biological\_concept or polypeptide\_region

Theme: instance of protein or biological\_process

ToLoc: instance of biological\_entity



# Methods



# Methods

| Event type         | Site  | AtLoc  | ToLoc                   |
|--------------------|---|--|-------------------------|
| Binding            | protein domain (SO), binding site (SO), DNA (SO), chromosome (SO) |  |                         |
| Gene expression    | gene (SO), biological entity (CCO)                                | tissue (BTO), cell type (CTO), cellular component (GO) |                         |
| Localization       |   | cellular component (GO)                                | cellular component (GO) |
| Phosphorylation    | amino acid (FMA), polypeptide region (SO)                         |  |                         |
| Protein catabolism | cellular component (GO)   |  |                         |
| Transcription      | gene (SO), biological entity (CCO)                                |  |                         |

BTO: BRENDA Tissue Ontology  
 CCO: Cell Cycle Ontology  
 CTO: Cell Type Ontology  
 GO: Gene Ontology  
 SO: Sequence Ontology

# Methods

- Manual pattern-writing
  - Before availability of training data: based on native speaker intuitions, examples from PubMed, and variations on same, as in Cohen et al. (2004)
  - After release of training data: based on examination of corpus data, targeting high-frequency predicates only
  - Nominalizations predominated; used insights from Cohen et al. (2008) regarding Theme placement
  - Protein binding rules re-used from BioCreative II protein-protein interaction task
  - Eschewed use of wildcards

# Results

|        | Our system   |       |       | Best team    |              |              | Best P/R/F   |              |              |
|--------|--------------|-------|-------|--------------|--------------|--------------|--------------|--------------|--------------|
|        | P            | R     | F     | P            | R            | F            | P            | R            | F            |
| Task 1 | <b>71.81</b> | 13.45 | 22.66 | 58.48        | <b>46.73</b> | <b>51.95</b> | <b>71.81</b> | <b>46.73</b> | <b>51.95</b> |
| Task 2 | <b>70.97</b> | 13.25 | 43.12 | 54.08        | <b>35.86</b> | <b>43.12</b> | <b>70.97</b> | <b>35.86</b> | <b>43.12</b> |
| Task 3 | 57.40        | 12.33 | 20.30 | <b>60.83</b> | <b>32.68</b> | <b>42.52</b> | <b>60.83</b> | <b>32.68</b> | <b>42.52</b> |

Task 1: P 10 points higher than second-highest

Task 2: P 14 points higher than second-highest

Task 3: P 3.4 points lower than highest (3/6)

# Results

## Unofficial results: contribution of bug repairs

|                  | P     | R     | F     |
|------------------|-------|-------|-------|
| Official results | 71.81 | 13.45 | 22.66 |
| With bug fixes   | 67.19 | 17.38 | 27.10 |

Still the highest precision  
(#2 was 62.21)

# Results

- Contribution of coördination-handling
  - Bug-fixed results: F 27.62 (Task 1)
  - Without coördination-handling: F 24.72
  - Decrease in F of 2.9 without coördination-handling



# Results

- Error analysis: false negatives
  - Intervening material between trigger words and arguments
  - Coördination of things other than protein names
  - Low coverage of trigger words
  - Anaphora and coreference
  - Appositive gene names and symbols

# Results

- Intervening material between trigger words and arguments

*to efficiently [express] in developing thymocytes  
a mutant form of the [NF-kappa B inhibitor]  
(PMID 10092801)*

- Solutions: use dependency syntax, multiword expressions

# Results

- Coördination of things other than protein names

*[transcription] and subsequent synthesis and secretion  
of [galectin-3] (PMID 8623933)*

# Results

- Low coverage of trigger words
  - Sharing of some trigger words across multiple event types
  - Time constraints of the shared task
- Similarly, limited attempts to deal with intervening material

# Results

- Anaphora and coreference

*Although 2 early lytic transcripts, [BZLF1] and [BHRF1], were also detected in 13 and 10 cases, respectively, the lack of ZEBRA staining in any case indicates that these lytic transcripts are most likely [expressed] by rare cells in the biopsies entering lytic cycle (PMID 8903467)*

# Results

- Appositive gene names and symbols

Rule: {gene\_expression} :=  
expression of {det}? [Theme] ;

*[expression] of **Fas ligand** ([FasL]) (PMID 10092076)*

Intervening  
material

# False Positive Analysis

- Training data
- Biologist examined all events scored as FP in 10 documents per event type
- 42% judged as actually TP

|                     | Analyzed | TP | FP |
|---------------------|----------|----|----|
| Gene Expression     | 11       | 7  | 4  |
| Binding             | 21       | 15 | 6  |
| Transcription       | 4        | 1  | 3  |
| Protein Catabolism  | 1        | 1  | 0  |
| Localization        | 4        | 1  | 3  |
| Phosphorylation     | 3        | 2  | 1  |
| Regulation          | 13       | 1  | 12 |
| Positive Regulation | 18       | 7  | 11 |
| Negative Regulation | 13       | 2  | 11 |

# False Positive Analysis

|   |
|---|
| Key:<br><u>trigger word</u><br><b>Theme</b> |
|---|

- Incorrect Theme chosen:
  - “IL-10 production by **gp41**”<sup>[10089566]</sup>
  - “induction of **I kappa B alpha** phosphorylation”<sup>[7499266]</sup>
- Missed a required entity:
  - “induction of **IL-10 production** by gp41”<sup>[10089566]</sup>
- Misrecognized trigger word:
  - “upstream of the **GM-CSF** transcription initiation site”<sup>[7478534]</sup>
  - “effects of **IL-11** were associated with reduced **NF-kappaB** activation”<sup>[10411003]</sup>
  - “up-regulation of **CD80 Ag**”<sup>[8690900]</sup>
- Sentence parse error:
  - “was suppressed by **alpha B2**. Coexpression of alpha B1”<sup>[7605990]</sup>



# Discussion and Conclusions

- Recall is a function of the rule set, not of the approach
- Can this high-precision, apparently low-recall approach scale to practical performance levels?
  - Redundancy in individual papers
  - These results make no use of syntactic analysis
  - Exploring rapid adaptation to new tasks
    - Rule/template inheritance
    - Extend coördination handling beyond proteins
    - Using dependency parses
    - Mappings from frames to WordNet for better trigger word coverage

# Discussion and Conclusions

- Not just building system for shared tasks—  
input to other tasks, such as high-throughput  
data analysis, which it has proven useful for

# Acknowledgements

- Michael Bada for help with loading the Sequence Ontology into Protégé
- Alias-I for LingPipe
- Anonymous ConceptMapper creator
- Shared task organizers for organizing the task and answering many questions
- Work supported by:
  - NIH grants R01LM009254, R01GM083649, and R01LM008111 to Larry Hunter
  - NIH grant T15LM009451 to Philip Ogren

# Availability

- OpenDMAP semantic parser and rule set:

`bionlp.sourceforge.net`

# Questions to be prepared for

- What does it mean to have an ontology be the central organizing structure of a semantic parser?
- Why didn't we get the top precision on Task 3?
- You didn't say anything about how you tackled negation and speculation

# Ontology as the central organizing principle

Ontology specifies the targets of information extraction

The screenshot displays a software interface with two main panels: 'CLASS BROWSER' on the left and 'CLASS EDITOR' on the right. The 'CLASS BROWSER' shows a hierarchical tree of classes under the project 'bionlp09-ontology'. The 'CLASS EDITOR' shows details for the class 'negative\_regulation', including its name, role, and a table of template slots.

**CLASS BROWSER**  
For Project: ● bionlp09-ontology

**Class Hierarchy**

- :THING
  - ▶ ● :SYSTEM-CLASS
  - ▶ ● entrez\_record
  - ▼ ● biological concept
    - ▶ ● biological\_entity
    - ▼ ● biological process
      - ▼ ● regulation
        - positive\_regulation
        - negative\_regulation
      - gene\_expression
      - transcription
      - protein\_catabolism
      - localization
      - ▶ ● binding
      - phosphorylation
      - detection\_method

CLASS EDITOR

For Class: ● negative\_regulation (instance of STANDARD-CLASS-WITH-PATTERNS)

Name

negative\_regulation

Role

Concrete ●

Template Slots

| Name         | Cardinality     | Type  |
|--------------|-----------------|---|
| AtLoc        | single          | Instance of biological_entity                           |
| Cause        | single          | Instance of protein or protein_conjunction              |
| CSite        | single          | Instance of biological concept or polypeptide_region    |
| event_action | single          | Instance of trigger_words or detection_method           |
| Site         | single          | Instance of biological_entity or polypeptide_region     |
| Theme        | required single | Instance of protein_conjunction or biological proces... |
| ToLoc        | single          | Instance of biological_entity                           |

Documentation

Any process that stops, prevents or reduces the frequency, rate or extent of a biological process. Biological processes are regulated by many means; examples include the control of gene expression, protein modification or interaction with a protein or substrate molecule. (GO:0048519)

# Ontology as the central organizing principle

The screenshot displays an ontology editor interface with two main panels: a Class Browser on the left and a Class Editor on the right.

**Class Browser:** Shows a class hierarchy for the project 'bionlp09-ontology'. The hierarchy is as follows:

- :THING
  - :SYSTEM-CLASS
  - entrez\_record
  - biological concept
    - biological\_entity
    - biological process
      - regulation** (highlighted)
      - gene\_expression
      - transcription
      - protein\_catabolism
      - localization
      - binding
      - phosphorylation
      - detection\_method
    - linguistic concept

**Class Editor:** Shows details for the class 'regulation' (instance of :STANDARD-CLASS).

- Name:** regulation
- Role:** Concrete
- Documentation:** Any process that modu process. Biological proc include the control of g with a protein or substr
- Template Slots:**

| Name         | Cardinality     | Type   |
|--------------|-----------------|--|
| AtLoc        | single          | Instance of biological_entity                                    |
| Cause        | single          | Instance of protein or protein_conjunction                       |
| CSite        | single          | Instance of biological concept or polypeptide_region             |
| event_action | single          | Instance of trigger_words or detection_method                    |
| Site         | single          | Instance of biological_entity or polypeptide_region              |
| Theme        | required single | Instance of protein_conjunction or biological process or protein |
| ToLoc        | single          | Instance of biological_entity                                    |

# Ontology as the central organizing principle

Ontology constrains slot fillers in rules

The screenshot displays a software interface with two main panels: 'CLASS BROWSER' and 'CLASS EDITOR'.

**CLASS BROWSER:** Shows a class hierarchy for the project 'bionlp09-ontology'. The hierarchy starts with ':THING' and includes several levels of biological entities, such as ':SYSTEM-CLASS', 'entrez\_record', 'biological concept', 'biological\_entity', 'gene or gene product', 'gene', 'macromolecule', 'protein', 'nucleic acid', 'DNA', 'RNA', 'anatomical part', 'cellular component', and 'cell surface'. A blue arrow points from the text 'Ontology constrains slot fillers in rules' to the 'gene or gene product' class in this hierarchy.

**CLASS EDITOR:** Shows details for the class 'regulation' (instance of ':STANDARD-CLASS'). The 'Name' field contains 'regulation'. The 'Role' is set to 'Concrete'. A blue arrow points from the text to the 'Role' dropdown menu. Below this is a 'Template Slots' table.

| Name         | Cardinality     | Type   |
|--------------|-----------------|--|
| AtLoc        | single          | Instance of biological_entity                                    |
| Cause        | single          | Instance of protein or protein_conjunction                       |
| CSite        | single          | Instance of biological concept or polypeptide_region             |
| event_action | single          | Instance of trigger_words or detection_method                    |
| Site         | single          | Instance of biological_entity or polypeptide_region              |
| Theme        | required single | Instance of protein_conjunction or biological process or protein |
| ToLoc        | single          | Instance of biological_entity                                    |

The 'Documentation' field on the right side of the Class Editor contains the text: 'Any process that modulate biological process. Biological processes include the control of gene expression with a protein or substrate.'



# Methods

- Negation and speculation handled by string literals
- Negation: *absence of, failure to, not KEYWORD*
- Speculation: *research, study, examine, investigate, suggests, unknown, seems*
- Heavy use of wildcards