

Proposal: Enhancing SBOannotator with LLM Integration & Dynamic Term Retrieval

Personal Information

Name: Jiahui Hu

Email: hu.jiahu@northeastern.edu

Alternative Email: lareinahuwork@gmail.com

GitHub: <https://github.com/lareinahu-2023>

University: Northeastern University, Master's in Computer Software Engineering

Prior Contribution: <https://github.com/draeger-lab/SBOannotator/pull/2>

Project Synopsis

This proposal aims to transform SBOannotator from a static, hard-coded tool into a dynamic, intelligent system for annotating SBML models with Systems Biology Ontology terms. After initial discussions with mentor [Nantia Leonidou](#), who provided valuable feedback on my approach, I've developed a comprehensive plan and initial implementation (PR #2 <https://github.com/draeger-lab/SBOannotator/pull/2>).

By implementing real-time SBO term retrieval, integrating multiple enzymatic data sources, and incorporating LLM-based annotation assistance, this project will significantly enhance both the accuracy and usability of SBOannotator while maintaining its core classification strengths. The addition of a standalone desktop GUI with interactive visualizations will make these powerful annotation capabilities accessible to a broader range of systems biology researchers.

Project Link: <https://github.com/nrnbc/GoogleSummerOfCode/issues/261>

Potential Mentor: Nantia Leonidou - nantia.leonidou@dkfz-heidelberg.de

Project Description

Background

SBOannotator is the first standalone tool that automatically assigns Systems Biology Ontology (SBO) terms to multiple entities of a given SBML (Systems Biology Markup Language) model. Its main strength lies in annotating biochemical reactions within metabolic models, as the correct assignment of precise SBO annotations requires extensive classification. SBO terms play a crucial role in precisely defining the functions of various components within biological models. By assigning these terms, the interpretability and overall understanding of the model are significantly improved.

However, in large-scale models containing thousands of reactions and chemical species, manually assigning SBO terms becomes a highly complex and time-consuming task. Currently, SBOannotator assigns terms that are hard-coded, making it unable to integrate newly introduced ontologies. Additionally, it derives enzymatic information needed for the correct classification of reactions from the BiGG database or a predefined static SQL database.

Problem Statement

The current SBOannotator has three key limitations:

1. **Static SBO Terms:** Hard-coded terms prevent the integration of newly introduced ontologies from the OLS server.
2. **Limited Data Sources:** Reliance on BiGG database restricts access to comprehensive enzymatic information.
3. **Complex User Experience:** The lack of an intuitive interface and visualization tools limits accessibility.

Project Goals

1. Develop a dynamic system to automatically fetch and update SBO terms from the OLS server
2. Expand enzymatic data integration to include KEGG and BRENDA databases
3. Implement an on-premise LLM-based annotation assistant to intelligently suggest SBO terms
4. Create interactive visualizations for reaction networks and annotation statistics
5. Build a standalone desktop GUI with an intuitive user interface

Benefits to the Community

This project will benefit the systems biology community by:

- **Improving Annotation Accuracy:** Access to up-to-date SBO terms and multiple enzymatic data sources will enable more precise annotations.
- **Enhancing Efficiency:** LLM-based suggestions and an intuitive GUI will significantly reduce the time required for annotation.
- **Expanding Accessibility:** The standalone application will make advanced annotation capabilities available to researchers with varying technical expertise.
- **Facilitating Model Comparison:** Visualization tools will allow for better understanding and comparison of model components.

- **Ensuring Sustainability:** Dynamic term retrieval will future-proof the tool against ontology changes.

Initial Implementation Progress

1. Description of What I Have Accomplished:

I've already laid the groundwork for integrating LLM capabilities into SBOannotator by creating the essential foundation components, as demonstrated in my PR #2:

<https://github.com/draeger-lab/SBOannotator/pull/2>.

My initial implementation includes:

1. Comprehensive Implementation Plan (LLM_ANNOTATION_PLAN.md):

- Defined a clear five-phase implementation strategy
- Created detailed architectural diagram showing component relationships
- Outlined specific technical approaches for reaction feature extraction
- Developed prompt engineering strategies for biochemical classification
- Established evaluation metrics for assessing annotation quality

2. Modular LLM Provider Interface (llm_interface.py):

- Built abstract base class for provider-agnostic LLM integration
- Implemented provider-specific classes for OpenAI and Anthropic
- Created reaction feature extraction framework with placeholder methods
- Designed SBOAnnotationAssistant class to orchestrate the annotation process
- Developed methods for batch processing entire SBML models

3. Template Management System (template_manager.py):

- Created a Jinja2-based template engine for dynamic prompt generation
- Implemented reaction-specific template selection logic
- Developed four specialized templates for different reaction types:
 - General biochemical reactions
 - Transport reactions
 - Enzymatic reactions
 - Exchange/boundary reactions
- Built template caching system for improved performance

These components establish a solid architectural foundation for the LLM integration, following software engineering best practices including:

- Clear separation of concerns
- Extensible interfaces
- Comprehensive documentation
- Type hints throughout the codebase

2. ScreenShots of My Key Files:

1. LLM_ANNOTATION_PLAN.md

```
1  # LLM-Based SBO Term Annotation: Implementation Plan
2
3  > h2: ## Overview
4
5
6
7  ## Project Objectives
8
9
10 1. Create a modular, extensible interface for interacting with different LLM providers
11 2. Develop intelligent reaction feature extraction for improved term suggestions
12 3. Design effective prompting strategies for accurate SBO term assignment
13 4. Integrate the LLM assistant with the existing SBOannotator workflow
14 5. Evaluate and benchmark the accuracy of LLM-based annotations against existing methods
```

```
14
15 ## Architecture
16
17 The solution follows a modular design with clear separation of concerns:
18
19 ...
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
```

```

49  ## Implementation Phases
50
51  ### Phase 1: Foundation (Current PR)
52
53  - [x] Design the abstract LLM provider interface
54  - [x] Create the SB0AnnotationAssistant class with core functionality
55  - [x] Define placeholder methods for reaction feature extraction and prompt generation
56  - [x] Outline initial provider implementations for OpenAI and Anthropic
57
58  ### Phase 2: Feature Extraction
59
60  - [ ] Implement comprehensive reaction feature extraction:
61    - [ ] Identify metabolite patterns (e.g., ATP/ADP for phosphorylation)
62    - [ ] Detect compartment changes for transport reactions
63    - [ ] Extract EC numbers from reaction annotations
64    - [ ] Analyze reaction reversibility and stoichiometry
65    - [ ] Process reaction names for keyword indicators
66
67  ### Phase 3: Prompt Engineering
68
69  - [ ] Develop effective prompting strategies:
70    - [ ] Research optimal prompts for biochemical reaction classification
71    - [ ] Include relevant SB0 term definitions in context
72    - [ ] Structure prompts to encourage specific, accurate responses
73    - [ ] Optimize for token efficiency
74    - [ ] Format instructions for structured output
75
76  ### Phase 4: Response Processing
77
78  - [ ] Create robust response parsing:
79    - [ ] Extract recommended SB0 terms from LLM output
80    - [ ] Handle various response formats gracefully
81    - [ ] Parse confidence scores and explanations
82    - [ ] Implement fallback logic for ambiguous responses
83

```

```

84  ### Phase 5: Integration & Validation
85
86  - [ ] Integrate with SB0annotator:
87    - [ ] Add LLM-based annotation as optional enhancement
88    - [ ] Implement configuration options for LLM settings
89    - [ ] Create hybrid approach that combines rule-based and LLM methods
90    - [ ] Develop validation metrics to compare LLM suggestions against existing annotations
91    - [ ] Implement benchmark tooling to evaluate different LLM models and prompts
92
93  ## Technical Details
94
95  ### Reaction Feature Extraction
96
97  The quality of LLM suggestions depends heavily on effective feature extraction. Key features include:
98
99  - Metabolite Patterns: Identifying characteristic metabolites (ATP/ADP, NAD/NADH, etc.)
100 - Compartment Analysis: Detecting cross-compartment transport
101 - EC Number Integration: Using existing EC annotations to inform suggestions
102 - Reaction Properties: Analyzing reversibility, stoichiometry, and other properties
103 - Naming Analysis: Extracting insights from reaction names (e.g., "kinase" indicating phosphorylation)
104
105  ### Prompt Engineering
106
107  Effective prompts must:
108
109  1. Provide sufficient context about SB0 terms and their meanings
110  2. Present reaction details in a structured format
111  3. Guide the LLM toward precise ontological classification
112  4. Specify response format for easier parsing
113  5. Include few-shot examples of correct classifications
114

```

```

115 Example prompt structure:
116 ```
117 You are a Systems Biology expert tasked with assigning SBO terms to biochemical reactions.
118
119 Reaction details:
120 - ID: {reaction_id}
121 - Name: {reaction_name}
122 - Reactants: {reactants}
123 - Products: {products}
124 - [Additional extracted features]
125
126 Relevant SBO terms include:
127 - SBO:0000176 (Biochemical reaction): General biochemical transformation
128 - SBO:0000200 (Redox reaction): Involves electron transfer (e.g., NAD/NADH)
129 - [Additional relevant terms]
130
131 Based on these details, what is the most appropriate SBO term for this reaction?
132 Provide your answer in JSON format with fields: sbo_term, confidence, explanation.
133 ```
134
135 ### Provider Implementation
136
137 For each LLM provider, we need to:
138
139 1. Implement authentication and API handling
140 2. Optimize request parameters (temperature, max tokens, etc.)
141 3. Manage rate limiting and error handling
142 4. Process provider-specific response formats
143
144 ## Future Enhancements
145
146 1. Fine-tuning: Create specialized models fine-tuned on SBO classification
147 2. Batch Processing: Optimize for efficient batch annotation
148 3. Confidence Thresholds: Implement adaptive confidence thresholds
149 4. User Feedback Loop: Incorporate user corrections to improve suggestions
150 5. Term Expansion: Expand beyond reactions to other model components
151
152 ## Evaluation
153
154 The LLM-based annotation will be evaluated on:
155 💡
156 1. Accuracy: Agreement with expert-annotated models
157 2. Coverage: Percentage of reactions receiving specific (non-generic) SBO terms
158 3. Consistency: Consistency of annotations across similar reactions
159 4. Performance: Annotation speed and resource usage
160 5. Explainability: Quality of explanations for suggested terms
161
162 ## Conclusion
163
164 This LLM-based annotation assistant represents a significant enhancement to SBOannotator, enabling more accurate and
165 comprehensive annotation of SBML models. The modular design ensures adaptability to different LLM providers
166 and future ontological developments.
167
168 The current PR establishes the foundation for this work, demonstrating an understanding of
169 both the SBOannotator codebase and the requirements for effective LLM integration.

```

2. llm_interface.py - Next show the core interface implementation

```

1  """LLM Interface for SBO term annotation assistance.
2
3  This module provides an abstract interface for using Large Language Models (LLMs)
4  to assist in the annotation of SBML models with appropriate SBO terms.
5  """
6
7
8
9  import abc
10 from typing import Dict, List, Tuple, Optional, Union, Any
11 from libsbml import SBMLDocument, Model, Reaction
12
13

```

```

14 @ class LLMProvider(abc.ABC): 3 usages 1 lareinahu-2023
15     """Abstract base class for LLM provider implementations."""
16
17     @abc.abstractmethod 1 lareinahu-2023
18     def initialize(self, **kwargs) -> None:
19         """Initialize the LLM provider with necessary credentials and settings.
20
21         Args:
22             **kwargs: Provider-specific initialization parameters
23         """
24         pass
25
26     @abc.abstractmethod 1 usage 1 lareinahu-2023
27     def generate_completion(self, prompt: str, **kwargs) -> str:
28         """Generate a completion from the LLM based on the given prompt.
29
30         Args:
31             prompt: The input prompt for the LLM
32             **kwargs: Additional provider-specific parameters
33
34         Returns:
35             The generated completion text
36         """
37         pass
38
39     @property 1 lareinahu-2023
40     @abc.abstractmethod
41     def provider_name(self) -> str:
42         """Return the name of the LLM provider."""
43         pass
44
45

```

```

46 class SBOAnnotationAssistant: 1 lareinahu-2023
47     """High-level interface for LLM-based SBO term annotation assistance."""
48
49     def __init__(self, llm_provider: LLMProvider): 1 lareinahu-2023
50         """Initialize the annotation assistant with an LLM provider.
51
52         Args:
53             llm_provider: An implementation of the LLMProvider interface
54         """
55         self.llm_provider = llm_provider
56         self.sbo_cache = {} # Cache for SBO terms info
57
58     def extract_reaction_features(self, reaction: Reaction) -> Dict[str, Any]: 1 usage 1 lareinahu-2023
59         """Extract relevant features from a reaction for annotation.
60
61         Args:
62             reaction: The SBML reaction to extract features from
63
64         Returns:
65             Dictionary of extracted features
66         """
67         # TODO: Implement feature extraction from reaction
68         # This should include:
69         # - Reaction ID and name
70         # - Reactants and products (including compartments)
71         # - Presence of specific metabolites (ATP, NAD, etc.)
72         # - Compartment changes (for transport reactions)
73         # - EC numbers from annotations
74         # - Other relevant metadata
75
76         features = {
77             "id": reaction.getId(),
78             "name": reaction.getName() or reaction.getId(),
79             "reversible": reaction.getReversible(),
80             # Add more features
81         }

```

```

82
83     return features
84
85     def format_prompt(self, reaction_features: Dict[str, Any]) -> str: 1 usage 1 lareinahu-2023
86         """Format a prompt for the LLM based on reaction features.
87
88         Args:
89             reaction_features: Dictionary of reaction features
90
91         Returns:
92             Formatted prompt string
93         """
94         # TODO: Implement prompt formatting
95         # This should create a detailed prompt that:
96         # - Provides context about SBO and its importance
97         # - Describes the reaction with relevant features
98         # - Asks for specific SBO term recommendations
99         # - Specifies response format
100
101         prompt = f"""
102         [Prompt would be constructed based on reaction_features]
103         """
104
105     return prompt
106

```



```

106
107 def parse_llm_response(self, response: str) -> Dict[str, Any]: 1 usage 1 lareinahu-2023
108     """Parse the LLM response into structured annotation suggestions.
109
110     Args:
111         response: Raw response from the LLM
112
113     Returns:
114         Dictionary containing structured annotation suggestions
115     """
116     # TODO: Implement response parsing
117     # This should extract:
118     # - Recommended SBO term ID(s)
119     # - Confidence scores
120     # - Explanations/reasoning
121     # - Alternative suggestions
122
123     result = {
124         "sbo_term": "SBO:0000000", # Placeholder
125         "confidence": 0.0,
126         "alternatives": [],
127         "explanation": ""
128     }
129
130     return result
131

```

```

132 def suggest_sbo_term(self, reaction: Reaction) -> Dict[str, Any]: 1 usage 1 lareinahu-2023
133     """Suggest an SBO term for a reaction using the LLM.
134
135     Args:
136         reaction: The SBML reaction to annotate
137
138     Returns:
139         Dictionary with SBO term suggestion and metadata
140     """
141     # Extract features
142     features = self.extract_reaction_features(reaction)
143
144     # Format prompt
145     prompt = self.format_prompt(features)
146
147     # Generate completion
148     response = self.llm_provider.generate_completion(prompt)
149
150     # Parse response
151     result = self.parse_llm_response(response)
152
153     return result
154

```

```

155 def batch_process_model(self, model: Model) -> Dict[str, Dict[str, Any]]: 1larelnahu-2023
156     """Process all reactions in a model and suggest SBO terms.
157
158     Args:
159         model: The SBML model to process
160
161     Returns:
162         Dictionary mapping reaction IDs to annotation suggestions
163     """
164     results = {}
165
166     for i in range(model.getNumReactions()):
167         reaction = model.getReaction(i)
168         suggestion = self.suggest_sbo_term(reaction)
169         results[reaction.getId()] = suggestion
170
171     return results
172
173 def validate_suggestion(self, sbo_term: str) -> bool: 1larelnahu-2023
174     """Validate that an SBO term ID is valid.
175
176     Args:
177         sbo_term: SBO term ID to validate
178
179     Returns:
180         Boolean indicating if the term is valid
181     """
182     # TODO: Implement validation against SBO ontology
183     return sbo_term.startswith("SBO:")
184
185

```

```

186 # Example provider implementations (to be implemented in separate modules)
187
188 class OpenAIProvider(LLMProvider):  # lareinahu-2023
189     """OpenAI LLM provider implementation."""
190
191     def initialize(self, api_key: Optional[str] = None, model: str = "gpt-3.5-turbo", **kwargs) -> None:  # lareinahu-2023
192         """Initialize OpenAI provider.
193
194         Args:
195             api_key: OpenAI API key
196             model: Model name to use
197             **kwargs: Additional parameters
198         """
199         # TODO: Implement OpenAI provider initialization
200         self.api_key = api_key
201         self.model = model
202
203     def generate_completion(self, prompt: str, **kwargs) -> str:  # lareinahu-2023
204         """Generate completion using OpenAI API.
205
206         Args:
207             prompt: The input prompt
208             **kwargs: Additional parameters
209
210         Returns:
211             Generated completion text
212         """
213         # TODO: Implement OpenAI API call
214         return "[OpenAI response placeholder]"
215
216     @property  # lareinahu-2023
217     def provider_name(self) -> str:
218         """Return the provider name."""
219         return "OpenAI"
220
221

```

```

222 class AnthropicProvider(LLMProvider):  # lareinahu-2023
223     """Anthropic Claude provider implementation."""
224
225     def initialize(self, api_key: Optional[str] = None, model: str = "claude-3.5", **kwargs) -> None:  # lareinahu-2023
226         """Initialize Anthropic provider.
227
228         Args:
229             api_key: Anthropic API key
230             model: Model name to use
231             **kwargs: Additional parameters
232         """
233         # TODO: Implement Anthropic provider initialization
234         self.api_key = api_key
235         self.model = model
236
237     def generate_completion(self, prompt: str, **kwargs) -> str:  # lareinahu-2023
238         """Generate completion using Anthropic API.
239
240         Args:
241             prompt: The input prompt
242             **kwargs: Additional parameters
243
244         Returns:
245             Generated completion text
246         """
247         # TODO: Implement Anthropic API call
248         return "[Anthropic response placeholder]"
249
250     @property  # lareinahu-2023
251     def provider_name(self) -> str:
252         """Return the provider name."""
253         return "Anthropic"
254
255

```

```

256 # Example usage
257 if __name__ == "__main__":
258     print("LLM Interface for SBO term annotation")
259     print("This module defines interfaces for LLM-assisted annotation")
260
261     # Example usage would be:
262     # provider = OpenAIProvider()
263     # provider.initialize(api_key="your-api-key")
264     # assistant = SBOAnnotationAssistant(provider)
265     # model = readSBML("path/to/model.xml").getModel()
266     # results = assistant.batch_process_model(model)

```

3. template_manager.py - Finally show the template system

```

1  """Template Manager for LLM-based SBO Annotation.
2
3  This module provides a template management system for loading and rendering
4  prompt templates used in LLM-based annotation.
5  """
6
7  __author__ = 'Jiahui Hu'
8
9  import os
10 import json
11 from typing import Dict, Any, Optional
12 from pathlib import Path
13 import jinja2
14
15
16 class TemplateManager: 1 usage  ± lareinahu-2023
17     """Manager for loading and rendering templates for LLM prompts."""
18
19     def __init__(self, templates_dir: Optional[str] = None) -> None:  ± lareinahu-2023
20         """Initialize the template manager.
21
22         Args:
23             templates_dir: Directory containing template files (default: 'templates')
24         """
25         # Use provided directory or default to a 'templates' directory
26         self.templates_dir = templates_dir or os.path.join(
27             os.path.dirname(os.path.abspath(__file__)), 'templates'
28         )
29
30         # Ensure templates directory exists
31         os.makedirs(self.templates_dir, exist_ok=True)
32
33         # Set up Jinja environment
34         self.env = jinja2.Environment(
35             loader=jinja2.FileSystemLoader(self.templates_dir),
36             trim_blocks=True,
37             lstrip_blocks=True,
38             autoescape=False
39         )
40
41         # Cache for loaded templates
42         self._template_cache = {}
43

```

```

44     def list_templates(self) -> list: 1 usage  ±lareinahu-2023
45         """List available templates in the templates directory.
46
47         Returns:
48             List of template names
49         """
50         return [f for f in os.listdir(self.templates_dir)
51                 if f.endswith((''.txt', '.j2', '.jinja', '.tpl'))]]
52
53     def render_template(self, template_name: str, context: Dict[str, Any]) -> str:  ±lareinahu-2023
54         """Render a template with the given context.
55
56         Args:
57             template_name: Name of the template file
58             context: Dictionary of variables to use in template rendering
59
60         Returns:
61             The rendered template string
62
63         Raises:
64             FileNotFoundError: If the template doesn't exist
65         """
66         # Load template (with caching)
67         if template_name not in self._template_cache:
68             try:
69                 self._template_cache[template_name] = self.env.get_template(template_name)
70             except jinja2.exceptions.TemplateNotFound:
71                 raise FileNotFoundError(f"Template not found: {template_name}")
72
73         template = self._template_cache[template_name]
74
75         # Render with provided context
76         return template.render(**context)
77

```

```

78     def load_template_from_string(self, template_string: str) -> jinja2.Template: 1 usage  ±lareinahu-2023
79         """Load a template from a string.
80
81         Args:
82             template_string: The template string to load
83
84         Returns:
85             Jinja Template object
86         """
87         return self.env.from_string(template_string)
88
89     def render_string_template(self, template_string: str, context: Dict[str, Any]) -> str:  ±lareinahu-2023
90         """Render a template string with the given context.
91
92         Args:
93             template_string: Template string to render
94             context: Dictionary of variables to use in template rendering
95
96         Returns:
97             The rendered template string
98         """
99         template = self.load_template_from_string(template_string)
100         return template.render(**context)
101

```

```

102     def save_template(self, template_name: str, content: str) -> str: #lareinahu-2023
103         """Save a template to the templates directory.
104
105         Args:
106             template_name: Name to save the template as
107             content: Template content
108
109         Returns:
110             Path to the saved template
111         """
112         # Add extension if not present
113         if not any(template_name.endswith(ext) for ext in ('.txt', '.j2', '.jinja', '.tpl')):
114             template_name += '.j2'
115
116         template_path = os.path.join(self.templates_dir, template_name)
117
118         with open(template_path, 'w') as f:
119             f.write(content)
120
121         # Clear cache for this template if it exists
122         if template_name in self._template_cache:
123             del self._template_cache[template_name]
124
125         return template_path
126
127     def create_default_templates(self) -> None: #usage #lareinahu-2023
128         """Create default templates if they don't exist."""
129         # Define default templates for different reaction types
130         default_templates = {
131             'reaction_base.j2': """
132 You are a Systems Biology expert tasked with assigning SBO (Systems Biology Ontology) terms to biochemical reactions.
133
134 Reaction details:
135 - ID: {{ reaction.id }}
136 - Name: {{ reaction.name }}
137 - Reversible: {{ reaction.reversible }}
138 - Reactants: {{ reaction.reactants|join(', ') }}
139 - Products: {{ reaction.products|join(', ') }}
140 {% if reaction.ec_numbers %}
141 - EC Numbers: {{ reaction.ec_numbers|join(', ') }}
142 {% endif %}
143 {% if reaction.compartments %}
144 - Compartments: {{ reaction.compartments|join(', ') }}
145 {% endif %}
146
147 Relevant SBO terms include:
148 - SBO:0000176 (Biochemical reaction): General biochemical transformation
149 - SBO:0000200 (Redox reaction): Involves electron transfer (e.g., NAD/NADH)
150 - SBO:0000216 (Phosphorylation): Transfer of phosphate groups (look for ATP/ADP)
151 - SBO:0000655 (Transport reaction): Movement across compartments
152 - SBO:0000627 (Exchange reaction): Exchange with environment
153 - SBO:0000629 (Biomass production): Overall cell growth
154
155 Based on these details, what is the most appropriate SBO term for this reaction?
156 Provide your answer in JSON format with fields: sbo_term, confidence, explanation.
157 """,
158             'transport_reaction.j2': """
159 You are a Systems Biology expert tasked with classifying transport reactions in the Systems Biology Ontology.
160
161

```

```

162 Transport Reaction details:
163 - ID: {{ reaction.id }}
164 - Name: {{ reaction.name }}
165 - Reversible: {{ reaction.reversible }}
166 - Reactants: {{ reaction.reactants|join(', ') }}
167 - Products: {{ reaction.products|join(', ') }}
168 - Source compartment: {{ reaction.source_compartment }}
169 - Target compartment: {{ reaction.target_compartment }}
170 {% if reaction.metabolites_transported %}
171 - Metabolites transported: {{ reaction.metabolites_transported|join(', ') }}
172 {% endif %}
173
174 Transport reaction SBO terms include:
175 - SBO:0000655 (Transport reaction): General transport
176 - SBO:0000657 (Active transport): Energy-requiring transport (ATP, GTP involved)
177 - SBO:0000658 (Passive transport): No energy required
178 - SBO:0000659 (Antiporter): Exchange of molecules in opposite directions
179 - SBO:0000660 (Symporter): Transport of molecules in same direction
180 - SBO:0000654 (Co-transport): Transport of multiple species
181
182 Based on these details, what is the most appropriate SBO term for this reaction?
183 Provide your answer in JSON format with fields: sbo_term, confidence, explanation.
184 """
185
186         'enzymatic_reaction.j2': """
187 You are a Systems Biology expert tasked with classifying enzymatic reactions using the Systems Biology Ontology.
188

```

```

189 Enzymatic Reaction details:
190 - ID: {{ reaction.id }}
191 - Name: {{ reaction.name }}
192 - EC Numbers: {{ reaction.ec_numbers|join(', ') }}
193 - Reversible: {{ reaction.reversible }}
194 - Reactants: {{ reaction.reactants|join(', ') }}
195 - Products: {{ reaction.products|join(', ') }}
196
197 Relevant SBO terms based on EC classification:
198 - SBO:0000200 (Oxidoreductase): EC 1.*
199 - SBO:0000402 (Transferase): EC 2.*
200 - SBO:0000376 (Hydrolase): EC 3.*
201 - SBO:0000211 (Lyase): EC 4.*
202 - SBO:0000377 (Isomerase): EC 5.*
203 - SBO:0000695 (Ligase): EC 6.*
204 - SBO:0000185 (Translocase): EC 7.*
205
206 More specific SBO terms include:
207 - SBO:0000216 (Phosphorylation): Transfer of phosphate groups
208 - SBO:0000217 (Glycosylation): Addition of glycosyl groups
209 - SBO:0000399 (Decarboxylation): Removal of carboxyl group
210 - SBO:0000400 (Decarbonylation): Removal of carbonyl group
211 - SBO:0000401 (Deamination): Removal of amino group
212
213 Based on these details, what is the most appropriate SBO term for this reaction?
214 Provide your answer in JSON format with fields: sbo_term, confidence, explanation.
215 """
216
217         'exchange_reaction.j2': """
218 You are a Systems Biology expert tasked with classifying exchange and boundary reactions in the Systems Biology Ontology.
219

```

```

219
220 Exchange Reaction details:
221 - ID: {{ reaction.id }}
222 - Name: {{ reaction.name }}
223 - Reversible: {{ reaction.reversible }}
224 - Reactants: {{ reaction.reactants|join(', ') }}
225 - Products: {{ reaction.products|join(', ') }}
226
227 Boundary reaction SBO terms include:
228 - SBO:0000627 (Exchange reaction): Exchange with environment (often prefixed with EX_)
229 - SBO:0000628 (Demand reaction): Pure consumption (often prefixed with DM_)
230 - SBO:0000632 (Sink reaction): Pure production (often prefixed with SK_)
231 - SBO:0000629 (Biomass production): Overall cell growth (contains "biomass" in name)
232 - SBO:0000630 (ATP energy): Non-growth associated maintenance
233
234 Based on these details, what is the most appropriate SBO term for this reaction?
235 Provide your answer in JSON format with fields: sbo_term, confidence, explanation.
236 """
237     }
238
239     for template_name, content in default_templates.items():
240         template_path = os.path.join(self.templates_dir, template_name)
241         if not os.path.exists(template_path):
242             with open(template_path, 'w') as f:
243                 f.write(content.strip())
244

```



```

244
245 def get_template_for_reaction(self, reaction_features: Dict[str, Any]) -> str: 1 usage  ▲ lareinahu-2023
246     """Select the most appropriate template for a reaction based on its features.
247
248     Args:
249         reaction_features: Dictionary of reaction features
250
251     Returns:
252         Template name to use for this reaction
253     """
254     # Determine the appropriate template based on reaction features
255     if 'exchange' in reaction_features.get('id', '').lower() or 'ex_' in reaction_features.get('id', '').lower():
256         return 'exchange_reaction.j2'
257     elif len(reaction_features.get('compartments', [])) > 1:
258         return 'transport_reaction.j2'
259     elif reaction_features.get('ec_numbers'):
260         return 'enzymatic_reaction.j2'
261     else:
262         return 'reaction_base.j2'
263
264

```

```

265 # Example usage
266 ▶ if __name__ == "__main__":
267     # Example usage
268     template_manager = TemplateManager()
269
270     # Create default templates if they don't exist
271     template_manager.create_default_templates()
272
273     print("Available templates:")
274     for template in template_manager.list_templates():
275         print(f"- {template}")
276
277     # Example reaction features
278     example_reaction = {
279         "id": "R_GAPD",
280         "name": "Glyceraldehyde-3-phosphate dehydrogenase",
281         "reversible": True,
282         "reactants": ["M_g3p_c", "M_nad_c", "M_pi_c"],
283         "products": ["M_13dpg_c", "M_h_c", "M_nadh_c"],
284         "ec_numbers": ["1.2.1.12"],
285         "compartments": ["c"]
286     }
287
288     # Get appropriate template for this reaction
289     template_name = template_manager.get_template_for_reaction(example_reaction)
290     print(f"\nSelected template for reaction {example_reaction['id']}: {template_name}")
291
292     # Example rendering (commented out since files might not exist yet)
293     # rendered = template_manager.render_template(template_name, {"reaction": example_reaction})
294     # print("\nRendered template:")
295     # print(rendered)

```

Technical Implementation

1. Dynamic SBO Term Retrieval System

- Implement Python API client using the requests library to fetch SBO terms from OLS server
- Develop SQLite-based caching mechanism for offline access with automatic update checking at startup
- Create differential update system to minimize bandwidth and processing overhead
- Implement local cache fallback for reliability when server access is unavailable

2. Enzymatic Data Integration Framework

- Build modular data provider system using the factory pattern for extensibility
- Develop adapters for KEGG and BRENDA databases
- Implement unified query interface via adapter pattern to standardize data access
- Create data normalization layer to resolve inconsistencies between sources

3. LLM-based Annotation Assistant

- Integrate pre-trained DistilBERT model (distilbert-base-uncased) using Hugging Face's Transformers library
- Fine-tune on a dataset of manually annotated SBML models provided by the mentor team
- Implement confidence scoring system for annotation suggestions
- Provide dual modes: automatic application of high-confidence suggestions and user selection from multiple options

4. Visualization Component

- Implement interactive reaction network visualization using Plotly and NetworkX
- Develop before/after annotation comparison views
- Create coverage analysis dashboard for model quality assessment

5. Standalone Desktop GUI

- Develop cross-platform application using PyQt6
- Design tabbed interface with drag-and-drop functionality
- Create annotation review interface with filtering capabilities

Detailed Technical Approach

LLM Implementation Details

I will implement the annotation assistant using DistilBERT (distilbert-base-uncased), a lightweight yet powerful language model that balances performance with resource requirements. Key implementation details include:

1. Training Data:

- Use a corpus of 20-30 manually annotated SBML models (to be provided by mentor)
- Supplement with existing SBO term definitions and relationships

- Implement data augmentation techniques to enhance training set size
2. **Fine-tuning Approach:**
- Focus on text classification rather than full model retraining
 - Use a sequence classification head on top of DistilBERT
 - Train on a modest-sized GPU (e.g., Google Colab or university resources)
3. **Resource Management:**
- Implement model quantization to reduce memory footprint
 - Use gradient accumulation to enable training with limited GPU memory
 - Cache inference results to improve runtime performance
4. **Fallback Strategy:**
- If fine-tuning proves too resource-intensive, implement a simpler semantic matching approach using pre-trained embeddings
 - Use cosine similarity between SBO term descriptions and reaction characteristics

Testing Methodology

I will implement a comprehensive testing approach using:

1. **Test Data:**
- Cross-validation set of 5-10 manually annotated SBML models (separate from training data)
 - Synthetic test cases for edge-case handling
 - Large-scale models to test performance and scalability
2. **Validation Metrics:**
- Precision, recall, and F1 score for annotation accuracy
 - Coverage percentage of model elements
 - Processing time for various model sizes
 - User satisfaction metrics from beta testers
3. **Continuous Testing:**
- Unit tests for each component
 - Integration tests for the complete workflow
 - UI testing for the desktop application

Risk Assessment and Mitigation

Risk	Probability	Impact	Mitigation Strategy
LLM training requires excessive resources	Medium	High	Use smaller pre-trained models; implement simpler semantic matching as fallback

OLS API integration is more complex than anticipated	Low	Medium	Focus on core term retrieval first; implement incremental features
Integration of multiple enzymatic databases creates conflicting data	Medium	Medium	Develop robust normalization rules; provide clear conflict resolution UI
GUI development exceeds time allocation	Medium	Medium	Start with minimal viable UI and enhance incrementally; prioritize functionality over aesthetics
Poor performance with large-scale models	Low	High	Implement lazy loading and processing; optimize critical path algorithms

Implementation Timeline (12 Weeks)

Community Bonding Period (2 weeks before official start)

- Set up development environment and communication channels
- Review existing codebase in detail
- Collect and analyze test SBML models
- Define API specifications in collaboration with mentor

Week 1-2: Dynamic SBO Term Retrieval System

- Implement Python API client for OLS server
- Develop SQLite caching mechanism
- Create automatic update checker
- **Milestone:** Working term retrieval with caching

Week 3-4: Enzymatic Data Integration Framework

- Build modular data provider architecture
- Implement KEGG database adapter
- Develop BRENDA database adapter
- **Milestone:** Integrated enzymatic data retrieval
- **Mentor Checkpoint:** Review data integration approach and results

Week 5-6: LLM-based Annotation Assistant (Phase 1)

- Set up Hugging Face Transformers integration
- Prepare training data with mentor assistance
- Implement model training pipeline
- **Milestone:** Initial LLM model with basic prediction capability
- **Mentor Checkpoint:** Evaluate initial model performance

Week 7-8: LLM-based Annotation Assistant (Phase 2) & Visualization

- Refine model based on mentor feedback
- Implement confidence scoring system
- Create basic reaction network visualization
- **Milestone:** Working LLM suggestion system with visualization
- **Mentor Checkpoint:** Review annotation accuracy and visualization

Week 9-10: Standalone Desktop GUI

- Develop cross-platform application shell
- Implement tabbed interface
- Create annotation review interface
- **Milestone:** Functional GUI with core features
- **Mentor Checkpoint:** Usability review and feedback

Week 11-12: Integration, Testing & Documentation

- Integrate all components
- Implement comprehensive testing
- Create user documentation
- Fix bugs and optimize performance
- **Final Milestone:** Complete working application
- **Final Mentor Review:** Comprehensive project evaluation

Mentor Collaboration Plan

I plan to maintain regular communication with my mentor throughout the project:

- **Weekly Progress Reports:** Brief written updates on completed tasks and current challenges
- **Bi-weekly Video Meetings:** Detailed discussion of progress, demonstrations, and technical guidance
- **Code Review Process:** Submit PRs for each completed component for mentor review
- **Critical Decision Points:** Specific checkpoints (noted in timeline) where mentor input is crucial for next steps

Contingency Planning

If time constraints become an issue, I will prioritize features in this order:

1. **Core Functionality:**
 - Dynamic SBO term retrieval (highest priority)
 - Basic enzymatic data integration
 - Simple GUI for accessibility

2. Enhanced Features:

- Advanced LLM integration
- Multiple data source integration
- Comprehensive visualization

If the LLM component proves particularly challenging, I will implement a simplified version using pre-trained embeddings and semantic matching rather than full model fine-tuning.

Deliverables

Required Deliverables (Will Definitely Complete)

1. Dynamic SBO term retrieval system with caching mechanism
2. KEGG and BRENDA database adapters for enzymatic data
3. Basic standalone GUI application
4. Comprehensive documentation and tests

Stretch Deliverables (Will Complete If Time Permits)

1. Advanced LLM-based annotation with high-confidence prediction
2. Interactive visualization dashboard
3. Performance optimizations for large-scale models
4. Video tutorials and demonstrations

Qualifications

As a Master's student in Computer Software Engineering at Northeastern University, I bring relevant expertise to this project:

- Proven Contribution to SBOannotator: I've already submitted a functional PR (#2: <https://github.com/draeger-lab/SBOannotator/pull/2>) implementing the LLM integration foundation
- Experience in NLP system development with fuzzy matching techniques
- Proficiency in API integration and distributed data processing
- UI development skills with PyQt
- Coursework in "Theory & Practice of AI Gen Model" providing insights into effective LLM implementation
- Strong Python programming skills with experience in scientific computing libraries

Related Work

My proposed enhancements build upon existing work in systems biology annotation tools while addressing specific gaps:

- Existing SBO term annotation tools lack dynamic updating and intelligent suggestions

- Current SBML editors provide limited visualization capabilities for annotation coverage
- Few tools integrate multiple enzymatic data sources for comprehensive annotation

Conclusion

This project represents a significant advancement for SBOannotator, transforming it from a tool with static, hard-coded functionality into a dynamic, intelligent system capable of adapting to evolving ontologies and leveraging multiple data sources. I am committed to delivering a realistic, well-tested implementation that maintains SBOannotator's core classification strengths while addressing its current limitations, with clear prioritization and contingency plans to ensure successful completion within the GSoC timeframe.