

# Technical Support Center: Automated IAB15 Classification of Scientific Papers

**Author:** BenchChem Technical Support Team. **Date:** December 2025

## Compound of Interest

Compound Name: IAB15

Cat. No.: B15619004

[Get Quote](#)

This technical support center provides troubleshooting guidance and frequently asked questions (FAQs) for researchers, scientists, and drug development professionals experimenting with automated **IAB15** classification of scientific papers.

## Frequently Asked Questions (FAQs)

### Q1: What is IAB15 classification and why is it challenging for scientific papers?

The Interactive Advertising Bureau (IAB) content taxonomy is a standardized hierarchy for categorizing content. The **IAB15** category represents "Science," with subcategories such as Biology (**IAB15-2**), Chemistry (**IAB15-3**), and Physics (**IAB15-6**). Automating the classification of scientific papers into these categories is challenging due to several factors:

- **Semantic Ambiguity and Jargon:** Scientific papers often contain highly specialized terminology and jargon that may have different meanings in different contexts, making it difficult for automated systems to interpret correctly. The overuse of acronyms and complex sentence structures further complicates text comprehension.<sup>[1][2][3][4][5]</sup>
- **Hierarchical Complexity:** The **IAB15** taxonomy is hierarchical. A paper on "biochemical signaling in cancer cells" could potentially fit under Biology (**IAB15-2**) and Chemistry (**IAB15-3**), requiring a hierarchical classification approach that can handle multiple levels of granularity.<sup>[6][7][8][9][10][11]</sup>

- **Multi-Label Nature:** A single scientific paper can often be relevant to multiple **IAB15** subcategories. For instance, a study on the geological impact of climate change could fall under Geology (**IAB15-4**) and Weather (**IAB15-10**). This necessitates multi-label classification models.[\[12\]](#)[\[13\]](#)[\[14\]](#)[\[15\]](#)
- **Interdisciplinary Research:** Modern scientific research is often interdisciplinary, blurring the lines between traditional categories. This makes it difficult to assign a single, definitive category to a paper.

## Q2: My model is performing poorly on classifying papers into specific **IAB15** subcategories. What are the common causes?

Poor performance in specific subcategories can often be attributed to:

- **Class Imbalance:** Some **IAB15** subcategories may have a much larger representation in your training dataset than others. This can lead to a model that is biased towards the majority classes.
- **Insufficient Training Data:** For niche scientific fields, there may not be enough labeled examples for the model to learn the distinguishing features of that category.
- **Inadequate Preprocessing:** Scientific text requires specialized preprocessing to handle its unique characteristics. Failure to properly address jargon, special characters, and complex sentence structures can lead to poor model performance.
- **Inappropriate Model Choice:** A simple text classification model may not be sufficient for the complexities of scientific literature. More advanced models, such as those based on transformer architectures (e.g., SciBERT, BioBERT), are often better suited for this task.[\[6\]](#)[\[8\]](#)[\[16\]](#)[\[17\]](#)

## Q3: How do I choose the right machine learning model for **IAB15** classification?

The choice of model depends on the complexity of your task and the resources available.

- Traditional Models: For a baseline, models like Support Vector Machines (SVM) and Naive Bayes can be effective, especially with proper feature engineering (e.g., TF-IDF).
- Deep Learning Models: For more nuanced understanding of scientific text, deep learning models are generally superior.
  - BERT (Bidirectional Encoder Representations from Transformers): A powerful general-purpose language model.
  - SciBERT: A BERT model pre-trained on a large corpus of scientific publications, making it particularly effective for this domain.[\[6\]](#)[\[8\]](#)[\[16\]](#)[\[17\]](#)
  - BioBERT: A BERT model pre-trained on biomedical literature, ideal for papers in the life sciences.[\[6\]](#)[\[16\]](#)[\[18\]](#)

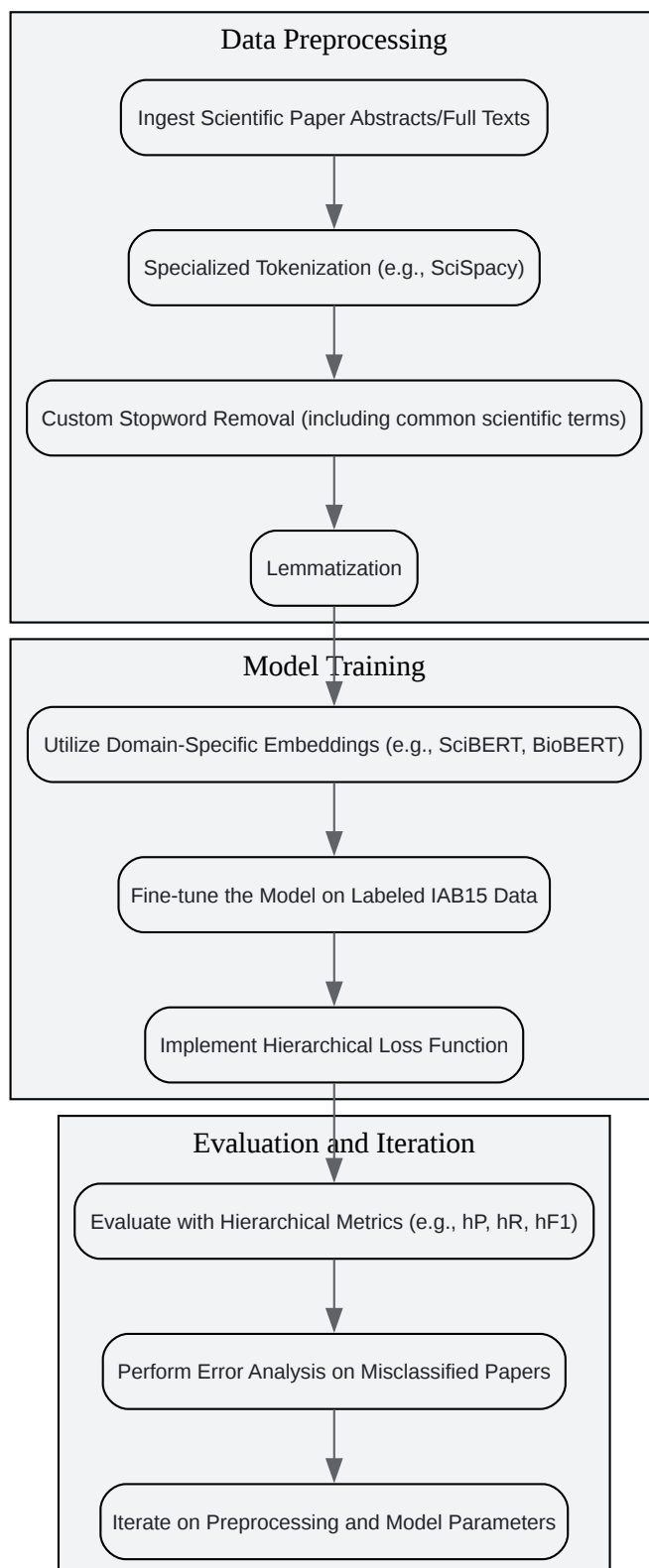
Recent studies have shown that SciBERT and BioBERT often outperform general BERT models on scientific text classification tasks.[\[6\]](#)[\[8\]](#)[\[16\]](#)[\[17\]](#)[\[18\]](#)

## Troubleshooting Guides

### Guide 1: Improving Model Accuracy for Ambiguous Terminology

**Problem:** Your model struggles to differentiate between closely related **IAB15** subcategories, such as Biology (**IAB15-2**) and Botany (**IAB15-9**), due to overlapping terminology.

**Solution Workflow:**



[Click to download full resolution via product page](#)

Workflow for improving model accuracy with ambiguous terminology.

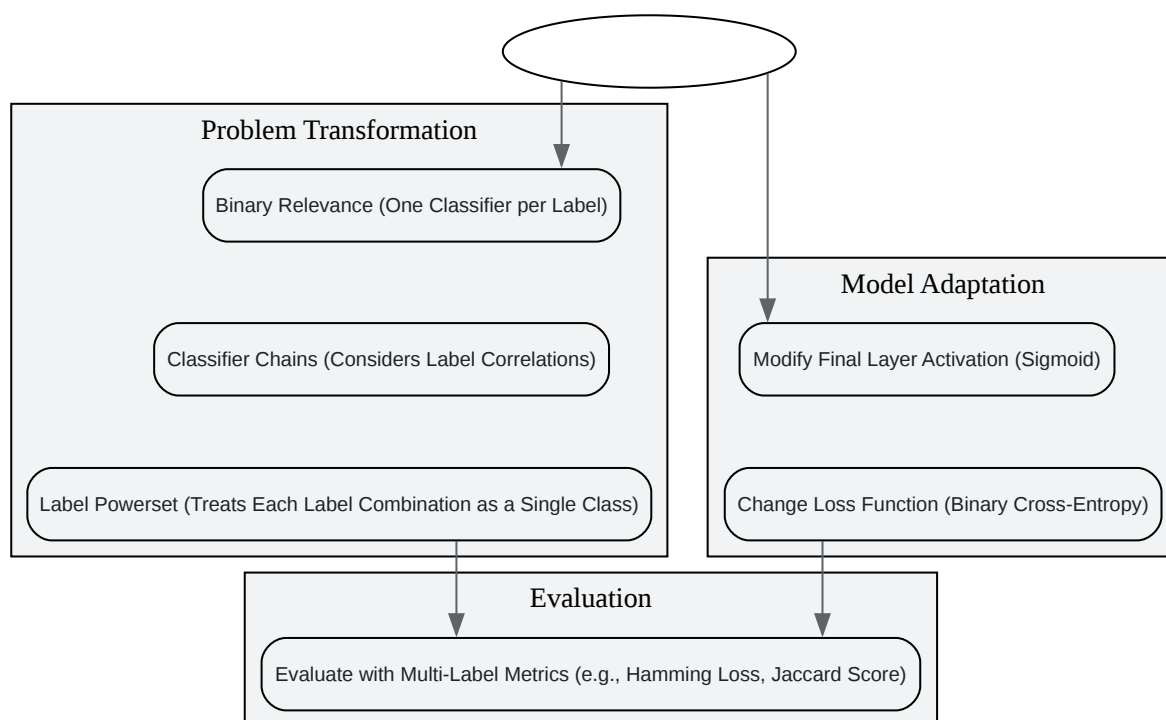
#### Detailed Steps:

- Specialized Preprocessing:
  - Tokenization: Use tokenizers designed for scientific text, such as SciSpacy, which can handle hyphens, chemical formulas, and other domain-specific syntax.
  - Custom Stopwords: Create a custom list of stopwords that includes common but non-discriminatory scientific terms (e.g., "study," "results," "method").
  - Lemmatization: Reduce words to their base form to group related terms.
- Domain-Specific Models:
  - Employ pre-trained language models like SciBERT or BioBERT that have been trained on large corpora of scientific and biomedical texts, respectively.<sup>[6][8][16][17][18]</sup> These models have a better understanding of scientific jargon and context.
- Hierarchical Classification Techniques:
  - Implement a hierarchical classification approach where the model first predicts the top-level category (**IAB15** - Science) and then the appropriate subcategories. This can be achieved using hierarchical loss functions that penalize errors at different levels of the hierarchy differently.

## Guide 2: Addressing Multi-Label Classification Challenges

**Problem:** Your model only assigns one **IAB15** subcategory to a paper, even when multiple categories are relevant.

**Solution Workflow:**



[Click to download full resolution via product page](#)

Approaches for handling multi-label **IAB15** classification.

Detailed Steps:

- Model Architecture Modification:
  - Final Activation Layer: Change the final activation function of your neural network from softmax (which is used for single-label classification) to sigmoid. This will output a probability for each class independently.
  - Loss Function: Use a binary cross-entropy loss function instead of categorical cross-entropy.
- Problem Transformation Methods:

- Binary Relevance: Train a separate binary classifier for each **IA B15** subcategory. This is a simple and often effective approach.
- Classifier Chains: Similar to binary relevance, but each classifier in the chain also receives the predictions of the previous classifiers as features, which can help model label correlations.
- Label Powerset: Transform the problem into a multi-class classification problem where each unique combination of labels is treated as a single class.
- Evaluation Metrics:
  - Use metrics designed for multi-label classification, such as Hamming Loss, Jaccard Score, and F1-score (macro, micro, and samples averaging).[\[12\]](#)[\[14\]](#)[\[15\]](#)

## Quantitative Data Summary

### Table 1: Performance Comparison of Models on Scientific Text Classification

| Model   | Dataset              | Accuracy | Micro F1-Score | Macro F1-Score | Reference            |
|---------|----------------------|----------|----------------|----------------|----------------------|
| BERT    | WoS-11967 (keywords) | 84%      | 0.85           | -              | <a href="#">[6]</a>  |
| SciBERT | WoS-11967 (keywords) | 87%      | 0.87           | -              | <a href="#">[6]</a>  |
| BioBERT | WoS-11967 (keywords) | 86%      | 0.85           | -              | <a href="#">[6]</a>  |
| BERT    | WoS-5736 (abstracts) | 97%      | 0.98           | -              | <a href="#">[6]</a>  |
| SciBERT | WoS-5736 (abstracts) | 98%      | 0.97           | -              | <a href="#">[6]</a>  |
| BioBERT | WoS-5736 (abstracts) | 98%      | 0.99           | -              | <a href="#">[6]</a>  |
| CovBERT | COVID-19 PubMed      | 94%      | -              | -              | <a href="#">[16]</a> |

WoS: Web of Science dataset

## Table 2: Common Multi-Label Classification Metrics



| Metric                                  | Description   | Interpretation  |
|---|---|---|
| Hamming Loss                            | The fraction of labels that are incorrectly predicted.  | Lower is better.                                      |
| Jaccard Score (Intersection over Union) | The size of the intersection of predicted and true labels divided by the size of their union.   | Higher is better.                                     |
| F1-Score (Micro Average)                | Calculated globally by counting the total true positives, false negatives, and false positives. | Higher is better; gives more weight to common labels. |
| F1-Score (Macro Average)                | Calculated for each label and then averaged.  | Higher is better; treats all labels equally.          |
| F1-Score (Samples Average)              | Calculated for each instance and then averaged.   | Higher is better.                                     |

## Experimental Protocols

### Protocol 1: A Step-by-Step Guide to Training a SciBERT Model for IAB15 Classification

This protocol outlines the key steps for fine-tuning a SciBERT model for multi-label classification of scientific papers into the **IAB15** taxonomy.

#### 1. Data Preparation:

- **Data Collection:** Gather a labeled dataset of scientific papers with their corresponding **IAB15** subcategories.
- **Data Cleaning:** Remove irrelevant information such as HTML tags, special characters, and author affiliations.[\[9\]](#)[\[19\]](#)
- **Train-Validation-Test Split:** Split your dataset into training, validation, and testing sets (e.g., 80-10-10 split).[\[20\]](#)

## 2. Preprocessing Pipeline:

- **Tokenization:** Use the SciBERT tokenizer to convert the text into tokens that the model can understand.
- **Input Formatting:** Create input tensors for the model, including `input_ids`, `attention_mask`, and labels.

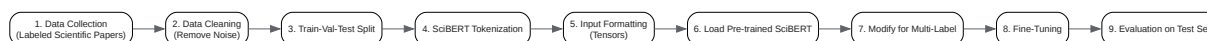
## 3. Model Training:

- **Load Pre-trained Model:** Load the pre-trained SciBERT model.
- **Modify for Multi-Label:** Add a classification head with a sigmoid activation function for multi-label classification.
- **Define Optimizer and Loss Function:** Use an Adam optimizer and binary cross-entropy loss. [\[21\]](#)
- **Fine-Tuning:** Train the model on your labeled dataset. Monitor the validation loss to prevent overfitting. [\[21\]](#)

## 4. Evaluation:

- **Prediction:** Use the fine-tuned model to make predictions on the test set.
- **Metrics Calculation:** Evaluate the model's performance using appropriate multi-label classification metrics (see Table 2).

## Experimental Workflow Diagram:



[Click to download full resolution via product page](#)

Experimental workflow for training a SciBERT classification model.

**Need Custom Synthesis?**

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: [info@benchchem.com](mailto:info@benchchem.com) or [Request Quote Online](#).

## References

- 1. Image Classification With Small Datasets: Overview and Benchmark | IEEE Journals & Magazine | IEEE Xplore [[ieeexplore.ieee.org](https://ieeexplore.ieee.org)]
- 2. Automatic jargon identifier for scientists engaging with the public and science communication educators - PubMed [[pubmed.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)]
- 3. researchgate.net [[researchgate.net](https://researchgate.net)]
- 4. Automatic jargon identifier for scientists engaging with the public and science communication educators | PLOS One [[journals.plos.org](https://journals.plos.org)]
- 5. Automatic jargon identifier for scientists engaging with the public and science communication educators - PMC [[pmc.ncbi.nlm.nih.gov](https://pmc.ncbi.nlm.nih.gov)]
- 6. Fine-Tuning Large Language Models for Scientific Text Classification: A Comparative Study [[arxiv.org](https://arxiv.org)]
- 7. mdpi.com [[mdpi.com](https://mdpi.com)]
- 8. sh-tsang.medium.com [[sh-tsang.medium.com](https://sh-tsang.medium.com)]
- 9. blog.devgenius.io [[blog.devgenius.io](https://blog.devgenius.io)]
- 10. direct.mit.edu [[direct.mit.edu](https://direct.mit.edu)]
- 11. aclanthology.org [[aclanthology.org](https://aclanthology.org)]
- 12. causeweb.org [[causeweb.org](https://causeweb.org)]
- 13. arxiv.org [[arxiv.org](https://arxiv.org)]
- 14. arxiv.org [[arxiv.org](https://arxiv.org)]
- 15. towardsdatascience.com [[towardsdatascience.com](https://towardsdatascience.com)]
- 16. mdpi.com [[mdpi.com](https://mdpi.com)]
- 17. kyleclo.com [[kyleclo.com](https://kyleclo.com)]
- 18. researchgate.net [[researchgate.net](https://researchgate.net)]
- 19. drlee.io [[drlee.io](https://drlee.io)]

- 20. [betterprogramming.pub \[betterprogramming.pub\]](#)
- 21. Step 4: Build, Train, and Evaluate Your Model | Machine Learning | Google for Developers [[developers.google.com](#)]
- To cite this document: BenchChem. [Technical Support Center: Automated IAB15 Classification of Scientific Papers]. BenchChem, [2025]. [Online PDF]. Available at: [<https://www.benchchem.com/product/b15619004#challenges-in-automated-iab15-classification-of-scientific-papers>]

---

### Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:** The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

**Need Industrial/Bulk Grade?** [Request Custom Synthesis Quote](#)

## BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

### Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: [info@benchchem.com](mailto:info@benchchem.com)