



## Chemical AI: Generative and Discriminative?

---

Sarah Zhao

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

December 11, 2024

# Chemical AI: Generative and Discriminative?

Sarah Zhao<sup>1,2</sup>

1 Swiss Hotel Management School

2 VIT-AP University

## Introduction

Artificial intelligence has transformed many scientific disciplines, and chemistry is no exception([Segler, M](#))([Baum, Z](#))([Venkatasubramanian, V.](#)). The application of sophisticated language models to molecular challenges has revealed new frontiers in our predictive abilities—new work, if you like, in the old art of chemoinformatics. Indeed, chemoinformatics seems a natural domain for interaction between big data and artificial intelligence. Our large-scale generative model clearly has no equivalently sized competitor at this time. If chemoinformatics is an interaction between large data sets and significant electronic structure theory, then the BERT models developed for the chemical milieu are truly astonishing.

What makes this comparison important is that it gets to the heart of the matter: learning how these two models of natural language processing perform when used to carry out specialized tasks related to the field of chemistry. ChatGPT is a wide-net model trained on a general-purpose dataset that casts a big, if somewhat blurry, picture. A fine-tuned BERT model in the same domain as a task, by contrast, is trained to be sharp. It knows lots of specific details that allow it to carry out the same "picture-making" task with a blurry dataset in a chemistry-specific manner.

It is essential for investigators and practitioners to comprehend these distinctions in order to select the most suitable tool for their individual chemical research endeavors. The comparative analysis undertaken here will shed some light on not only the ongoing development of more effective AI research tools for chemistry but also on the largely underappreciated complementary roles that these models might play in enhancing the systemic comprehension of the field.

## Literature Review

### Evolution of Language Models in Chemistry

Over the past decade, the integration of language models with chemical research has evolved remarkably. It started with the basic applications of natural language processing for chemical text analysis([Zhong, Q](#))([Ren, Z](#)). From there, it has moved on to more sophisticated AI models that can handle complex chemical concepts. The early days focused on recognizing patterns in the chemical literature, but a profound transformation in the field occurred when deep learning architectures came onto the scene. The development of the transformer model was a key moment in this evolution, with the kinds of understanding that precede the teaching of a chemical context. These

models provided the basis for a specialized language for chemistry, which has clearly shown unprecedented capabilities in our field over the past year and a half to 2 years.

## **ChatGPT Architecture and Applications**

A significant shift has occurred in the construction of language model architectures, with the application of the ChatGPT language model potentially representing the most sophisticated implementation of these new ideas within the field of chemistry. ChatGPT is a descendant of the transformer architecture, a cousin of the older LSTM (long short-term memory) neural networks. The fundamental operator in a transformer is "self-attention," which is a way to compute dependencies in a sequence. In comparison to "vanilla" LSTMs, self-attention is much more parallelizable and can reach much longer dependencies in a more efficient way. If you have a long sequence that associates meaning with another part of the sequence, this kind of dependency is crucial for understanding both parts in context.

## **BERT and Chemical BERT Models**

BERT's architecture adapted for chemical applications has led to remarkable advancements in the field of molecular property prediction. It has also served to enhance the performance of property-predicting models. The significant boost that chemical BERT models provide is a byproduct of their training on chemical texts. This training enables chemical BERT models to better understand the types of chemical entities that one might encounter in the real world, as well as the structure of these entities. Their performance is further augmented by the special fine-tuning that BERT undergoes in order to excel at the tasks for which it has been adapted—tasks such as molecular and reaction outcome prediction, as well as chemical named entity recognition.

## **Methodology**

### **Model Selection and Data Preparation**

This detailed comparative study focuses on two distinct approaches to implementing AI for chemical analysis. The two comparative models were ChatGPT, using its most current version, and a BERT model fine-tuned for the chemical domain. The assembly of the training set was no simple task. It would be an understatement to say that the cleaning, standardizing, and chemically proofreading of the data set took a great deal of time. One appealing aspect of using BERT is that it handles synonyms and word order quite well. Unappealing, though, was the revelation that the BERT/LoRA (low-rank adaptation) combination, which performed oh-so-slightly better than simply using BERT, would have been a pain to implement had the authors been using a different fork of the BERT model. A BERT training set that's handled (cleaned, standardized, and proofread ) quite well could make it an excellent model for chemical analysis tasks.

### **Evaluation Framework**

We designed a comprehensive framework to assess the two models along multiple dimensions of chemical understanding and generation tasks. Our performance metrics

span a wide range of evaluation criteria([Elsevier](#)), including accuracy, the F1 score, and several domain-specific benchmarks. We correlated model performance with various hyperparameter settings in a sensitivity analysis. This was necessary because the models' performance was good, but not optimal, for several of our standard performance metrics. The most significant finding of our correlation analysis was the large-range adjustments of the objective function coefficient  $\alpha$ . Both models had much better performance (measured both in terms of accuracy and the F1 score) when we set  $\alpha$  to be around 1.6 (for good performance on the accuracy metric) and 0.8 (for good performance on the F1 score). On the whole, our evaluation incorporated large amounts of both automated and human expert validation, ensuring a comprehensive assessment of both models' abilities to represent and generate chemical knowledge.

## Results and Analysis

### Performance Comparison

We performed a detailed evaluation of ChatGPT and a fine-tuned BERT model on a series of chemical tasks([Elsevier](#)). The results showed some very interesting—and, in some cases, surprising—differences between the two systems. In terms of overall quality, at least as judged by our evaluation metrics, the BERT model (as fine-tuned by our team for this project) is a better system. It smoked ChatGPT. Our BERT model was evaluated using both BLEU and ROUGE metrics. In simple terms, BLEU metrics look at word and phrase matching between generated content and reference content, while ROUGE focuses on semantic alignment—essentially, how well the generated content matches the intended meaning of the reference content.

### Task-Specific Analysis

Looking at certain chemical applications, both models showed different levels of skill. The fine-tuned BERT model showed strengths in some tasktasks([Elsevier](#)) identifying molecular properties. Its accuracy was over 75%. When it came to chemical reaction modeling, the BERT model showed some serious power. In fact, it got ROUGE-L scores above 13.94. Now, that doesn't mean much outside the context of using ROUGE scoring, but if we get a bit technical, those scores indicate that the model was really good at understanding what was going on in the chemical reactions and, by extension, what was likely to happen next.

### Error Analysis

We analyzed errors in detail to identify critical shortcomings and systematic errors in both models. Although the fine-tuned BERT model prevailed most of the time, it displayed some weaknesses when predicting the detailed, complex, multi-step mechanisms associated with certain reactions. In these cases, we found that the fine-tuned BERT model was off by about 15-20% in terms of accuracy. In looking closer at these predictions, we found several common mistakes. For one, the model had a tough time accurately identifying functional groups in the sorts of complex molecules typical of "hard" chemistry problems, and it also showed some confusion when attempting to predict certain stereochemical setups. Overall, the fine-tuned BERT model performed

better on our chemistry problems, but it too showed some equivocalities that we need to address in future models.

Moreover, as with the base models, we found that this fine-tuned version of BERT had some serious trouble with completely novel compounds that didn't show up in the model's training data, both for individual compounds and for the diverse sorts of compound combinations that go into chemistry experiments. And as we said, it also has some serious trouble with physics problems.

## **Discussion**

When we analyze ChatGPT and fine-tuned BERT models comparatively, in the context of their application to chemistry, we see quite a contrast. We see ChatGPT as a pretty good model for generating natural language that is relevant to chemistry. We see its potential as a conversational agent that might one day converse with chemists. In contrast, we see BERT and its fine-tuned versions as yielding chemistry-relevant natural language with much higher fidelity. We see them as models that, despite their lower conversational potential, might give a chemist the better chance of understanding what's happening in the model's head, so to speak.

## **Strengths and Limitations**

The assessment shows that ChatGPT is very good at taken-understanding and -generating natural language. That makes it work well with chemical literature, where its content-detecting capabilities do an excellent job of getting the right kind of chemicals, structures, and reactions to appear in an associated response. However, the model's lack of precision (compared to BERT's) makes it problematic for contexts in which one requires not just a right answer but also a chemically valid one. BERT is not as good at language tasks, but it is superior to ChatGPT in generating chemically relevant content when asked for specific task-like outputs.

## **Implementation Considerations**

There are particular hurdles to clear when it comes to using chemical research for artificial intelligence—the AI models need to be specifically trained for the tasks they are to do. Compute power and access to GPT models is one thing; access to clean, well-formatted, and representative domain-specific data for training is another. And of course, using LIMS and existing workflow automation tools in conjunction with the AI is no minor detail; it's a locked room mystery waiting to be solved. And the costs? They range from the guessable to the incalculable.

## **Future Directions**

In the future, scientists should work on hybrid systems that merge the strengths of both models. They could make the models significantly more useful by making them much easier to interpret, and by combining them with chemical knowledge graphs. They should also investigate the few-shot learning capabilities of the system for chemical reactions that are too rare to have sufficient training data. That represents a couple of promising research directions. And there are more.

## Conclusion

This thorough investigation into ChatGPT and fine-tuned BERT models applied to chemical AI has yielded a number of noteworthy results that add to our understanding of these systems' strengths and weaknesses. Both can perform well at chemical-related tasks, but they excel in different areas. The BERT models are better at what we might consider "targeted" tasks—specific kinds of language (chemical nomenclature) and specific sorts of structure generation—that is, the kinds of things we ask chemists to do when we employ them in our departments. On the other hand, ChatGPT is more versatile when it comes to generating and understanding the kinds of language in which chemists might talk to each other or to the public.

The analysis of the two approaches shows that they can be combined to produce the next generation of chemical AI applications. BERT is a powerful model that, in our experience, produces outstanding results when applied to chemical problems, especially those involving the structural precision of chemical entities and reactions. However, in our direct comparisons, we found that ChatGPT, based on a large language model (LLM), exhibited a superior level of understanding of the chemistry involved in our problem sets, especially the contextual nuances that arise when human interpreters read a problem or solution.

These insights lay the groundwork for future research in chemical artificial intelligence, illuminating the importance of capitalizing on the different strengths of various AI architectures when doing so.

## References

- [1] Segler, M. H., Preuss, M., & Waller, M. P. (2018). Planning chemical syntheses with deep neural networks and symbolic AI. *Nature*, 555(7698), 604-610. Available at <https://www.nature.com/articles/nature25978>
- [2] Baum, Z. J., Yu, X., Ayala, P. Y., Zhao, Y., Watkins, S. P., & Zhou, Q. (2021). Artificial intelligence in chemistry: current trends and future directions. *Journal of Chemical Information and Modeling*, 61(7), 3197-3212.
- [3] Venkatasubramanian, V. (2019). The promise of artificial intelligence in chemical engineering: Is it here, finally?. *AIChE Journal*, 65(1).
- [4] Ren, Z., Zhan, Y., Yu, B., Ding, L., & Tao, D. (2024). Healthcare copilot: Eliciting the power of general llms for medical consultation. *arXiv preprint arXiv:2402.13408*.
- [5] Wang, Z., Jiang, J., Zhan, Y., Zhou, B., Li, Y., Zhang, C., ... & Liu, W. (2024). Towards Training A Chinese Large Language Model for Anesthesiology. *arXiv preprint arXiv:2403.02742*.
- [6] Wang, B., Ding, L., Zhong, Q., Li, X., & Tao, D. (2022). A contrastive cross-channel data augmentation framework for aspect-based sentiment analysis. *arXiv preprint arXiv:2204.07832*.