



**HAL**  
open science

## An Instruction Dataset for Extracting Quantum Cascade Laser Properties from Scientific Text Authors

Deperias Kerre, Anne Laurent, Kenneth Maussang, Dickson Odhiambo Owuor

### ► To cite this version:

Deperias Kerre, Anne Laurent, Kenneth Maussang, Dickson Odhiambo Owuor. An Instruction Dataset for Extracting Quantum Cascade Laser Properties from Scientific Text Authors. *Data in Brief*, 2025, 58 (111255), 10.1016/j.dib.2024.111255 . hal-04866715

**HAL Id: hal-04866715**

**<https://hal.science/hal-04866715v1>**

Submitted on 6 Jan 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

1 **Article Information**

2 **Article title**

3 An Instruction Dataset for Extracting Quantum Cascade Laser Properties from Scientific Text

4 **Authors**

5 Deperias Kerre<sup>a,c</sup> \*, Anne Laurent<sup>a</sup>, Kenneth Maussang<sup>b</sup>, Dickson Owuor<sup>c</sup>

6 **Affiliations**

7 <sup>a</sup> LIRMM, Univ Montpellier, CNRS, Montpellier, France

8 <sup>b</sup> Institut d'Electronique et des Systemes, UMR 5214, Univ Montpellier, CNRS, Montpellier, France

9 <sup>c</sup> SCES, Strathmore University, Nairobi, Kenya

10 **Corresponding author's email address and Twitter handle**

11 \*Corresponding author: Deperias Kerre, Email: dkerre@strathmore.edu

12 **Keywords**

13 Information Extraction, Large Language Models, Machine Learning, Quantum Cascade Lasers

14 **Abstract**

15 Quantum Cascade Lasers (QCL) are promising semiconductor lasers, compact and powerful, but of  
16 complex design. Availability of structured data of the QCL properties can support data mining activities  
17 that seek to understand the relationship between these properties, for instance between the design  
18 and performance features. The main open source of QCL data is in scientific text which in most cases  
19 is usually unstructured. One of the ways to extract and organize this data is by utilizing Information  
20 Extraction techniques. These techniques can accelerate the process of curating QCL properties data  
21 from scientific articles for further analysis. One of the main challenges in developing machine learning  
22 algorithms for extraction of QCL properties from text is lack of quality training data for these  
23 algorithms. Large Language Models (LLMs) have demonstrated great capabilities in materials property  
24 extraction from text. They however experience challenges with domain specific properties, for  
25 instance the heterostructure and design types in the QCL domain hence for adaptation. In this paper,  
26 we present an original instruction dataset for training and evaluation of large language models (LLMs)  
27 for QCL properties extraction from text. The data is generated by augmenting sample sentences from  
28 scientific articles with GPT-3.5 instruct with a few shot strategy. The dataset then is manually annotated  
29 with the help of QCL experts and is composed of 1300 rows of training examples consisting of an  
30 Instruction, Input Text and the Output.

31 **SPECIFICATIONS TABLE**

|                              |   |
|------------------------------|---|
| <b>Subject</b>               | Computer Science, Materials Science   |
| <b>Specific subject area</b> | Artificial Intelligence (Information Extraction and Text Mining)- Extraction of Quantum Cascade Laser properties from text. |

|                             |   |
|-----------------------------|---|
| <b>Type of data</b>         | Processed data in CSV format.   |
| <b>Data collection</b>      | The dataset was generated by augmenting sample sentences from gold open access articles in the Quantum Cascade Laser domain. The sentences are augmented using the GPT-3.5 instruct model in order to generate similar sentences with various transformations such as rephrasing, new values of properties and generation of new sentences capturing same concepts but with different examples. The data is then labelled with an instruction to the large language model and the expected output for the property to be extracted. |
| <b>Data source location</b> | Sentences from 12 Gold open access articles from the AIP (American Institute of Physics) and IOP (Institute of Physics) publishers in accordance with the text and data mining policies. The articles document proposed Quantum Cascade Laser devices and their properties.   |
| <b>Data accessibility</b>   | Repository name: Recherche Data Gouv<br>Data identification number: 10.57745/U3U7XR<br>Direct URL to data: <a href="https://doi.org/10.57745/U3U7XR">https://doi.org/10.57745/U3U7XR</a><br>The dataset is licensed under the CC-BY license to enable re-use.   |

32

### 33 VALUE OF THE DATA

- 34 • The dataset provides a background in understanding the relationship between the various  
35 Quantum Cascade Laser properties. This is by enabling training of machine learning models to  
36 extract these properties from text to generate data that can be used for analysis for instance,  
37 in the prediction of the laser performance properties based on the design features.
- 38 • The dataset provides a benchmark dataset for evaluating large language models' performance  
39 for the task of information extraction of properties in the Quantum Cascade Laser domain. It  
40 covers the following properties: working temperature, frequency, power, the heterostructure  
41 stacking materials and design types. LLMs experience challenges in identifying QCL domain  
42 specific properties with direct prompting.
- 43 • The dataset gives a wide range of diverse sentences that capture the Quantum Cascade Laser  
44 properties in various ways hence enabling models to learn the various ways of expressing the  
45 properties. Most of the QCL properties are summarized in the abstract with the properties  
46 mentioned at sentence level. The dataset is therefore suitable for the extraction of properties  
47 from the abstracts. The focus is on articles describing a single novel QCL device hence enabling  
48 the generation of structured QCL properties data for various devices discussed in different  
49 papers.
- 50 • The dataset serves as central resource for researchers in the field of Information Extraction  
51 with an interest of developing efficient methodologies for materials science data extraction  
52 from scientific text. The dataset can be extended by further augmentation or with additional  
53 Quantum Cascade Laser properties to provide a diverse dataset for properties data extraction  
54 from text. The format of the dataset can also be adopted in developing information extraction  
55 datasets for other properties in materials science.

## 56 BACKGROUND

57 The main motivation behind the generation of this dataset is for understanding the relationship  
58 between the various properties of Quantum Cascade Lasers (QCLs). A Quantum Cascade Laser (QCL)  
59 is a semiconductor laser device made up of different layers of semiconductor materials and with a  
60 spectral emission range in the mid-to far infrared [1]. A QCL device with certain design characteristics  
61 (the materials layers and the design types) has corresponding working properties such as the  
62 temperature, power and frequency among others. This implies that a particular design has an influence  
63 on the working properties of the proposed QCL device. Understanding these relationships between  
64 properties can give insights on the fabrication of QCL devices with target properties. One of the ways  
65 to explore and understand the QCL data is by capturing it in knowledge representation systems such  
66 as Ontologies and Knowledge Graphs. There already exists ontologies for concepts such as units of  
67 materials properties such as the QUDT Ontology [2] and Ontologies for other characteristics of  
68 materials data such as measurable properties for instance the Materials Design Ontology [3]. There is  
69 also a proposal for the Ontology for the QCL properties in order to formally represent these properties  
70 [4]. With already such existing resources for representing the QCL properties, access to structured QCL  
71 properties data will accelerate the process of developing platforms for analysing the relationships  
72 between the QCL properties.

73 In most cases, the QCL design and working properties are captured in scientific literature. An example  
74 of a description of QCL laser properties can be given as “A Terahertz quantum cascade laser with a  
75 rather high injection coupling strength based on an indirectly pumped scheme is designed and  
76 experimentally implemented. To effectively suppress leakage current, the chosen quantum cascade  
77 module of the device is based on a five-well GaAs/Al<sub>0.25</sub>Ga<sub>0.75</sub>As structure. The device lases up to  
78 151 K with a lasing frequency of 2.67 THz.”[5]. This text description contains various QCL properties  
79 for instance the heterostructure, working temperature, frequency and design type. Rule-based  
80 methods have been proposed for QCL properties extraction from text [6]. These are however limited  
81 to the reporting style in text and requires a lot of expertise and effort in rewriting of the rules when  
82 there is even a slight change in the text structure or being applied to other QCL properties.

83  
84 The existence of QCL properties in scientific text necessitates the need for efficient methods to extract  
85 these properties to collect data for analysis. The QCL properties data constitutes domain specific data  
86 that cannot be efficiently mined by information extraction methods without adaptation. This calls for  
87 quality datasets that can be used in training these algorithms to efficiently extract these properties  
88 from scientific text. Based on this motivation, we propose an instruction dataset that provides an initial  
89 attempt in generating quality data for developing efficient algorithms for QCL property extraction from  
90 text [7].

## 91 DATA DESCRIPTION

92 The dataset consists of a CSV file with 1300 rows of sentences generated from sample sentences from  
93 scientific articles. The data is availed a single file without any partitions for instance the training, test  
94 or validation sets. The CSV file has three columns: Instruction, Input Text and Output. The Instruction  
95 column consists of sentences that entails an instruction to the large language model to extract a  
96 specific QCL property from the text. The input text column consists of sentences containing the QCL  
97 properties mentioned in them. The output column contains the expected output of extracted  
98 properties based on the user inquiry in the instruction column. From the machine learning perspective,  
99 the output column may be seen as the label for the input text and the instruction.

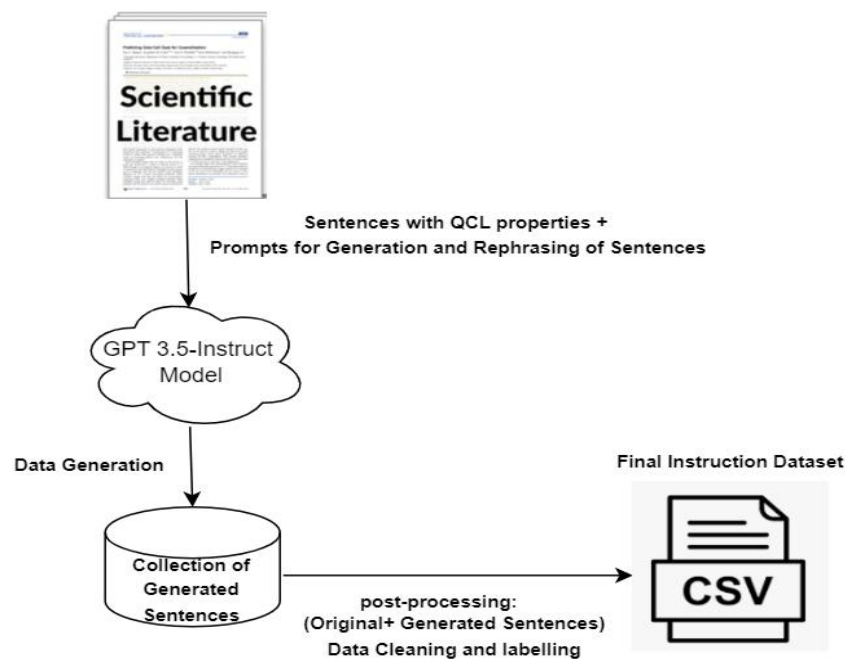
## 100 MATERIALS AND METHODS

### 101 Sentence Sampling

102 In order to obtain sample sentences for augmentation, we focus on scientific open access articles that  
103 solely describes originally proposed QCL devices and the corresponding properties . We don't consider  
104 articles that proposes applications of these devices in other applications. We also limit this on articles  
105 whose publishers permit text and data mining for research purposes. Through this process, we were  
106 able to extract 36 sentences from 12 scientific articles to be used in the augmentation process [5], [8-  
107 18]. The sentences contain the various QCL properties of interest: working temperature, power, laser  
108 frequency and the laser design types.

### 109 Dataset Generation

110 Large Language Models have shown great progress in text data augmentation [19,20]. Inspired by this,  
111 we employ a few shot strategy to generate an augmented dataset of QCL properties from sample  
112 sentences extracted from the scientific articles using GPT 3.5 instruct model. The sentences are parsed  
113 to the GPT-3.5 instruct model in a few short manner to paraphrase and generate similar sentences  
114 with varying expression styles and values for the QCL properties. The process is repeated in iterations  
115 hence obtaining different sentences with various values for the laser properties. The generated  
116 sentences are validated by an expert to make sure those that make sense are the only one included in  
117 the final dataset. The generated sentences are combined with the original sentences to make a total  
118 of 1300 rows. Figure 1. Dataset Generation Process shows the overall process of the dataset  
119 generation. The LLM prompts and the original sentences are publicly availed at the Gitlab repository:  
120 <https://gite.lirmm.fr/dkerre/qplInstruct> .



121

122

Figure 1. Dataset Generation Process

### 123 Data Annotation

124 The final stage entails annotating the data with various kinds of information to be extracted from the  
125 text. Two columns are added i.e. the Instruction column and the output column. In the instruction  
126 column, the instructions to the LLM models are captured. The output column consists of the expected

127 properties to be extracted as per the instruction in the Instruction column. The output column is  
128 labelled with assistance from the QCL domain experts.

## 129 LIMITATIONS

130 The dataset only captures the following QCL properties: the working temperature, power, frequency  
131 , design type and the heterostructure materials. It was also a challenge to generate more sentences  
132 due to the cost implications of using GPT and the required labelling and post-processing efforts.  
133 Researchers can however use this dataset to augment and generate larger instruction datasets for QCL  
134 properties extraction from text. The quality of the generated sentences may not fully represent the  
135 nature of descriptions given in scientific articles detailing the QCL properties.

## 136 ETHICS STATEMENT

137 The authors have read and followed the ethical requirements for publication in Data in Brief. The  
138 current work does not involve human subjects, animal experiments, or any data collected from social  
139 media platforms. The scientific articles were sourced with permission from the publishers and with  
140 adherence to the text and data mining policies of the publishers.

## 141 CREDIT AUTHOR STATEMENT

142 **Deperias Kerre:** Resources, Investigation, Methodology, Data Preparation, Writing-Original Draft,  
143 Writing-Review and Editing

144 **Anne Laurent:** Methodology, Writing- Review and Editing, Supervision

145 **Kenneth Maussang:** Methodology, Writing- Review and Editing, Supervision

146 **Dickson Owuor:** Methodology, Writing- Review and Editing, Supervision

## 147 ACKNOWLEDGEMENTS

148 This study has been funded by the French Embassy in Kenya (Scientific and Academic Cooperation  
149 Department) and the CNRS (under the framework “Dispositif de Soutien aux Collaborations avec  
150 l’Afrique sub-saharienne”). The authors would also like to appreciate Strathmore university, School of  
151 Computing and Engineering sciences and the Doctoral academy for creating an opportunity for this  
152 work to be realized.

## 153 DECLARATION OF COMPETING INTERESTS

154 • The authors declare that they have no known competing financial interests or personal  
155 relationships that could have appeared to influence the work reported in this paper.

## 156 REFERENCES

- 157 1. Faist, J., Capasso, F., Sivco, D. L., Hutchinson, A. L., Sirtori, C., & Cho, A. Y. (1995). Quantum  
158 cascade laser: a new optical source in the mid-infrared. *Infrared Physics & Technology*, 36(1),  
159 99-103.
- 160 2. Haas, R., Keller, P., Hodges, J., and Spivak, J. (2014). Quantities, units, dimensions and data  
161 types ontologies (qudt). Technical report, Top Quadrant and NASA.
- 162 3. Li, H., Armiento, R., and Lambrix, P. (2020). An ontology for the materials design domain.  
163 *In The Semantic Web – ISWC 2020*, page 212–227.
- 164 4. Kerre, D., Laurent, A., Maussang, K., and Owuor, D. (2023). A concise ontological model of the  
165 design and optoelectronic properties in the quantum cascade laser domain. preprint.

- 166 5. Razavipour, S. G., Dupont, E., Chan, C. W. I., Xu, C., Wasilewski, Z. R., Laframboise, S. R., ... &  
167 Ban, D. (2014). A high carrier injection terahertz quantum cascade laser based on indirectly  
168 pumped scheme. *Applied Physics Letters*, 104(4).
- 169 6. Kerre, D., Laurent, A., Maussang, K., & Owuor, D., (2023, August). A text mining pipeline for  
170 mining the quantum cascade laser properties. In European Conference on Advances in  
171 Databases and Information Systems (pp. 393-406). Cham: Springer Nature Switzerland.
- 172 7. Kerre, D., Laurent, A., Maussang, K., & Owuor, D., (2024). An Instruction Dataset for Extracting  
173 Quantum Cascade Laser Properties from Scientific Text. Recherche Data Gouv: DOI:  
174 10.57745/U3U7XR.
- 175 8. Kumar, S., Hu, Q., & Reno, J. L. (2009). 186 K operation of terahertz quantum-cascade lasers  
176 based on a diagonal design. *Applied Physics Letters*, 94(13).
- 177 9. Williams, B. S., Callebaut, H., Kumar, S., Hu, Q., & Reno, J. L. (2003). 3.4-THz quantum  
178 cascade laser based on longitudinal-optical-phonon scattering for depopulation. *Applied*  
179 *Physics Letters*, 82(7), 1015-1017.
- 180 10. Hempel, M., Röben, B., Niehle, M., Schrottke, L., Trampert, A., & Grahn, H. T. (2017).  
181 Continuous tuning of two-section, single-mode terahertz quantum-cascade lasers by fiber-  
182 coupled, near-infrared illumination. *AIP Advances*, 7(5).
- 183 11. Wang, F., Slivken, S., Wu, D. H., Lu, Q. Y., & Razeghi, M. (2020). Continuous wave quantum  
184 cascade lasers with 5.6 W output power at room temperature and 41% wall-plug efficiency in  
185 cryogenic operation. *AIP Advances*, 10(5).
- 186 12. Wang, T., Liu, J. Q., Chen, J. Y., Liu, Y. H., Liu, F. Q., Wang, L. J., & Wang, Z. G. (2013).  
187 Continuous-wave operation of terahertz quantum cascade lasers at 3.2 THz. *Chinese Physics*  
188 *Letters*, 30(6), 064201.
- 189 13. Lü, X., Röben, B., Schrottke, L., Biermann, K., & Grahn, H. T. (2021). Correlation between  
190 frequency and location on the wafer for terahertz quantum-cascade lasers. *Semiconductor*  
191 *Science and Technology*, 36(3), 035012.
- 192 14. Khabibullin, R. A., Shchavruk, N. V., Pavlov, A. Y., Klochkov, A. N., Glinskiy, I. A., Tomosh, K.  
193 N., ... & Zhukov, A. E. (2019). Design and fabrication of terahertz quantum cascade laser with  
194 double metal waveguide based on multilayer GaAs/AlGaAs heterostructures. In *IOP*  
195 *Conference Series: Materials Science and Engineering* (Vol. 475, No. 1, p. 012020). IOP  
196 Publishing.
- 197 15. Ohtani, K., Turčinková, D., Bonzon, C., Benea-Chelmus, I. C., Beck, M., Faist, J., ... & Stutzki,  
198 J. (2016). High performance 4.7 THz GaAs quantum cascade lasers based on four quantum  
199 wells. *New Journal of Physics*, 18(12), 123004.
- 200 16. Wang, X., Shen, C., Jiang, T., Zhan, Z., Deng, Q., Li, W., ... & Duan, S. (2016). High-power  
201 terahertz quantum cascade lasers with ~ 0.23 W in continuous wave mode. *Aip Advances*, 6(7).
- 202 17. Valmorra, F., Scalari, G., Ohtani, K., Beck, M., & Faist, J. (2015). InGaAs/AlInGaAs THz  
203 quantum cascade lasers operating up to 195 K in strong magnetic field. *New Journal of*  
204 *Physics*, 17(2), 023050.
- 205 18. Kumar, S., Williams, B. S., Hu, Q., & Reno, J. L. (2006). 1.9 THz quantum-cascade lasers with  
206 one-well injector. *Applied Physics Letters*, 88(12).
- 207 19. Fang, L., Lee, G. G., & Zhai, X. (2023). Using gpt-4 to augment unbalanced data for automatic  
208 scoring. arXiv preprint arXiv:2310.18365.
- 209 20. Dai, H., Liu, Z., Liao, W., Huang, X., Cao, Y., Wu, Z., ... & Li, X. (2023). Auggpt: Leveraging chatgpt  
210 for text data augmentation. arXiv preprint arXiv:2302.13007.