

# Lerche: Generating data file processors in Julia from EBNF grammars

James R. Hester<sup>\*1</sup> and Erez Shinan<sup>2</sup>

<sup>1</sup> Australian Nuclear Science and Technology Organisation, Sydney, Australia <sup>2</sup> Independent researcher

DOI: [10.21105/joss.03497](https://doi.org/10.21105/joss.03497)

## Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Sebastian Benthall](#) ↗

## Reviewers:

- [@ziotom78](#)
- [@eschnett](#)

Submitted: 03 May 2021

Published: 20 July 2021

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

In a scientific context, structured text is commonly encountered in data files and domain-specific languages (DSLs) for handling data. Extended Backhaus-Naur Format (EBNF) ([ISO Central Secretary, 1996](#)) is a well-established standard for describing the syntax of such structured text. A useful subset of such grammars is known as LALR(1), meaning that the grammars describe text that can be unambiguously parsed based only on the tokens already seen and the next token of input ([DeRemer, 1969](#)). LALR(1) stands out for being able to parse most programming languages, while guaranteeing  $O(n)$  run-time complexity and very light memory use. Julia package *Lerche* automatically generates a parser that processes any data file or domain-specific language that can be described using a LALR(1) EBNF. The parse tree can be immediately transformed into application-specific data structures using user-supplied rules. This parser generator fills a gap in standards-based scientific work for the Julia ecosystem.

## Statement of Need

Standards for scientific data formats are essential for correct interpretation of data produced and consumed by parties that are not in direct contact. Standardisation is also becoming increasingly important as requirements for accessible and reusable data increase among funding bodies and publishers. Use of formal description languages, such as EBNF, remove ambiguity that may arise when using natural language descriptions for these standards. Furthermore, a parser that is machine-generated from an EBNF description is guaranteed to generate a correct parse tree from conformant data, unlike hand-written parsers. In addition, development of a formal scientific standard written using EBNF is aided by the availability of EBNF parsers that identify ambiguities and errors.

Scientific data readers typically form a small component of larger projects, which will dictate the programming environment and policies. This environment restricts the choice of EBNF parser generators; an EBNF parser generator that executes in the usual project environment is preferable to one that introduces new build-time dependencies and environments. *Lerche* was translated into Julia from the Python Lark project ([Shinan, 2021](#)) in order to remove the need for any extra language dependencies or build steps when developing standards-conformant data format readers for Julia projects. The generated parsers run within the Julia language environment, and are straightforward to integrate into larger Julia projects.

Other native parser tools available for Julia projects include *ParserCombinator* ([Cooke, 2021](#)), *Pegparser* ([Schneider, 2020](#)), and the built-in Julia macro system. None of these

<sup>\*</sup>Corresponding author

39 take EBNF as input. Reliance on EBNF parser generators external to the Julia environment,  
40 for example by calling Lark or pre-generating parsers, complicate distribution and create a  
41 barrier to user contribution to packages.

42 dREL (Spadaccini et al., 2012) is a relational data processing language for specifying crystallo-  
43 graphic algorithms. The dREL LALR(1) EBNF (Hester, 2021a) was developed using Lerche  
44 to verify its correctness and conformance to LALR(1) requirements. Package DrelTools (Hes-  
45 ter, 2021b) is built around the parser automatically generated by Lerche from this EBNF.  
46 Lerche is also used by CrystalInfoFramework (Hester, 2021c) to generate the parser for  
47 data files written in the Crystallographic Information Framework (CIF) format (Bernstein et  
48 al., 2016).

## 49 Features

50 Lerche optionally augments the range of applicable EBNF grammars by using contextual  
51 lexing. In this mode, only those tokens that are possible in the current parsing state are  
52 matched by the lexer, which can be useful for grammars in which certain character sequences  
53 match more than one type of token; for example, a keyword may also be a possible value for  
54 a plain sequence of characters in certain contexts, in which case a non-contextual lexer might  
55 wrongly fail to recognise the keyword.

56 The Lark grammars recognised by Lerche extend the EBNF standard in several ways. To  
57 aid in composability, they support templating, and importing rules and terminals from other  
58 grammars. To aid in refactoring, they support expressing rule semantics, which are translatable  
59 to common tree operations. For example, rules starting with `_` are considered to be auxiliary,  
60 and don't produce their own node. Terminals starting with `_` aren't included in the tree, which  
61 is often desired for punctuation such as commas and parentheses. In order to resolve possible  
62 ambiguities or conflicts, there is support for specifying priority in terminals and rules.

## 63 References

- 64 Bernstein, H. J., Bollinger, J. C., Brown, I. D., Gražulis, S., Hester, J. R., McMahon, B.,  
65 Spadaccini, N., Westbrook, J. D., & Westrip, S. P. (2016). Specification of the Crystal-  
66 lographic Information File format, version 2.0. *Journal of Applied Crystallography*, 49(1),  
67 277–284. <https://doi.org/10.1107/S1600576715021871>
- 68 Cooke, A. (2021). In *Github repository*. Github. [https://github.com/andrewcooke/](https://github.com/andrewcooke/ParserCombinator.jl)  
69 [ParserCombinator.jl](https://github.com/andrewcooke/ParserCombinator.jl)
- 70 DeRemer, F. L. (1969). *Practical translators for LR(k) languages* [PhD thesis]. MIT.
- 71 Hester, J. R. (2021c). In *Github repository*. Github. [https://github.com/jamesrhester/](https://github.com/jamesrhester/CrystalInfoFramework.jl)  
72 [CrystalInfoFramework.jl](https://github.com/jamesrhester/CrystalInfoFramework.jl)
- 73 Hester, J. R. (2021a). In *Github repository*. Github. [https://github.com/COMCIFS/dREL/](https://github.com/COMCIFS/dREL/blob/master/annotated-grammar.rst)  
74 [blob/master/annotated-grammar.rst](https://github.com/COMCIFS/dREL/blob/master/annotated-grammar.rst)
- 75 Hester, J. R. (2021b). In *Github repository*. Github. [https://github.com/jamesrhester/](https://github.com/jamesrhester/DrelTools.jl)  
76 [DrelTools.jl](https://github.com/jamesrhester/DrelTools.jl)
- 77 ISO Central Secretary. (1996). *Information technology — Syntactic metalanguage — Ex-*  
78 *tended BNF* (Standard ISO/IEC 14977:1996). International Organization for Standard-  
79 ization. <https://www.iso.org/standard/26153.html>
- 80 Schneider, A. (2020). In *Github repository*. Github. [https://github.com/abeschneider/](https://github.com/abeschneider/PEGParser.jl)  
81 [PEGParser.jl](https://github.com/abeschneider/PEGParser.jl)

- <sup>82</sup> Shinan, E. (2021). In *Github repository*. Github. <https://github.com/lark-parser/Lark>
- <sup>83</sup> Spadaccini, N., Castleden, I. R., Boulay, D. du, & Hall, S. R. (2012). dREL: A relational ex-  
<sup>84</sup> pression language for dictionary methods. *Journal of Chemical Information and Modeling*,  
<sup>85</sup> 52(8), 1917–1925. <https://doi.org/10.1021/ci300076w>

DRAFT