

UNIVERSITY OF COPENHAGEN
Computer Science Department
Data-Parallel Compilation
Lexical analysis & Syntax Tree Construction

Author: William Henrich Due (mcj284)
Advisor: Troels Henriksen
Submitted: 5th of April 2024

Abstract

The theory of a lexer generator creating data-parallel lexers is presented. Such a lexer generator is implemented and used to create a Lisp lexer which performance is compared with other Lisp lexers. Furthermore a description of how to extend an LLP parser generator such that it can construct concrete syntax trees is given. This is also implemented together with the lexer generator as part of Alpacc¹.

Contents

1	Introduction	2
2	Theory	2
2.1	Data-parallel Lexical Analysis	2
2.1.1	Data-parallel DFA Traversal	2
2.1.2	Data-parallel Tokenization	5
2.2	Concrete Syntax Trees	8
2.2.1	Concrete Syntax Tree Construction	9
2.2.2	Grammar Extension	13
3	Implementation	16
3.1	Tokenization	16
4	Tests	17
4.1	Tokenization	17
4.1.1	General Testing	17
4.1.2	Large LISP Input	18
4.2	Parent Vector	18
5	Benchmarks	19
5.1	Lisp lexer	19
6	Discussion	21

¹<https://github.com/diku-dk/alpacc>

1 Introduction

2 Theory

Hills paper “Parallel lexical analysis and parsing on the AMT distributed array processor” [2] describes a method to obtain the path in a deterministic finite automaton (DFA) given a input string. It does not describe how to create such a DFA but it describes how to simulate it in a work-efficient data-parallel manner. This section will describe the theory of this method and describe a method of using it for tokenization.

2.1 Data-parallel Lexical Analysis

2.1.1 Data-parallel DFA Traversal

To explain the theory of parallel lexical analysis we first remind the reader of the definition of a deterministic finite automaton.

Definition 2.1 (DFA). A deterministic finite automaton [3] [7] is given by a 5-tuple $(Q, \Sigma, \delta, q_0, F)$ where.

1. Q is the set of states where $|Q| < \infty$.
2. Σ is the set of symbols where $|\Sigma| < \infty$.
3. $\delta : \Sigma \times Q \rightarrow Q$ is the transition function.
4. $q_0 \in Q$ is the initial state.
5. $F \subseteq Q$ is the set of accepting states.

This definition is fine as is but we will need to reformualte it to develop data-parallel lexical analysis. We would want the definition to use a curried transition function. But for this to hold then the DFA would also have to be total.

Definition 2.2 (Total DFA). A DFA $(Q, \Sigma, \delta, q_0, F)$ is said to be total if and only if

$$\forall(a, q) \in \Sigma \times Q : \delta(a, q) \in Q$$

If a DFA is total we may use a curried transition function $\delta : \Sigma \rightarrow Q \rightarrow Q$.

This is needed since else the the function would not be fully defined in the domains Σ and Q .

The reason for doing so is because if we have any two functions $g = \delta(a)$ and $f = \delta(a')$ then it follows from composition that.

$$g(f(q)) = (g \circ f)(q)$$

This allows for an alternative way of determining if a string can be produced by an DFA. Instead of first evaluating $f(q)$, then $g(f(q))$ and then checking if this state is a member of F . We could instead partially apply δ to the symbols and then compose them to a single function which could be used to determine if a string is valid. This sets the stage for data-parallel lexing, we want to find a way to make the problem into a **map-reduce**. We want to do this because it can be computed using a data-parallel implementation unlike the normal way of traversing a DFA.

For the ability to use a data-parallel **map-reduce** we must have a monoidal structure. Here Δ is the set of all the composed partially applied delta functions needs to be closed under function composition.

Remark. A partially applied delta function is shortened to *PADF*. If multiple PADF's are composed then they are a *composed PADF*.

Proposition 2.1 (DFA Composition Closure). Given a total DFA then the set of composed PADF's $\Delta : Q^Q$ will be closed under composition. The set Δ is the set Δ_i in the recurrence relation with the smallest i such that $\Delta_i = \Delta_{i+1}$.

$$\begin{aligned}\Delta_1 &= \{\delta(a) : a \in \Sigma\} \\ \Delta_{i+1} &= \Delta_i \cup \{f \circ g : f, g \in \Delta_i\}\end{aligned}$$

Proof. We will circit by showing that a solution Δ exists. First note that the cardinality is monotonically increasing i.e. $\Delta_i \subseteq \Delta_{i+1}$ since Δ_{i+1} is a union of Δ_i and another set. Secondly note that since $|Q| < \infty$ then a finite amount of functions of the form $Q \rightarrow Q$ can exists. Since the set is bounded and increasing then at some point $\Delta_i = \Delta_{i+1}$ and the smallest i where it holds is the solution Δ .

For Δ to be closed under composition, then for arbitray $f, g \in \Delta$ it must hold that $f \circ g \in \Delta$. Since Δ_1 is the set of PADF's that constructs Δ and composition is associative then all elements of Δ can be expressed of the form.

$$\delta(a_1) \circ \dots \circ \delta(a_n) \in \Delta$$

Given an expression of the aforementioned form we could construct a sequence of PADF's.

$$\delta(a_1), \dots, \delta(a_n)$$

So we would have a surjective map from sequences to Δ , since some sequences would map to the same composed PADF. So if all permutations with replacement of Δ_1 of any sequence length are members of Δ then Δ would be closed under composition. Furthermore, since it was shown that Δ is finite so the sequences at some point $\Delta_i = \Delta_{i+1}$ would only add new sequences but no new composed PADF's. Therefore it suffices to show that if all sequences of length k where $1 \leq k \leq i$ is a subset of Δ_i then Δ is closed under composition. This can be shown using a proof by induction.

Base: Δ_1 trivially holds since it only contains sequences of length one and they are the initial PADF's.

Step: Given Δ_i contains every sequence of length i or less then we to show this implies that Δ_{i+1} will contain every sequence of length $i+1$ or less.

By the induction hypothesis Δ_{i+1} must contain every sequence of length i or less due to $\Delta_i \subseteq \Delta_{i+1}$. It remains to show that every sequence of length $i+1$ is a member of Δ_{i+1} . It is known that a direct product of Δ_i is used in the definition of Δ_{i+1} so $\{f \circ g : f, g \in \Delta_i\} \subseteq \Delta_{i+1}$. A direct product between sequences of length 1 and i will create every sequence of length $i+1$ and therefore every sequence of length $i+1$ is a member of Δ_{i+1} . Thereby Δ is closed under composition. \square

Now since Δ is closed under composition then it follows that Δ and function composition induces a monoidal structure.

Corollary 2.1 (DFA Composition Monoid). DFA composition closure induces a semigroup which in turn induces the monoid $(\Delta \cup \{id\}, \circ)$ where $id : Q \rightarrow Q$ and $id(q) = q$.

Knowing this we can establish the following algorithm

Algorithm 2.1 (Data-parallel String Match). It can be determined in $O(n)$ work and $O(\log n)$ span if a string can be produced by a DFA. First construct the total DFA $(Q, \Sigma, \delta, q_0, F)$ from the DFA.

1. Partially apply δ to every symbol in the input string such that it becomes an array of PADF's.

$$\mathbf{map} \, \delta [a_1, a_2, \dots, a_{n-1}, a_n] = [\delta(a_1), \delta(a_2), \dots, \delta(a_{n-1}), \delta(a_n)]$$

-
2. Reduce the PADFs into a single composed PADF $\delta' : Q \rightarrow Q$.

$$\text{reduce}(\circ) \text{ id } [\delta(a_1), \delta(a_2), \dots, \delta(a_{n-1}), \delta(a_n)] = \delta'$$

3. Evaluate $\delta'(q_0)$ and determine if $\delta'(q_0) \in F$.

2.1.2 Data-parallel Tokenization

For data-parallel tokenization we need to extend data-parallel algorithm 2.1 will be needed to be extended. The idea will be to use a data-parallel `map-scan` instead since it will give all the states. This is also the methods described in Hills [2] paper. The problem is we need to be able to recognize the longest stretch of symbols that results in a token. And we also need to restart the traversal of DFA if a final state is hit while no options to traverse further. To do so we first need a function to define a function to recognize tokens.

Definition 2.3 (Token Function). Given a DFA and a set of tokens T . The token function $\mathcal{T} : F \rightarrow T$ is a function that maps accepting states to some token.

We will also need a single state to point to which is the dead state. This will become useful when the DFA needs to be traversed multiple times. Since we will need to be able to recognize when the end of a traversal is reached and we have to restart.

Definition 2.4 (Total DFA with a Dead State). Given a DFA it is made total with a dead state by defining a new set of states $Q' = Q \cup \{q_\perp\}$ where q_\perp is the dead state. Additionally a new transition function δ' is defined.

$$\delta'(a, q) = \begin{cases} q_\perp & (a, q) \notin \text{dom}(\delta) \\ \delta(a, q) & \text{otherwise} \end{cases}$$

Now that we have a definition of a DFA where the dead state is known another problem is needed to be solved. The problem is as mentioned before that the traversal of the DFA has to be restarted if a dead state is reached after an accepting state. This is done using the following binary operation.

Definition 2.5 (Make PADF Safe Under Composition). Given a DFA with two PADFs $f = \delta(a)$ and $g = \delta(b)$. The operation \oplus makes f safe to compose such that it will continue traversing the DFA.

$$(f \oplus g)(x) = \begin{cases} q_0 & f(g(x)) = q_\perp \wedge g(x) \in F \\ g(x) & \text{otherwise} \end{cases}$$

This definition will make every possible final state become the initial state. This forgets the final state but it allows for the traversal to continue. The forgetfulness is not a problem this can be solved by looking at the previous PADF. This puts a limit on the lexer which is it only allows for going back to the previous state if a dead state is hit. Using this and previous definitions can be put together to the following algorithm.

Algorithm 2.2 (Data-parallel Tokenization). Given a total DFA with a dead state where $q_0 \notin F$ and a token function $\mathcal{T} : F \rightarrow T$. A string can be tokenized in $O(n)$ work and $O(\log n)$ span.

1. Let $s = [a_1, a_2, \dots, a_{n-1}, a_n]$ be a string that will be tokenized then partially apply δ to every symbol.

$$\mathbf{map}(\delta) s = x$$

2. Make every PADF safe for composition beside for the last PADF.

$$\mathbf{map}(\oplus) (\mathbf{init} x) (\mathbf{tail} x) ++ [\delta_n] = [\delta_1 \oplus \delta_2, \dots, \delta_{n-1} \oplus \delta_n, \delta_n] = y$$

3. Do a scan to get every composition.

$$\mathbf{scan}(\circ) id y = z$$

4. Compute the actual state and determine weather it is an end state².

```

let  $f = \lambda i \rightarrow$ 
  let  $s = \mathbf{if} i = 0 \mathbf{then} x[i](q_0) \mathbf{else} (z[i - 1] \circ x[i])(q_0)$ 
  in  $(i = n - 1 \vee (z[i](q_0) = q_0 \wedge s \in F), s)$ 
in map  $f$  (iota  $n$ )

```

5. Remove every state that is not an ending state.

$$\mathbf{filter}(\lambda(b, s) \rightarrow b)$$

6. Assert that the ending state is an accepting state.

$$(\lambda(b, s) \rightarrow s \in F) \circ \mathbf{last}$$

²If you were to also keep track of the index then the span of each token could also be found.

7. Produce the token sequence.

$$\mathbf{map} (\lambda(b, s) \rightarrow \mathcal{T}(s))$$

We will give an explanation of why the algorithm produces an array of paths where each path starts in q_0 and is not included in the path. Knowing where these paths ends it can then be trivially used to construct the token sequence using the token function \mathcal{T} . We will also show that it can be correctly determine whether the DFA paths taken are correct.

We know in step 1-2. every PADF is created and is made safe so we know at this point either it maps to some original states or some maps to q_0 because the composition are safe. At step 3. a prefix sum using function composition is computed of these composed safe PADFs. Which by Algorithm 2.1 can give us all states traversed in the DFA. But due to \oplus is used then an accepting state followed by an dead states will become the initial state q_0 . This almost gives all the paths, since now when a final state in a path is reach it becomes the initial state and the traversal is restarted.

Then at step 4. the paths taken in the DFA from state q_0 is found. We know for a given state s in a path it will can be one of these cases.

- If $i = 0$ then $x[i](q_0)$ will be the first state visited after q_0 by definition of the transition function.
- If $i \neq 0$ then s will be $(z[i - 1] \circ x[i])(q_0)$. $z[i - 1](q_0)$ could wrongly map to q_0 by definition of making the composition safe.
 - If $z[i - 1](q_0)$ maps wrongly to q_0 then the traversal of the DFA will be reset and $(z[i - 1] \circ x[i])(q_0)$ will map to the correct state of a new path.
 - If $z[i - 1](q_0)$ maps to the correct state then $(z[i - 1] \circ x[i])(q_0)$ will also map to the correct state.

It can be determined where each path ends by the following predicate.

$$i = n - 1 \vee (z[i](q_0) = q_0 \wedge s \in F)$$

We know if $i = n - 1$ then it must be the last state of the last path. The second case is if $z[i](q_0) = q_0$ and $s \in F$ holds. We know if $z[i](q_0) = q_0$ holds then it is possible a final state was mapped to q_0 by composition being safe. And if $s \in F$ then since $q_0 \notin F$ we know $z[i](q_0) \neq s$ meaning the reason $z[i](q_0) = q_0$ must be because of a reset by composition being safe meaning then it must be a final state in a path that is not the last path.

Step 5-7. will first keep every last state by this predicate. Afterwards it is asserted that the last state $s \in F$. Only the last state of the last path is needed since all previous path would had been valid else it would had not been reset by composition being safe. Finally just map every state to the token it has.

Example 2.1 (Problems with this method). If we had the following three patterns.

$$\tau_1 = a(a \quad \tau_2 = a+ \quad \tau_3 = a)a$$

Then the algorithm would not be ble to determine the tokens for the following input.

$$a(aaa)a$$

The token sequence of this input could be determined to be τ_1, τ_2, τ_3 with spans $(0, 3), (3, 4), (4, 7)$ if you reset at the last acceptings state. But since the algorithm resets at the last state from an accepting state then it fails. Meaning it would be able to recongnize τ_1 . The problem is at the token τ_2 would get the span $(3, 5)$ and then the algorithm would not recongnize any token with the pattern $)a$.

There exists an alternative to this method given in Robin Voetters masterthesis [6, pp. 13–16]. This is done by constructing a NFA for each tokens regular expression pattern. Then create a new NFA where there is an epsilon transition from the initial state to every tokens NFAs inital state. Then you would convert this NFA into a DFA where you would have a function \mathcal{T} which maps final state to a token. Then for every outgoing transition $\delta(q_0, t) = q$ of the DFA then we add a transition from a final state $q_s \in F$ such that if $(q_s, t) \notin \text{dom}(\delta)$ then $\delta(q_s, t) = q$. If any of these new transitions are visited during the traversal of the DFA then a token is produced. We also have to account for the first transition visited which also will be a producing a token.

2.2 Concrete Syntax Trees

The alpacc parser generator is a parser generator which creates parser for the LLP grammer class. The LLP parsers described in the LLP parsing paper [5] only touches upon the verification and production sequence construction of parsing not the creation of concrete the syntax tree. This would also had to be done in a data-parallel manner else we would run into a bottle-neck. A method for this is described in Voetters master thesis [6, pp. 32–33] about compilation on the GPU. Here it is described how to create a concrete syntax

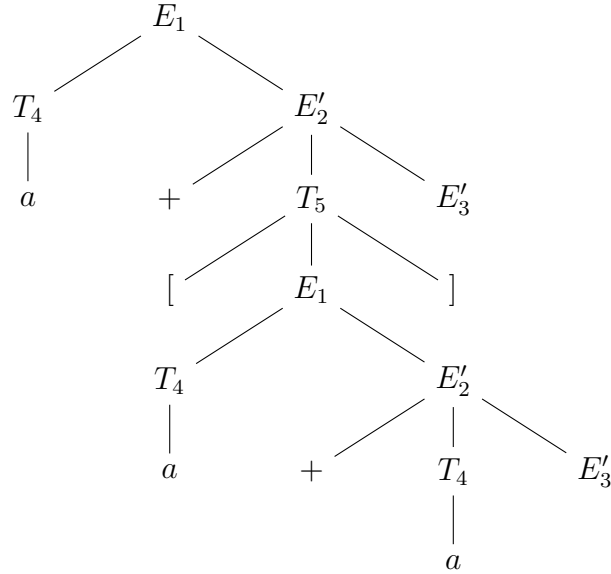
tree using the representation from Hsu PhD dissertation [4, pp. 77–81] where you have a parent that can be used for multiple transformations of a syntax tree. The representation of the tree will be an array of nodes where each node has the index to its parent node. The nodes will also be given in a preorder traversal of the tree. This kind of representation is much more fit for data parallelism than some recursive datatype.

2.2.1 Concrete Syntax Tree Construction

So as an example we wish to construct syntax tree for the grammar.

- 1) $E \rightarrow TE'$ 2) $E' \rightarrow +TE'$ 3) $E' \rightarrow \varepsilon$ 4) $T \rightarrow a$ 5) $T \rightarrow [E]$

Then for the input $a + [a + a]$ we would want to construct the following tree.



The first problem to tackle is terminals are not included in the left parse returned by the parser, currently it only returns the productions.

$$[E_1, T_4, E'_2, T_5, E_1, T_4, E'_2, T_4, E'_3, E'_3]$$

To create the syntax trees the first problem is how to weave together the production sequence with the sequence of terminals. This is done by extending the grammar such that each terminal is associated with a unique production. That way a production in the production sequence will have a direct correspondent to a terminal. This can then be use when constructing the syntax tree.

$$[E_1, T_4, A_a, E'_2, A_+, T_5, A_[, E_1, T_4, A_a, E'_2, A_+, T_4, A_a, E'_3, A_], E'_3]$$

Here A_t is a production that produces the single terminal it represents. Then we can just map everyone of these productions to its terminal like so.

$$[E_1, T_4, a, E'_2, +, T_5, [, E_1, T_4, a, E'_2, +, T_4, a, E'_3,], E'_3]$$

A concern with this method is if this changes the number of grammars that can be parsed by LLP or if it changes the size of the table used for LLP parsing. These problem does not exists for this extension and it is later proven these problems does not exists.

The parent vector can then be constructed with the following algorithm.

Algorithm 2.3 (Parent Vector). Given a left parse a_0 then we can construct the the parent vector as followed.

1. For each production map it to its arity i.e. the number of nonterminals in the right-hand side of the production.

$$a_1 = \mathbf{map} \text{ arity } a_0$$

Here *arity* is a function that maps productions to its arity.

2. Substract one from each arity.

$$a_2 = \mathbf{map} (-1) a_1$$

3. Do an exclusive scan using addition on the arities.

$$a_3 = \mathbf{scan} (+) a_2$$

4. Find the index of the previous smaller or equal element for each element from the previous step.

$$a_4 = \mathbf{map} \text{ prev } a_3$$

Here *prev* is a function that finds the index of the previous smaller or equal element in a_3 .

5. Set the first elements parents to be it self. $a_5 = a_4[0] = 0$

Using this algorithm we can show the computation of parent vector for the example.

Arity	Stack Change	Stack Depth	Parent Index	Node	Index
2	1	0	0	E_1	0
1	0	1	0	T_4	1
0	-1	1	1	a	2
3	2	0	0	E'_2	3
0	-1	2	3	$+$	4
3	2	1	3	T_5	5
0	-1	3	5	$[$	6
2	1	2	5	E_1	7
1	0	3	7	T_4	8
0	-1	3	8	a	9
3	2	2	7	E'_2	10
0	-1	4	10	$+$	11
1	0	3	10	T_4	12
0	-1	3	12	a	13
0	-1	2	10	E'_3	14
0	-1	1	5	$]$	15
0	-1	0	3	E'_3	16

Table 1: Overview of how the parent vector is constructed.

Note how the arity for the original productions end up being the number of symbols in the right-hand side while each terminal symbol ends up being zero. This is due to the extension done to the original grammar.

Step 4. in algorithm 2.3 could be solve in different ways. A possible way to do it in a data-parallel manner is doing a linear backwards search. This would lead to $O(n^2)$ work and $O(n)$ span. This would be a pretty huge bottle neck-compared to the $O(n)$ work and $O(\log n)$ span for both tokenization and parsing. So a method for finding the previous smaller or equal element is used. This can be done in $O(n \log n)$ work and $O(\log n)$ span as described in Voettters masterthesis [6, pp. 32–33].

Algorithm 2.4 (Previous Smaller or Equal Element). Given an array of a_0 where $|a_0| = 2^n$ and the array contains some set of elements A where (\leq, A) is a total order.

-
1. Construct tree of minima.

```

let create arr =
  if  $|arr| = 1$ 
  then arr
  else let arr' = map ( $\lambda i \rightarrow \min arr[2i] arr[2i + 1]$ ) (iota  $|arr|$ )
    in create arr' ++ arr

```

This will compute a binary tree where index 0 of the array is the root. Furthermore for a node at index i you can find its parent using $\lfloor \frac{i-1}{2} \rfloor$, its left child $2i + 1$ and its right child $2i + 2$ assuming they exists.

let *tree* = *create a₀*

2. To find the previous smaller or equal element of a value v with index i . We will first find the common ancestor of the elements from i some previous smaller or equal element that may exists.

```

let ascent j =
  if  $j \neq 0 \wedge (left\ j \vee tree[sibling\ j] > v)$ 
  else ascent (parent j)
  then j

```

Here *left* is a predicate that determines if j is the index of a left node i.e. if it is odd. While *sibling* determines the neighbouring node. Then compute the index $k = ascent(2^n + i)$, if $k = 0$ then the previous smaller or equal element does not exists.

3. Now descent to the right most element that is smaller or equal to v and adjust for the offset.

```

let descent j =
  if  $j < 2^n$ 
  else if  $tree[2j + 2] \leq v$  then descent (2j + 2) else descent (2j + 1)
  then  $j - 2^n$ 

```

Note that the use of recursion can be easily rewritten to loops.

2.2.2 Grammar Extension

We will first need a precise definition of the grammar extension. This is needed to show it does not lead to fewer grammars being LLP or larger LLP tables.

Given a grammar (N, T, P, S) where N is the set of nonterminals, T is the set of terminals, P is the set of productions, and S is the initial nonterminal. We first construct the following set of nonterminals.

$$N_T = \{A_t : t \in T\}, \text{ where } A_t \notin N \text{ for any } t \in T$$

And we also create an secondary set of productions.

$$P_T = \{A_t \rightarrow t : t \in T\}, \text{ where } A_t \rightarrow t \notin P \text{ for any } t \in T$$

Now we create a new set of productions by replacing every terminal in the right-hand side of a production by its corresponding nonterminal.

$$P_r = \{A \rightarrow \delta[A_{t_1}/t_1, \dots, A_{t_n}/t_n] : A \rightarrow \delta \in P\}$$

The final grammar that will be used for parsing is now.

$$G' = (N', T', P', S') = (N \cup N_T, T, P_r \cup P_T, S)$$

First we will need to show that any LL grammar extended can still be parsed using LL parsing. The definition of a $LL(k)$ ³ [1] parsing table is as followed.

Definition 2.6 ($LL(k)$ table). Let $G = (N, T, P, S)$ be a $LL(k)$ context-free grammar and $\tau : N \times T^* \rightarrow \mathbb{P}(\mathbb{N})$ denote the $LL(k)$ table. For a given production $A \rightarrow \delta = p_i \in P$ where $i \in \{0, \dots, |P| - 1\}$ is a unique index.

$$i \in \tau(A, s) \text{ where } s \in \text{FIRST}_k^G(\delta) \bullet^k \text{FOLLOW}_k^G(A)$$

If for a given grammar it holds that $|\tau(A, s)| \leq 1$ for all $(A, s) \in N \times T^*$ then the grammar is $LL(k)$.

Proof. Given a $LL(k)$ grammar G , then for a production $A \rightarrow \delta \in P$ and its corresponding extended production $A \rightarrow \delta[A_{t_1}/t_1, \dots, A_{t_n}/t_n] = A' \rightarrow \delta' \in P_r$ we have.

$$\text{FIRST}_k^G(\delta) \bullet^k \text{FOLLOW}_k^G(A) = \text{FIRST}_k^{G'}(\delta') \bullet^k \text{FOLLOW}_k^{G'}(A)$$

³I do not know where else to cite from. The definition is found in my Bachelor thesis and is partly based on other definitions.

The first and follow-sets must be unchanged for productions in P_r because all added nonterminals directly becomes their corresponding terminals in a derivation. Therefore none of these productions leads to a table conflict. Then for a production $A_t \rightarrow t \in P_T$ we have.

$$i \in \tau(A_t, s) \text{ where } s \in \text{FIRST}_k^{G'}(t) \bullet^k \text{FOLLOW}_k^{G'}(A_t)$$

Since A_t only has a single production rule then it will always hold that $|\tau(A_t, s)| \leq 1$. Meaning the extended grammar G' will also be a $\text{LL}(k)$ grammar. \square

Something to note is this extension does lead to larger LL tables but this is not the case for LLP. So now we will show this extension does not lead to larger LLP tables and does not change if a LLP grammar can be parsed. To do so we need to be reminded of the PSLS⁴ [5] definition.

Definition 2.7. PSLS Let $G = (N, T, P, S)$ be a context-free grammar. The function $\text{PSLS}(x, y)$ for a pair of strings $x, y \in T^*$ is defined as follows:

$$\begin{aligned} \text{PSLS}(x, y) = & \{ \alpha : \exists S \Rightarrow_{lm}^* wuA\beta \Rightarrow wxB\gamma \Rightarrow^* wxy\delta, \\ & w, u, y' \in T^*, a \in T, A, B \in N, \alpha, \beta, \gamma, \delta \in (N \cup T)^*, u \neq x, \\ & y = ay', \alpha \text{ is the shortest prefix of } B\gamma \\ & \text{such that } (y, \alpha, ()) \vdash^* (y', \omega, \pi) \} \\ \cup & \{ a : \exists S \Rightarrow^* wuA\beta \Rightarrow wxa\gamma \Rightarrow^* wxy\delta, \\ & a = \text{FIRST}_1^G(y), w, u \in T^*, \beta, \gamma, \delta \in (N \cup T)^*, u \neq x \} \end{aligned}$$

We will also need the following Lemma.

Lemma 2.1 (Grammar Extension Bijection). Let G be an LL grammar and let G' be the extended grammar of G . If we have the following two leftmost derivations from each grammar $S \Rightarrow_{lm}^* wA\delta \Rightarrow w\gamma$ and $S' \Rightarrow_{lm}^* wA'\delta \Rightarrow w\gamma'$ then the following derivations hold.

- $S \Rightarrow_{lm}^* w\gamma'[A_{t_1}/t_1, \dots, A_{t_n}/t_n]$
- $S' \Rightarrow_{lm}^* w\gamma[t_1/A_{t_1}, \dots, t_n/A_{t_n}]$

Hence there exists a bijection between these two derivations.

⁴This is a slightly modified definition by Vladislav Vagner because the old definition in the paper does not work.

Proof. We know that G and G' produces the same language and due to them both being LL then there is only one way of deriving their leftmost derivations. It is also known that since A is a nonterminal in G then it must not be a A_t nonterminal in G' else they would not share a common terminal string w . So the derivation of A' must use a production of the form $A' \rightarrow \alpha[A_{t_1}/t_1, \dots A_{t_n}/t_n] \in P_r$ where $A \rightarrow \alpha \in P$. Then we know that when the single derivation on A' is performed then not a single A_t nonterminal in γ' has been derived to its terminal t . Therefore it must be the case that $\gamma = \gamma'[A_{t_1}/t_1, \dots A_{t_n}/t_n]$. Likewise A' has just been derived so no A_t nonterminals have been derived meaning $\gamma' = \gamma[t_1/A_{t_1}, \dots t_n/A_{t_n}]$. \square

Now we will be able to show what we wanted to show.

Proof. The proof will be done by showing that given a LLP grammar G i.e. for any admissible pair (x, y) we have $|\text{PSLS}(x, y)| = 1$. Then it will follow that for its corresponding extended grammar G' we must also have $|\text{PSLS}(x, y)| = 1$.

- If we are given an pair (x, y) where $|\text{PSLS}(x, y)| = 1$ for some LLP grammar G and we have some derivation.

$$S \Rightarrow_{lm}^* wuA\beta \Rightarrow wxB\gamma \Rightarrow^* wxy\delta$$

Where there only exists a single shortest prefix α of $B\gamma$ such that $(y, \alpha, ()) \vdash^* (y', \omega, \pi)$ and for any a that may be in the second set we have $a = \alpha$.

Given the extended grammar G' , then for the pair (x, y) to be in the first set of the PSLS definition then it must be of the form.

$$S' \Rightarrow_{lm}^* wu'A_tB'\gamma' \Rightarrow wxB'\gamma' \Rightarrow^* wxy\delta'$$

Both of the leftmost derivations will have the common subderivation wx since they still produce the same language and they both must produce x .

By Lemma 2.1 we have a bijection between derivations of $wxB'\gamma'$ and $wxB\gamma$ then there must also exist a single shortest prefix α' of some $B'\gamma'$ such that $(y, \alpha', ()) \vdash^* (y', \omega', \pi')$. Since if there were multiple different shortest prefixes of some $B'\gamma'$ then so would $B\gamma$.

- If we are given an pair (x, y) where $|\text{PSLS}(x, y)| = 1$ for some LLP grammar G and we have some derivation.

$$S \Rightarrow^* wuA\beta \Rightarrow wxa\gamma \Rightarrow^* wxy\delta$$

Where there only exists a single $a = \text{FIRST}_1^G(y)$ where for any α that may be in the first set we have $a = \alpha$.

Then for the extended grammar G' then a derivation of the following form must also exist.

$$S' \Rightarrow_{lm}^* wu'A_tB'_a\gamma' \Rightarrow wxB'_a\gamma' \Rightarrow^* wxy\delta'$$

By Lemma 2.1 we have a bijection between derivations of $wxB'_a\gamma'$ and $wxa\gamma$. Knowing this and $B'_a \Rightarrow a$ then there only exist a single $a = \text{FIRST}_1^{G'}(B'_a)$ for this form of derivation.

We would now need to show that B'_a would exist in the first set of the PSLS definition since it cannot exist in the second set. We see the derivation is of the same form as a derivation as in the first set. We also see the shortest prefix α' of $B'_a\gamma'$ is such that $(y, \alpha', ()) \vdash^* (y', \varepsilon, \pi')$ is B'_a . By $a = \text{FIRST}_1^{G'}(B'_a) = \alpha'$ then α' is the only such prefix.

Thereby the extended grammar is also LLP and since the bijection grammar and its extended version then we know the PSLS table does not get larger. \square

3 Implementation

3.1 Tokenization

To implement tokenization we run into the problem that we want to represent an partially applied delta function in Futhark. Futhark does not allow for reposing functions as an array element. Even if it was possible then it might not be the most efficient representation. Instead if everyone of of states in the DFA is assigned a unique number then we could represent it as an array. So `delta[i] = j` would mean the partially applied function `delta` maps `i` to `j`. We could then also compose this function like so.

```
let compose a b = map(\i -> a[i]) b
```

This was at some time the implementation of `alpacc` but it is possible to go faster. You can enumerate every composed partially applied delta function and create a table of compositions. Then your function composition would become.

```
let compose a b = composition_table[a, b]
```

The way you would then enumerate every composed partially applied delta function is by first constructing the set of such functions and then assign each some number. A way to construct this set is by Proposition 2.1. You

would just have to keep track of the composed PADF's used to construct each composed PADF to make a table. Something one might notice is it is actually computationally wasteful since transitions and subpaths that are not neighbouring will result in a valid path in the DFA. So as an alternative one could only compose the composed PADF that are neighbouring each other. Then when you compose two transitions then you would need to add them as new transitions. Then for all the composed PADF's found you would have to define a dead composed PADF. It can then be used when some composition of existing composed PADF's are not defined in the table.

So if you had the following pattern you wanted to match abc then a is a neighbor to b and b is a neighbor to c . So initially we have a set of PADF's.

$$\{\delta(a), \delta(b), \delta(c)\}$$

Then we compose both $\delta(a)$ with $\delta(b)$ and $\delta(b)$ with $\delta(c)$ since both pairs are neighbouring transitions. Then the set of the new composition together with the old compositions results in the set.

$$\{\delta(a), \delta(b), \delta(c)\} \cup \{\delta(a) \circ \delta(b), \delta(b) \circ \delta(c)\}$$

Now $\delta(a) \circ \delta(b)$ is composed with $\delta(c)$ and $\delta(a)$ is composed with $\delta(b) \circ \delta(c)$. Here both $\delta(a) \circ \delta(b)$ and $\delta(b) \circ \delta(c)$ are recomputed. The set of these compositions are then unioned with the old set.

$$\{\delta(a), \delta(b), \delta(c), \delta(a) \circ \delta(b), \delta(b) \circ \delta(c)\} \cup \{\delta(a) \circ \delta(b) \circ \delta(c)\}$$

This ends up being the final set. This can also be done on the set of compositions that were made safe. Currently this implementation is extremely slow, it does some recomputations which could be mitigated using memoization or removing connections that have been computed before.

4 Tests

4.1 Tokenization

4.1.1 General Testing

To test the tokenization it is compared with a sequential version. The sequential version works much like the data-parallel version. It will traverse the DFA until it reaches a dead state. Afterwards it will determine if the previous state was accepting if so then reset the traversal else the input could not be tokenized.

The testing strategy is then to choose some set of patterns and strings of interest. Then compare if the sequence of tokens with their spans are the same. One of the set of patterns that was deemed interesting are the one given in Example 2.1.

$$\tau_1 = a(a \quad \tau_2 = a+ \quad \tau_3 = a)a \quad \tau_4 = \backslash s \backslash n \backslash t \backslash r$$

Here something like $a(aa)a$ can get recongnized while $a(aaa)a$ can not get recongnized. This test is to make sure concat, alternation, and the one-or-more pattern works as expected. This is tested multiple times using other test cases also. What is more important, it also assures that $a(aaa)a$ will fail. Another pattern of interest is.

$$\tau_1 = \mathbf{if} \quad \tau_2 = \mathbf{then} \quad \tau_3 = \mathbf{else} \quad \tau_4 = [a-z] + \quad \tau_5 = \backslash s \backslash n \backslash t \backslash r$$

Here we want the key words **if**, **then**, and **else** to be associated with a token before it is associated with τ_4 . Here the order the tokens are defined are the order they are prioritized from most important to least. There were also some other interesting patterns and inputs but these examples gives a good idea of cases tested for. ⁵

A problems with this method of testing is it is limited to single examples and what properties the author finds to be interesting. There might be edge cases that were missed and not accounted for since they did not seem to be of interest.

4.1.2 Large LISP Input

Since the general testing is only on small inputs of the DFA it would also be good to test on large randomized inputs. This would give greater confidence in the lexer generator works as expected. This is done by a LISP lexer is compared with a LISP lexer made with the Logos lexer library. Here a large randomized input that could be tokenized by a LISP lexer is created randomly.

4.2 Parent Vector

To assert the correctness of the parent vector construction the algorithm for finding the previous equal or smaller element is tested. The reason for only testing this is the steps 1-3. in algorithm 2.3 do offer much complexity. A problem with this testing method is it assumes the correctness of the parent vector algorithm. Also something to consider is why not testing that the

⁵The different cases tested for can be found [here](#).

concrete syntax trees correspond to a sequential version. This is because the parent vector is the last step in testing since the production sequence is already tested in Alpacc [1].

The strategy for testing is using property based testing. This is done by property based testing. The property tested for is given an array a of 32-bit integers then f maps algorithm 2.4 over a . Then the result of $f(a)$ should be the same as if g . Here g is an implementation that for every element in the the array does a backwards linear search to find the previous smaller or equal element. Then the follow property should be fulfilled.

$$\forall a. f(a) = g(a)$$

Since not all a can be tested then a is randomized a given amount of times. Also note if the previous or smaller element does not exists then f and g will map to -1 .

5 Benchmarks

5.1 Lisp lexer

To assert the performance of a data-parallel some benchmarks of Lisp lexers are compared. These Lisp lexers understand the following tokens patterns.

$$\tau_1 = (\quad \tau_2 =) \quad \tau_3 = [a-zA-Z0-9] + \quad \tau_4 = (\backslash s | \backslash n | \backslash t | \backslash r) +$$

Here τ_1 and τ_2 are the open and closing parenthesis in S-expressions. τ_3 are the atoms while τ_4 is the spacing between the important tokens for the grammar τ_1 , τ_2 , and τ_3 .

The benchmarks are of three lexers in different languages and described below

- Rust: A lexer made with the Logos library which saves every token in memory which must be done for the data-parallel lexer.
- C: A lexer written by Troels Henriksen.
- Futhark: A data-parallel lexer from Alpacc [1] using the method described.

The benchmarks are measured using a Ryzen 5 1600x CPU and a Nvidia RTX 3060 GPU. Which may not seem like a fair comparison since the CPU

was released in 2017 and the GPU was released in 2021.⁶ We will just have to assume this is comparable since single threaded performance has stalled the past 10 years.

Input Size (MiB)	Rust (μs)	C (μs)	Futhark (μs)
10	44,866	24,440	3,388
20	86,898	47,876	6,960
30	131,183	70,986	10,667
40	180,367	98,906	14,400
50	214,258	118,350	17,744
60	274,776	148,982	23,108
70	330,414	175,402	26,358
80	370,787	200,867	30,512
90	430,997	234,474	34,591
100	469,772	248,722	39,199

Table 2: Lisp Lexer Benchmarks

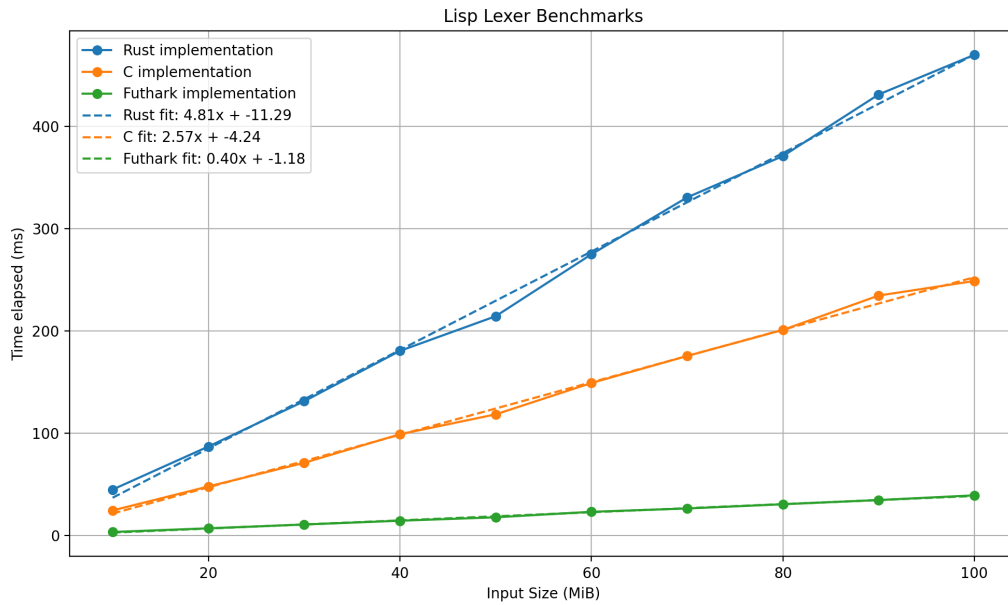


Figure 1: Lisp Lexer Benchmarks

As we see from the plot the times grows roughly linearly with the input size

⁶When taking these benchmarks I did not have access to other hardware since I could not get on the university VPN to access the compute cluster Hendrix.

for all the benchmarks. So we can use the linear equations found to solve for each lexers throughput in GiB/s.

	Rust (μs)	C (μs)	Futhark (μs)
Throughput (GiB/s)	0.2	0.4	2.5

Table 3: Lisp Lexers Throughput

6 Discussion

As we see the the throughput of

7 Conclusion

Conclusion.

References

- [1] William Henrich Due. *Parallel Parsing using Futhark*. 2023. URL: <https://futhark-lang.org/publications.html>.
- [2] Jonathan M.D Hill. “Parallel lexical analysis and parsing on the AMT distributed array processor”. In: *Parallel Computing* 18.6 (1992), pp. 699–714. ISSN: 0167-8191. DOI: [https://doi.org/10.1016/0167-8191\(92\)90008-U](https://doi.org/10.1016/0167-8191(92)90008-U). URL: <https://www.sciencedirect.com/science/article/pii/016781919290008U>.
- [3] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. USA: Addison-Wesley Longman Publishing Co., Inc., 2006. ISBN: 0321455363.
- [4] Aaron Wen-yao Hsu. “A Data Parallel Compiler Hosted on the GPU”. PhD thesis. Indiana University, 2019.
- [5] Ladislav Vagner and Bořivoj Melichar. “Parallel LL parsing”. In: *Acta Informatica* 44.1 (Apr. 2007), pp. 1–21. ISSN: 1432-0525. DOI: [10.1007/s00236-006-0031-y](https://doi.org/10.1007/s00236-006-0031-y). URL: <https://doi.org/10.1007/s00236-006-0031-y>.
- [6] Robin Voetter. “Parallel Lexing, Parsing and Semantic Analysis on the GPU”. MA thesis. LIACS, Leiden University, 2021.

-
- [7] Wikipedia contributors. *Deterministic finite automaton* — *Wikipedia, The Free Encyclopedia*. [Online; accessed 4-February-2024]. 2023. URL: https://en.wikipedia.org/w/index.php?title=Deterministic_finite_automaton&oldid=1192025610.