Fusing Lexing and Parsing

ANONYMOUS AUTHOR(S)

1 2 3

4

5

6

7

8

0

10 11

12

13

14

15

16

17

18

19

20

21

22 23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

Lexers and parsers are typically defined separately and connected by a token stream. This separate definition is important for modularity, but harmful for performance.

We show how to *fuse* together separately-defined lexers and parsers, drastically improving performance without compromising modularity. Our staged parser combinator library, flap, provides a standard parser combinator interface, but generates specialized token-free code that runs several times faster than ocamlyacc on a range of benchmarks.

INTRODUCTION 1

Software systems are easiest to understand when their components have clear interfaces that hide internal details. For example, a typical compiler includes separate lexer and parser components that communicate via a token stream.

Unfortunately, while interfaces improve clarity, they can harm performance, since hiding internal details reduces optimization opportunities. Parsers exemplify this tension: the token stream interface isolates parser definitions from character syntax details like whitespace, but it also carries overheads that reduce parsing speed.

Parsers built for efficiency avoid backtracking: only the initial token of the stream is typically needed at any time. However, even with this restriction, materializing and case-switching on tokens comes with a cost.

Contributions. This paper presents a transformation that significantly improves the performance of parsing by entirely eliminating tokens and fusing together lexers and parsers. Specifically, we present lexer-parser fusion using a parser combinator library, flap (*fused lexing and parsing*). The lexers and parsers are built using standard tools: Brzozowski's derivatives [Brzozowski 1964] for lexers (as reformulated by Owens et al. [2009]), and Krishnaswami and Yallop's typed algebraic parser combinators [Krishnaswami and Yallop 2019]. We review these standard tools in Section 2. We present the following contributions:

- We propose Deterministic Greibach Normal Form, a variant of Greibach Normal Form [Greibach 1965] for deterministic languages that captures syntactically the constraints enforced by the types in Krishnaswami and Yallop's typed context-free expressions (Section 3.1). We then formalize a translation from typed context-free expressions into Deterministic Greibach Normal Form, which serves as a basis for follow-up optimizations. We prove that the translation is well-defined and preserves the semantics of context-free expressions (Section 4).
- We present *lexer-parser fusion*, showing how to transform a separately-defined lexer and a normalized parser into a single piece of code that is specialized for calling contexts, entirely avoids materializing tokens, and case-switches only on individual characters, not on intermediate structures (Section 5).
- We implement our techniques as a parser combinator library flap and use multi-stage programming to generate efficient token-free code from the fused grammar (Section 6).
- We demonstrate the performance of flap, by showing that lexer-parser fusion results in efficient code that runs several times faster than code produced by standard tools such as ocamllex, ocamlyacc and menhir. We also assess other metrics, such as code size and compilation time (Section 7).

Finally, Section 8 surveys related work and Section 9 sets out some directions for further development.

Anon.





2 BACKGROUND: LEXER AND PARSER COMBINATORS

Figure 1 presents the novel code generation architecture of flap. As the figure shows, flap's interface is built from standard lexer [Owens et al. 2009] and parser combinators [Krishnaswami and Yallop 2019] drawn from existing work. This section gives an overview of those combinators.

2.1 Derivatives of regular expressions

Since lexical syntax is typically regular, lexers are typically defined using regular expressions (*regexes*). One particularly elegant formulation of regex matching, introduced almost six decades ago [Brzozowski 1964], is based on the idea of *derivatives*.

The *derivative* of a regex r with respect to a character c is another regex $\partial_c r$ that matches s exactly when r matches $c \cdot s$. For example, the regex (b|c)+ matches a sequence of one or more occurrences of b and c in any order. For a string that begins with c, either an empty suffix or some further sequence of b and c is acceptable, and so we have:

$$\partial_c (b|c) + = (b|c) *$$

The full rules, shown in Figure 2, are defined inductively on the syntax of regexes, with cases for the standard constructs \perp (which matches nothing), ϵ (which matches only the empty string), characters *b* and *c*, sequencing, alternation, and Kleene star, and for the less commonly-supported constructs intersection and negation. We refer the reader to Owens et al. [2009] for a fuller exposition.

There is a simple relationship between derivatives and automata for regular expressions: one way to construct an automaton is to take regular expressions *r* as states, and add a transition from r_i to r_j via character *c* whenever $\partial_c r_i = r_j$. For example, here is an automaton for (b|c)+:

start
$$\rightarrow (b|c)+ b (b|c)* c b c$$

In this example, the transitions all target the same state, since $\partial_b (b|c) + = \partial_c (b|c) + = \partial_b (b|c) = \partial_c (b$

Constructing an automaton is a common way to implement a regex matcher, and derivatives make it straightforward to built an automaton that is deterministic and compact. This process often involves representing the automaton as a graph or table. Alternatively, *multi-stage programming* [Taha 1999], makes it possible to directly generate code embodying the automaton.

Regex matching is an archetypal example of staged computation [Davies and Pfenning 1996] as found in languages like BER MetaOCaml [Kiselyov 2014]: although matching is a function of two inputs, regex and string, since the former is typically available first, it can be used to construct specialized code for processing the latter. In other words, while an unstaged matcher might have the following two-argument OCaml type

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:2

100	-
101	$\partial_c \perp =$
102	$\partial_c \epsilon =$
103	$\partial_c b =$
104	
105	Fig. 2. Derivative
106	otherwise).
107	
108	
109	val matchr :
110	a staged matche
111	of one argumen
112	val smatchr
113	For the regex
114	≪ let rec b
115	match s
116	
117	and bosta
118	match s
119	
120	in fun s
121	The quotation n
122	two fundament
123	r and string s fo

99

126

133

134

```
\begin{array}{rcl} \text{regular expression } r & \coloneqq & \perp \mid \epsilon \mid c \mid r \cdot s \mid (r \mid s) \mid r \ast \mid (r \& s) \mid \neg r \\ \partial_c \perp & = & \perp & \partial_c c = & \epsilon & & \partial_c r \ast = & \partial_c r \cdot r \ast \\ \partial_c \epsilon & = & \perp & \partial_c (r \cdot s) = & \partial_c r \cdot s \mid v(r) \cdot \partial_c s & \partial_c (r \& s) = & \partial_c r \& \partial_c s \\ \partial_c b & = & \perp & \partial_c (r \mid s) = & \partial_c r \mid \partial_c s & & \partial_c \neg r = & \neg(\partial_c r) \end{array}
```

Fig. 2. Derivatives of regular expressions (The *nullability* function v(r) expands to ϵ if r matches ϵ and \perp otherwise).

```
val matchr : regex \rightarrow string \rightarrow bool
```

a staged matcher is instead a function of one argument which generates code for another function of one argument:

val smatchr : regex \rightarrow (string \rightarrow bool) code

For the regex (b|c)+, smatchr might generate code found in the corresponding automaton:

The quotation notation $\ll e \gg$ in this example indicates a typed code value; it is one of MetaOCaml's two fundamental constructs (along with antiquotation \tilde{e}) for staging. The function for each regex r and string s follows a simple pattern, returning **true** if r is nullable and s is empty, and moving to the function for $\partial_c r$ if s begins with c.

2.2 Lexing with derivatives

Regexes built from derivatives are convenient for building lexers. A lexer is typically defined as an ordered mapping from regular expressions to *actions*, where an action might return a token, raise an error or invoke the lexer recursively to skip over some input. Figure 3a gives an example lexer with four actions: three of which return tokens atom, lpar and rpar, and one of which skips over whitespace.

Using MetaOCaml it is straightforward to express lexers as functions that accept lists of regexaction pairs and return code:

```
135type 'a action = Skip | Return of 'a code136val slex : (re * 'a action) list \rightarrow (string \rightarrow 'a) code
```

In practice, it is useful to extend these definitions with additional actions (such as Error) and additional information (such as the position of matched strings), which we omit here for succinctness.

Regex derivatives extend naturally to lexers by matching the input string against multiple regexes in parallel. Figure 3b is the automaton for matching one token with the sexp lexer, where each state corresponds to a vector of regexes, one for each lexer rule. The transition function ∂_c acts pointwise on the regex vector. **Return** rules correspond to labeled accepting states, and the **Skip** rule resets the vector to its initial state.

As Owens et al. [2009] show, lexers based on derivatives provide a practical basis for real-world lexing tools such as ml-ulex and the PLT Scheme scanner generator. One particularly useful feature for implementing lexers is the support derivatives provide for negation and disjunction, which make

1:3

Anon.





it straightforward to transform implicitly-ordered clauses for regexes *r* and *s* into order-independent disjoint clauses for *r* and $\neg(r)$ &s with the same semantics.

2.3 Parsing with typed context-free expressions

Parser combinators, introduced almost four decades ago by Wadler [1985], provide an elegant
 way to define parsers using functions. A parser combinator library provides functions denoting
 token-matching, sequencing, disjunction, and so on, allowing the library user to describe a parser
 by combining these functions in a way that reflects the structure of the corresponding grammar.
 Here are partial interfaces for constructing both regexes (type re) and parsers (type pa) in this way:

```
171(* Regex combinators *)(* Parser combinators *)171type retype 'a pa172val chr: char \rightarrow re(* token match *)val tok: 'a tok \rightarrow 'a pa173val (>>>): re \rightarrow re \rightarrow re (* sequence *) val (>>>): 'a pa \rightarrow 'b pa \rightarrow ('a * 'b) pa174val star: re \rightarrow re(* recursion *) val fix: ('a pa \rightarrow 'a pa) \rightarrow 'a pa175......
```

Both interfaces provide functions for token matching, sequencing and recursion. However, there
are some important differences: first, regexes act on characters, while parsers act on tokens (type
tok, a parameter of the library); second, parsers provide a general-purpose recursion operator fix,
while regexes offer only the more restrictive Kleene star; finally, the parser type is parameterized,
allowing parsers to construct and return suitably-typed syntax trees.

The earliest parser combinator libraries represented nondeterministic parsers, with support for 181 arbitrary backtracking and multiple results. Parsers defined in this way enjoyed various pleasant 182 properties (such as a rich equational theory), but suffered from potentially disastrous performance. 183 In a recent departure from the nondeterministic tradition, Krishnaswami and Yallop [2019] define 184 typed context-free expressions, whose types track properties such as FIRST sets and nullability in 185 order to preclude nondeterminism and ensure linear-time parsing using a single token of lookahead. 186 Krishnaswami and Yallop's design provides the standard set of parser combinators (as defined 187 above), but adds an additional type-checking step. They further apply multi-stage programming to 188 ensure that type-checking is completed before parsing begins, and to generate specialized parsing 189 code based on type information, leading to performance competitive with ocamlyacc. 190

Figure 4 summarizes Krishnaswami and Yallop's type system. A type is a triple recording nullability, the first set, and FLAST (analogous to the FOLLOW set). There is one typing rule for each combinator (e.g. sequencing $g_1 \cdot g_2$ and recursion $\mu\alpha : \tau \cdot g$); types are constructed using corresponding combinators (e.g. $\tau_1 \cdot \tau_2$). The two contexts Γ and Δ restrict the positions in which variables can occur to disallow left recursion, and the side conditions *separation* $\tau_1 \circledast \tau_2$ and *apartness* $\tau_1 # \tau_2$

159 160 161

162

197

216 217

218 219

220

221

222

223 224

225

226

227

228

229

230

231

232

233

234 235

236

198	Context-f	ree e	xpression $g := \epsilon \mid t \mid \perp \mid \alpha \mid g_1 \cdot g_2 \mid g_1 \lor g_2 \mid \mu\alpha : \tau. g$		
199 200	Types τ	€	{Null : 2; First : $\mathcal{P}(\Sigma)$; FLast : $\mathcal{P}(\Sigma)$ }	$\overline{\Gamma;\Delta\vdash\epsilon:\tau_\epsilon}$	$\overline{\Gamma;\Delta\vdash t:\tau_t}$
201 202 203	$ au_{\epsilon}$ $ au_{t}$	= = _	{NULL = true; FIRST = \emptyset ; FLAST = \emptyset } {NULL = false; FIRST = { t }; FLAST = \emptyset } {NULL = false; FIRST = \emptyset ; FLAST = \emptyset }	$\overline{\Gamma; \Delta \vdash \bot : \tau_{\bot}}$	$\frac{\alpha:\tau\in\Gamma}{\Gamma;\Delta\vdash\alpha:\tau}$
204 205 206	$\tau_1 \cdot \tau_2$	=	$\begin{cases} \text{NULL} = \tau_1.\text{NULL} \land \tau_2.\text{NULL} \\ \text{FIRST} = \tau_1.\text{FIRST} \cup \tau_1.\text{NULL}? \tau_2.\text{FIRST} \\ \text{FLAST} = \tau_2.\text{FLAST} \cup \tau_2.\text{NULL}? (\tau_2.\text{FIRST} \cup \tau_1.\text{FLAST}) \end{cases}$	$\frac{\Gamma; \Delta, \alpha:}{\Gamma; \Delta \vdash \mu \alpha}$	$\frac{\tau \vdash g : \tau}{z : \tau \cdot g : \tau}$
207 208 209	$\tau_1 \vee \tau_2$	=	$\begin{cases} \text{Null} = \tau_1.\text{Null} \lor \tau_2.\text{Null} \\ \text{First} = \tau_1.\text{First} \cup \tau_2.\text{First} \\ \text{FLast} = \tau_1.\text{FLast} \cup \tau_2.\text{FLast} \end{cases}$	$\begin{array}{c} \Gamma; \Delta \vdash \\ \Gamma, \Delta; \bullet \vdash \\ \tau_1 \in \end{array}$	$g_1: au_1 \ = g_2: au_2$ 9 $ au_2$
 210 211 212 213 214 215 	$\tau_1 \circledast \tau_2$ $\tau_1 \# \tau_2$ b?S	def = def = def =	τ_1 .FLAST $\cap \tau_2$.FIRST = $\emptyset \land \neg \tau_1$.NULL $(\tau_1$.FIRST $\cap \tau_2$.FIRST = $\emptyset) \land \neg(\tau_1$.NULL $\land \tau_2$.NULL) if <i>b</i> then <i>S</i> else \emptyset	$ \overline{\Gamma; \Delta \vdash g_1} \cdot \\ \Gamma; \Delta \vdash \\ \overline{\Gamma; \Delta \vdash g_2 : \tau_2} \\ \overline{\Gamma; \Delta \vdash g_1 \vee} $	$g_2 : \tau_1 \cdot \tau_2$ $g_1 : \tau_1$ $z \tau_1 \# \tau_2$ $g_2 : \tau_1 \lor \tau_2$

Fig. 4. Krishnaswami and Yallop's type system for context-free expressions

on the rules for sequencing and alternation reject ambiguous grammars, ensuring respectively that strings matched by sequenced parsers have a unique decomposition, and that the languages matched by alternated parsers do not overlap.

As an example, consider the following well-typed s-expression (we often omit τ in $\mu\alpha$: τ . g):

 $\mu \text{ sexp } .(\text{Lpar} \cdot (\mu \text{ sexps } . \epsilon \lor \text{ sexp } \cdot \text{ sexps }) \cdot \text{Rpar}) \lor \text{Atom}$

The following code shows how to use Krishnaswami and Yallop's parser combinators to define this grammar, using a token type with LPAR, RPAR and ATOM constructors, with the explicit fixed point represented using the Kleene star:

```
fix (fun sexp \rightarrow (tok LPAR >>> star sexp >>> tok RPAR)

$ fun p \rightarrow \ll Sexp (snd (fst ~p)) \gg

<|> tok ATOM $ fun s \rightarrow \ll Atom ~s \gg)
```

The fix, tok, >>> and star constructors are from the parser interfaces. The example additionally uses alternation <|> and map \$, which transforms the parsing result via a user-defined function, and MetaOCaml's quotation and antiquotation constructs $\ll exp \gg$ and $\sim exp$ to build code values.

3 OVERVIEW

The algorithm for parsing with typed context-free expressions introduced by Krishnaswami and Yallop is efficient at a high level, since it uses only a single token of lookahead and its execution time is linear in the length of its input. However, it is less efficient at a low level, since it examines each token multiple times: once at each alternation in the grammar, and then once again at the token combinator.

As an example, consider again μ sexp. (LPAR \cdot (μ sexps. $\epsilon \lor$ sexps.) \cdot RPAR) \lor ATOM. In the sub-grammar $\epsilon \lor$ sexp. \cdot sexps, neither alternative explicitly matches a token. Determining which branch to take therefore requires analysing the types to calculate whether the next tokens

1:5

in the input fall into the FIRST set of either ϵ or sexp \cdot sexps. Furthermore, once the token has 246 been examined and the branch taken, the parsing algorithm must examine it a second time. For 247 example, since LPAR \in FIRST(sexp \cdot sexps), the token LPAR causes the parsing algorithm to switch 248 to the sexp \cdot sexps branch. After the switch, parsing again uses the typing information for the 249 selected branch to determine which of sexp's sub-branches to take (i.e. LPAR or ATOM). Eventually, 250 the parsing algorithm encounters the LPAR node in the grammar, and the token is consumed. To 251 address these low-level inefficiencies, Krishnaswami and Yallop applied a variety of multi-stage 252 253 programming techniques, such as CPS conversion [Bondorf 1992; Nielsen and Sørensen 1995] to improve the results of staging, ultimately achieving performance that is competitive with ocamllex 254 and ocamlyacc. 255

In this work, we take a more systematic approach, making use of the guarantees offered by
 the types to transform grammars into a normal form that is amenable to a sequence of further
 optimizations. More precisely,

- (1) We first propose a novel normal form, *Deterministic Greibach Normal Form (DGNF)*, which gathers together the places in the grammar that involve branching on tokens, allowing tokens to be discarded immediately after inspection (Section 3.1).
 - (2) We then formalize a normalization algorithm that traverses a context-free expression and returns a DGNF grammar. Normalization works well for well-typed context free expressions, and the resulting DGNF grammar sets the basis for follow-up optimizations (Section 3.2).
 - (3) Based on the normal form, we present a fusion process that ultimately eliminates the need to materialize tokens altogether. The fusion algorithm starts with a separately-defined lexer and parser, connected together via tokens, and produces entirely token-free code, in which the only branches involve inspecting individual characters (Section 3.3).
 - (4) Finally, flap uses MetaOCaml's staging facilities to generate code for the fused grammar. Since the normalized grammar representation is already amenable to generating optimized code, flap does not need the sophisticated techniques used by Krishnaswami and Yallop [2019] (Section 3.4).

As we shall see, these optimizations make parsers built from typed context-free expressions significantly more efficient than both ocamlyacc and Krishnaswami and Yallop's system (Section 7).

The running example. This section illustrates flap's key ideas through a running example given in Figure 5. Figure 5a presents the grammars that will be used throughout this paper. Figure 5b and 5c repeat our previous s-expression lexer and example grammar respectively for better readability. From now on we use colors to distinguish different grammars. Regular expressions r include \perp for nothing, ϵ for the empty string, characters c, sequencing $r \cdot s$, alternation $r \mid s$, Kleene star r*, intersection r & s, and negation $\neg r$. Lexers L are a set of regex-action pairs where each action either returns a token or skips. Context-free expressions g are \perp for the empty language, ϵ for the language containing only the empty string, t which matches the language containing only the single-element string t, variables α , sequences $g_1 \cdot g_2$, unions $g_1 \lor g_2$, and the least fixed point operator $\mu \alpha : \tau .g$. We will introduce the normal form grammar G and fused grammar F later.

3.1 Deterministic Greibach Normal Form

Like Krishnaswami and Yallop's system, flap first ensures that input grammars are well-typed according to the rules in Figure 4. It then applies a normalization algorithm that transforms welltyped grammars into a novel normal form that avoids the need for repeated branching.

Specifically, to ensure that grammars can be used for deterministic parsing with a single token of
 lookahead, we introduce *Deterministic Greibach Normal Form* (DGNF), a variant of *Greibach Normal Form* (GNF) [Greibach 1965]. More precisely, in GNF, all the productions of a grammar take the

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285 286

Fusing Lexing and Parsing

343

295 $:= \perp |\epsilon| c |r \cdot s| (r |s) |r * | (r \& s) |\neg r$ regular expression r 296 lexer L ::= $\{ r \Rightarrow \text{Return } t \} \cup \{ r \Rightarrow \text{Skip} \}$ 297 $\perp |\epsilon| t |\alpha| g_1 \cdot g_2 |g_1 \vee g_2| \mu\alpha : \tau. g$ context-free expression ::= q 298 299 normal form Ν $::= \epsilon | t \overline{n} | \alpha \overline{n}$ 300 G $:= \{n \rightarrow N\}$ normal form grammar 301 302 $:= \{n \to r \overline{n}\} \cup \{n \to ?r\}$ F fused grammar 303 (a) Syntax of lexers, forms, and grammars in flap 304 305 id **Return** ATOM 306 id $\stackrel{\text{def}}{=}$ [a-z]+ space \Rightarrow Skip (b) . 307 space $\stackrel{\text{def}}{=}$ $| \ |$ fusing (e) \Rightarrow Return LPAR 308) Return RPAR → (d) 309 normalizing 310 (b) S-expression lexer (2.2) 311 312 $\mu \text{ sexp } .(\text{Lpar} \cdot (\mu \text{ sexps } . \epsilon \lor \text{ sexp } \cdot \text{ sexps }) \cdot \text{Rpar}) \lor \text{Atom}$ 313 314 (c) A well-typed s-expression grammar (2.3) 315 316 sexp LPAR sexps rpar RPAR sexps LPAR sexps rpar sexps rpar ::= ::= ::= 317 АТОМ ATOM sexps 318 ϵ 319 320 (d) The above s-expression grammar in Deterministic Greibach Normal Form (3.2) 321 322 id Return ATOM ~id~~~Return/~romid \Rightarrow **Return** ATOM 323 space Skip Skip space Skip space 324 VReturn ArAr Return LPAR Return LPAR 325 Return RPAR-**Return** RPAR Return BAAR -326 327 sexp (sexps rpar rpar) sexps ::= (sexps rpar sexps ::= ::= 328 id sexps id space rpar 329 space sexps space sexp 330 $?\neg$ (id | space | () 331 (e) Fusing drops lexing rules that return non-matchable tokens (top); the fused s-expr grammar (bottom) (3.3) 332 333 Fig. 5. flap: running example of an s-expression. Grammars in the example are written in BNF form. 334 335 336 form $n \to tn_1n_2 \dots n_k \ (k \ge 0)$ where n and n_i are nonterminals and t is a terminal. DGNF further 337 imposes the following syntactic constraints: 338 DEFINITION 1 (DETERMINISTIC GREIBACH NORMAL FORM (SYNTAX)). A grammar G is in Deter-339 ministic Greibach normal form if all productions are of form $n \to tn_1n_2 \dots n_k (k \ge 0)$, and moreover, 340 • (Determinism) for any pair of a nonterminal n and a terminal t, there is at most one production 341 beginning $n \rightarrow tn_1n_2 \dots n_k$; 342

(Guarded ε-productions) Additional, each nonterminal may have an optional ε-production n→ ε, with the proviso that the ε-production may only be used when no terminal symbol in the non-ε productions matches the input string.

Figure 5a includes the definition of normal form and normal form grammar, where normal form *N* is either an epsilon ϵ or a terminal followed by a list of nonterminals $t \ \overline{n}$, with \overline{n} denoting $n_1 n_2 \dots n_k (k \ge 0)$. Note that *N* also includes a special internal form $\alpha \ \overline{n}$, highlighted with gray, which is explained in Section 3.2. The normal form grammar *G* is a set of productions that map nonterminals to normal forms.

Intuitively speaking, the constraints on the DGNF grammar are a syntactic analogue of the constraints enforced by the types in Krishnaswami and Yallop's typed context-free expressions. The constraints have a simple practical motivation in parsers. That is, each $n \rightarrow tn_1n_2 \cdots n_k$ production represents one branch that matches a distinct terminal *t*, and guarded ϵ -productions represent an else branch that is taken if none of the active productive branches matches the input.

Examples. We consider a few examples. For readability, we write the grammar in BNF form, e.g., $n ::= An_1n_2 \mid B$ corresponds to $n \rightarrow An_1n_2$ and $n \rightarrow B$.

((1) n	::=	An_1n_2	(2)	n	::=	ABn_1	(3)	n	::=	A <i>n</i> ₁	(4)	n	::=	An_1n_2
			В		n_1	::=	С				An_2		n_1	::=	С
	n_1	::=	С						n_1	::=	С				ϵ
	n_2	::=	D						n_2	::=	D		n_2	::=	С

Here (1) is in DGNF, while (2) (3) are not. In (2), n starts with two terminals, while in (3), n has two productions starting with A, violating the determinism condition.

(4) is the most subtle case. Consider matching *n* with AC. First, *n* expands to An_1n_2 . But should *n*₁ then expand to c or ϵ ? In a general nondeterministic grammar, it is impossible to tell simply by looking ahead at the next token c: we may first consider $n_1 \rightarrow c$ and, finding that n_2 fails to match, backtrack to the other branch to consider $n_1 \rightarrow \epsilon$ and $n_2 \rightarrow c$ and succeed. However, the proviso in Definition 1 eliminates this choice: only $n_1 \rightarrow c$ applies, and so the grammar does *not* match Ac. In fact, the semantics of DGNF (Section 4.2) will rule out (4) as a DGNF grammar, ensuring that parsing is deterministic.

As the examples demonstrate, the determinism and guarded ϵ -production conditions ensure that there is never any ambiguity about whether a production rule applies during parsing.

3.2 Normalization

We formalize a normalization algorithm, which, given any well-typed context-free expression, turns
it into a DGNF grammar. As an example, Figure 5d presents the result in BNF form of normalizing
the s-expression grammar in Figure 5c. It is straightforward to check that the normalized grammar
represents exactly the same language as the original context-free expression.

Importantly, the normalized DGNF presentation addresses the problem of repeated branching discussed in the beginning of Section 3. With this normalized form, parsing a sexps involves reading the next token, and branching to the first, second or third branch depending on whether the token is LPAR, ATOM or something else. In the first two cases the token is consumed immediately, and parsing moves on to the next token in the input. The last case is somewhat more costly, since the token may be needed again immediately afterwards, and more care is needed to avoid wasted work.

Section 4 discusses the normalization algorithm, which traverses a context-free expression, and builds grammar productions according to the expression structure. As we will see, defining the algorithm poses significant challenges, particularly around fixed points. When normalizing $\mu\alpha$. *g*, although we do not yet know the normalized grammar for α , we must proceed with normalizing *g*

392

357 358

359

365

366

374

375 376

377

344

345

regardless. In this case, it is necessary to "tie the knot" when the result of normalizing *g* becomes available, which requires us to introduce an intermediate non-DGNF form $n \rightarrow \alpha \overline{n}$, causing complication and subtleties during normalization. This non-DGNF form is purely internal, and does not appear in the normalization results of closed expressions.

Finally, great care needs to be taken to guarantee that normalization produces indeed DGNF grammars. That requires us to ensure that the normalization captures the constraints enforced by the types in Krishnaswami and Yallop's system. We prove that normalization is correct, and that normalized grammars preserve the denotational semantics of context-free expressions (Section 4).

3.3 Fusion

Next, flap applies lexer-parser fusion, one of our central contributions. Fusion acts on a lexer and a normalized parser, connected together via tokens, and products a grammar representation that is entirely token-free, in which the only branches involve inspecting individual characters.

Figure 5a defines the syntax of fused grammars, where the fused form f is either a regex followed by a list of nonterminals $r \overline{n}$, or a single-token lookahead ?r for tokens matched by r. The fused grammar F is a set of productions $\{n \to f\}$.

We illustrate the key idea of fusion through Figure 5e, which fuses the s-expression lexer in Figure 5b and the normalized parser in Figure 5d. As the first step, the fusion algorithm implicitly specializes the lexer to each nonterminal *n* in the normalized grammar, and lexing rules that return tokens not in productions for the nonterminal *n* are discarded. We take rpar as an example: the rpar nonterminal has only a single production, which begins with the terminal RPAR. We then look at the lexing rules, and discard those rules that do not return RPAR. However, the skip rule is retained, since skipped characters can precede any token. Then, the algorithm fuses the lexing rules and the parsing rules, by substituting the tokens in the parsing rules by regexes in the lexing rules that return corresponding tokens. The bottom of Figure 5e presents the fused grammar for rpar, which has two branches. The first branch fuses lexing and parsing, by having the original token RPAR replaced with the regex). The second branch is an extra production corresponding to the skip rule in the lexer, allowing rpar to match an arbitrary number of space. Notably, after fusion rpar now directly matches space or), without referring to the token RPAR.

The bottom of Figure 5e presents the complete result of fusing the s-expression lexer and normalized grammar. Like the case for rpar, the tokens ATOM, LPAR and RPAR in the grammar have been replaced with the regular expressions id, (and) associated with those tokens in the lexer. Moreover, for each nonterminal *n* there is an extra production n ::= space n, corresponding to the skip rule in the lexer. Finally, the ϵ -production sexps $\rightarrow \epsilon$ has been replaced with a lookahead rule sexp \rightarrow ? \neg (id | space | (), consisting of the complement of the three regular expressions that appear at the start of the right hand side of the other productions.

Normalization allows the fusion algorithm to be defined concisely (Section 5). In particular, the constraints on the positions of terminals make it straightforward to fuse the lexing rules into the grammar without disrupting its structure. More importantly, the fused grammar inherits the properties of DGNF: the productions of a nonterminal start with distinct regexes, and an optional lookahead rule may only be used when no regexes in other productions match the input string.

In summary, fusion transforms a separately-defined lexer and a normalized parser into branches on individual characters, entirely eliminating intermediate tokens. As Section 7 shows, fusion significantly improves the performance of parsing. Furthermore, Section 7.3 suggests that the size increase resulting from normalization and fusion in flap is relatively modest, and Section 7.4 reports that compilation times are sufficiently low for interactive use.

1:10

Anon.

442			normal form $N := \epsilon t \overline{n} q \overline{n}$	
443			normal form grammar $G := \{n \to N\}$	
444		-		
445		N[[a]]	I returns $n \rightarrow C$ with a grammar C and the start nonterminal n	
446	E	$\mathcal{N} \llbracket \mathcal{G} \rrbracket$	η returns $n \to 0$, with a graninal 0 and the start nonterminal n	~
447	Ea	ich rule alloca	cates a fresh nonterminal n , except for rule $(f(x))$, which allocates a fresh c	٤
448	(epsilon)	$\mathcal{N} \llbracket \epsilon \rrbracket$	$= n \Longrightarrow \{ n \to \epsilon \}$	
449	(token)	$\mathcal{N} \llbracket t \rrbracket$	$= n \Longrightarrow \{ n \to t \}$	
450	(bot)	$\mathcal{N}\llbracket ot bracket$	$= n \Rightarrow \emptyset$	
451	(seg)	$\mathcal{N}[\![q_1 \cdot q_2]\!]$	$= n \Longrightarrow \{ n \to N_1 n_2 \mid n_1 \to N_1 \in G_1 \} \cup G_1 \cup G_2$	
452		191 951	where $\mathcal{N}[\![q_1]\!] = n_1 \Rightarrow G_1 \land \mathcal{N}[\![q_2]\!] = n_2 \Rightarrow G_2$	
453				0.110
454	(alt)	$N \llbracket g_1 \lor g_2 \rrbracket$	$ = n \Rightarrow \{n \rightarrow N_1 \mid n_1 \rightarrow N_1 \in G_1\} \cup \{n \rightarrow N_2 \mid n_2 \rightarrow N_2 \in G_2\} \cup $	$G_1 \cup G_2$
455			where $\mathcal{N} \llbracket g_1 \rrbracket = n_1 \Rightarrow G_1 \land \mathcal{N} \llbracket g_2 \rrbracket = n_2 \Rightarrow G_2$	
456	(fix)	N[[μα. g]]	$= \alpha \Longrightarrow \{ \alpha \to N \mid n \to N \in G \}$	(1)
457			$\cup \{ n' \to N \overline{n}' \mid n' \to \alpha \overline{n}' \in G \land n \to N \in G \}$	(2)
458			$\cup G \setminus_{n' \to \alpha} \overline{n'}$	(3)
459			where $\mathcal{N}[\![g]\!] = n \Rightarrow G$	
460			$G \setminus_{n' \to \alpha} \overline{n'}$ is <i>G</i> with all $n' \to \alpha \overline{n'}$ removed for any <i>n'</i> and	\overline{n}'
461	(var)	$N[\alpha]$	$= n \Rightarrow \{n \to \alpha\}$	
462	()	· · II II		
463				

Fig. 6. Normalization of well-typed context-free expressions.

3.4 Staging

464 465 466

467

468

469

470

471

472

473

474

475

476

477

478

479

480 481

482

483

484

485 486

487

488

489 490 In the last step, flap uses MetaOCaml's staging facilities to generate code for the fused grammar. The normalized grammar representation used in flap makes this process comparatively straightforward; it does not involve sophisticated optimization techniques such as the binding-time improvements applied by Krishnaswami and Yallop [2019]. Furthermore, flap does not rely on compiler optimizations to further simplify the code it generates; instead, it directly generates efficient code, containing no indirect calls, no higher-order functions and no allocation, except where these elements are inserted by the user of flap in semantic actions.

Specifically, the staging step in flap generates one function for each parser state (i.e. for each pair of a nonterminal and a regex vector), following a parsing algorithm with fused grammars, but eliminating information that is statically known, such as the nullability and derivatives of the regexes associated with each state.

Section 6 presents the algorithm underlying flap's staged parsing implementation, and Section 6.5 describes the implementation itself in more detail.

4 NORMALIZING CONTEXT-FREE EXPRESSIONS

This section presents a normalization algorithm that transforms context-free expressions into grammars in Deterministic Greibach Normal Form (DGNF). The normalization sets the basis for follow-up optimizations of fusion and staging.

4.1 Normalization to DGNF

Figure 6 defines a normalization algorithm for typed context-free expressions. We repeat here the syntax of normal forms and normal form grammars. As discussed in Section 3.2, normal forms

496

497

498

499

500

501

502 503

504

505

506

507

508

509

510 511 512

526

527

528

529

530

531 532 533

534

535 536

537

538 539

include a non-DGNF internal form $\alpha \overline{n}$ used when normalizing fixpoints, where α is interpreted as a special kind of nonterminals. Since α is a nonterminal, α itself may appear as part of a $t \overline{n}$ (e.g., *t* α). Because of this internal form, normal forms *N* are actually *not* in DGNF. But as we will show later, $\alpha \overline{n}$ is an intermediate form which will be entirely eliminated in the final result, eventually turning the grammar into DGNF.

The key to normalization is the function $\mathcal{N}[\![g]\!]$ that normalizes a context-free expression g, yielding a normalized grammar G and a distinguished start nonterminal n of the grammar.

There are seven cases for the seven context-free expression constructors. Each case involves allocating a fresh nonterminal (*n* or α) to use as the start symbol. The cases with sub-expressions $(g_1 \cdot g_2, g_1 \vee g_2 \text{ and } \mu \alpha, g)$ are defined compositionally in terms of the normalization of those sub-expressions. Since normalization simply merges together all the production sets resulting from sub-expressions, the situation frequently arises where some productions are not reachable from the start symbol; the definition here ignores this issue, since it is easy to trim unreachable productions in the implementation.

Rules (*epsilon*), (*token*), and (*bot*) are straightforward. For each of ϵ and c, normalization produces a grammar with a single production whose right-hand side is ϵ or c respectively. For \perp , normalization produces an empty grammar, with a start symbol and no productions.

Rule (*seq*) is defined compositionally in terms of the normalization of their sub-expressions. For sequencing $g_1 \cdot g_2$, normalizing g_1 and g_2 produces the start symbol n_1 and n_2 , respectively. Now what we want is $n \rightarrow n_1 n_2$. That is:

$$(seq1) \quad \mathcal{N}\llbracket g_1 \cdot g_2 \rrbracket = n \Rightarrow \{n \to n_1 \ n_2\} \cup G_1 \cup G_2 \qquad \text{where} \quad \mathcal{N}\llbracket g \rrbracket = n \Rightarrow G$$

However, while this is semantically correct, $n \rightarrow n_1 n_2$ is not in normal form. Therefore, normaliza-513 tion in rule (seq) copies each production N_1 of n_1 and appends to each the start symbol n_2 , producing 514 $N_1 n_2$. To see why that is correct, consider if N_1 is of form $t \overline{n}$ or $\alpha \overline{n}$. Then $N_1 n_2$, i.e., $t \overline{n} n_2$ or 515 $\alpha \overline{n} n_2$, is indeed of normal form. However, we need to prevent N_1 from being ϵ , or otherwise ϵn_2 516 would be ill-formed. The case for sequencing is one of several places where the correctness of 517 normalization depends on the types. In particular, if $q_1 \cdot q_2$ is well-typed, then N_1 cannot be ϵ . 518 Specifically, the typing rule for sequencing (Figure 4) depends on the separation relation $\tau_1 \circledast \tau_2$, 519 which guarantees that the NULL of q_1 is false. With that guarantee, we show that N_1 cannot be ϵ 520 (Lemma 4.2), ensuring that the result of the normalization function is in normal form. 521

Rule (*alt*) for alternation $g_1 \vee g_2$ is similar. In this case, normalization merges the productions for the start symbols n_1 of g_1 and n_2 of g_2 into the productions for the new start symbol n. During typing (Figure 4), the apartness relation $\tau_1 \# \tau_2$ ensures that the FIRST sets of g_1 and g_2 do not intersect, and that at most one of g_1 and g_2 is nullable. This guarantees that the result is well-formed.

The final two rules (*fix*) and (*var*) deal with the binding fixed point operator $\mu\alpha$. *g* and with bound variables α . In rule (*fix*), we assume we can always alpha-rename a fixed point so α is unique. Normalization for the fixed point operator $\mu\alpha$. *g* takes place in two stages. First, the body *g* is normalized. The normalization result suggests that the grammar of the body *g* has a start symbol *n*. Then, according to the semantics of fixed point, we should proceed to tie the knot by producing $\alpha \rightarrow n$ and return α as the start symbol. That is:

(fix1)
$$\mathcal{N}[\![\mu\alpha, g]\!] = \alpha \Rightarrow \{\alpha \to n\} \cup G$$
 where $\mathcal{N}[\![g]\!] = n \Rightarrow G$

However, the rule $\alpha \to n$ is not in normal form. Therefore, instead of directly returning $\alpha \to n$, we proceed to copy the productions for *n* into the rules for α :

(fix2)
$$\mathcal{N}[\![\mu\alpha, g]\!] = \alpha \Rightarrow \{\alpha \to N \mid n \to N \in G\} \cup G$$
 where $\mathcal{N}[\![g]\!] = n \Rightarrow G$

But there is some extra work before we return the result. In particular, productions in *G* might start with α (e.g., $n' \rightarrow \alpha \overline{n}'$). While such form is allowed by the syntax of *N*, our ultimate goal is to turn

1:12

5	4	0

540	
541	$\mathcal{N}\llbracket \mu \operatorname{sexps} \cdot \epsilon \vee \operatorname{sexp} \cdot \operatorname{sexps} \rrbracket = \operatorname{sexps} \Rightarrow \{\operatorname{sexps} \to \epsilon, \operatorname{sexps} \to \operatorname{sexps} \}$
542	$\mathcal{N}[\![\operatorname{LPAR} \cdot (\mu \operatorname{sexps} \cdot \epsilon \lor \operatorname{sexps})]\!] = n_1 \Rightarrow \{ n_1 \to \operatorname{LPAR} \operatorname{sexps} , \operatorname{sexps} \to \epsilon, \operatorname{sexps} \to \operatorname{sexp} \operatorname{sexps} \} \qquad \mathcal{N}[\![\operatorname{RPAR}]\!] = \cdots$
543	$\overline{\mathcal{N}[[\mathtt{LPAR} \cdot (\mu \text{ sexps } \cdot \epsilon \lor \text{ sexps }) \cdot \mathtt{RPAR}]]} = n_2 \Rightarrow \{ n_2 \to \mathtt{LPAR \text{ sexps } rpar }, \mathtt{sexps } \to \epsilon, \mathtt{sexps } \to \mathtt{sexp \text{ sexps }}, \mathtt{rpar } \to \mathtt{RPAR} \}$
544	$\mathcal{N}[\![(\text{LPAR} \cdot (\mu \text{ sexps } \cdot \epsilon \lor \text{ sexps }) \cdot \text{RPAR}) \lor \text{ATOM}]\!] = n_3 \Rightarrow \{ n_3 \rightarrow \text{LPAR sexps rpar}, n_3 \rightarrow \text{ATOM}, \text{sexps } \rightarrow \epsilon, \text{sexps } \rightarrow \text{sexp sexps, rpar} \rightarrow \text{RPAR} \}$
545 546	$\mathcal{N}[\![g]\!] = \operatorname{sexp} \Rightarrow \{ \operatorname{sexp} \to \operatorname{Lpar} \operatorname{sexps} \operatorname{rpar}, \operatorname{sexp} \to \operatorname{Atom}, \operatorname{sexps} \to \epsilon, \operatorname{sexps} \to \operatorname{Lpar} \operatorname{sexps} \operatorname{rpar} \operatorname{sexps}, \operatorname{sexps} \to \operatorname{Atom} \operatorname{sexps}, \operatorname{rpar} \to \operatorname{Rpar} \}$

547 548

549

566

567 568

577

582

583

584

585

586

587 588

Fig. 7. Normalizing s-expression $q = \mu \operatorname{sexp} (\operatorname{LPAR} \cdot (\mu \operatorname{sexps} \cdot \epsilon \lor \operatorname{sexps}) \cdot \operatorname{RPAR}) \lor \operatorname{ATOM}$

the productions into DGNF, where every nonterminal either starts with a terminal or is ϵ . Now 550 that we learn the rules of α , we can look up and substitute in G all productions that start with α . 551 For example, if $\alpha \to B$ and $n' \to \alpha \overline{n}$, then after substitution we have $n' \to B \overline{n}$. Note that we still 552 allow α , as a special kind of nonterminal, to appear inside G if it is not the start of a production. 553 That is, α in $n' \rightarrow t \alpha$ won't get substituted. It would actually be wrong to perform the substitution: 554 if $\alpha \to B$, then after substitution $n' \to tB$ is not in DGNF. 555

Rule (fix) in Figure 6 presents the final form of normalizing a fixed point. The normalization first 556 copies the productions for n into the rules for α (1), then substitutes in G all productions that start 557 with α (2), and finally adds back all productions in G that do not start with α (3). As we will see, 558 rule (*fix*) effectively guarantees that normalizing closed context-free expressions produces DGNF. 559

Lastly, in rule (*var*), we create a fresh start symbol *n* with a singleton production $n \to \alpha$. Combin-560 ing rule (fix) with rule (var), normalization essentially treats α as a placeholder for the productions 561 denoted by the fixed point. As soon as we know what α should be, we substitute α with its pro-562 ductions if necessary (as in rule (fix)). It may be tempting here to return $\alpha \rightarrow \emptyset$ with α as a start 563 symbol and no productions. That would be wrong, as $\alpha \Rightarrow \emptyset$ means an empty grammar, causing 564 problems when, for example, rule (alt) copies productions. 565

Example. Figure 7 presents the simplified derivation of normalizing

$$g = \mu \operatorname{sexp} .(\operatorname{Lpar} \cdot (\mu \operatorname{sexps} . \epsilon \lor \operatorname{sexps}) \cdot \operatorname{Rpar}) \lor \operatorname{Atom}$$

569 where we highlight grammar changes in light gray, and omit some details via · · · for space reasons 570 and also since normalizing tokens is straightforward. We include the complete derivation tree in the 571 appendix. Of particular interest is the last step, which normalizes a fixed point. In this case, sexp is 572 used as the variable bound by the fixed point, and we have a production sexps \rightarrow sexp sexps which is 573 not DGNF. First, sexp copies all productions from n_4 . Then, since the production sexps \rightarrow sexp sexps 574 starts with sexp, the production expands to two productions where sexp is replaced by its two 575 normal forms respectively, making the resulting grammar in DGNF. 576

4.2 Semantics of DGNF

578 Our Definition 1 of DGNF presented in Section 3.1 gives a high-level description of a DGNF grammar. 579 To prove that our normalization actually results in DGNF grammars, we first define the meaning of 580 DGNF in terms of our formalization. 581

We start with the expansion relation:

DEFINITION 2 (EXPANSION ($G \vdash \rightarrow$)). Given a grammar G, we define the expansion relation by (1) (Base) $G \vdash n \rightarrow n$; (2) (Step) if $G \vdash n \rightarrow \overline{t} n' \overline{n}$, and $(n' \rightarrow N' \in G)$, we have $G \vdash n \rightarrow \overline{t} N' \overline{n}$. We write $G \vdash n \rightsquigarrow w$ when n expands to a complete word w.

The expansion relation essentially captures what a nonterminal can expand to. For example, if $n \to B n_1 \in G$ and $n_1 \to C \in G$, then we have $G \vdash n \to BC$. We enforce a left-to-right expansion

Fusing Lexing and Parsing

593

594

595

596

597

598

599

600 601

602

603

604

605

606

607

608

609

610

611 612

613

614

615

616

617

618 619

620 621

622

623 624

625

626

627

628

629

630

631 632

633

order for clarity and to stay close to the parsing behavior, but that is not necessary: it is easy to imagine an arbitrary order expansion, but any order leads to the same set of words.

With the notion of expansion, we define what it means for a grammar to be in DGNF precisely.The definition below follows the syntax given in Definition 1:

DEFINITION 3 (DETERMINISTIC GREIBACH NORMAL FORM (SEMANTICS)). A grammar G is in Deterministic Greibach normal form if all productions are of form $n \rightarrow tn_1n_2 \dots n_k (k \ge 0)$, and moreover,

- (Determinism) for any pair of a nonterminal n and a terminal t, if there are two distinct productions (n → t₁ n
 ₁) ∈ G and (n → t₂ n
 ₂) ∈ G, we have t₁ ≠ t₂;
- (Guarded ϵ -productions) If $G \vdash n \rightsquigarrow \overline{t} n_1 n_2 \overline{n}$, if $(n_1 \to \epsilon) \in G$, then for any t either $(n_1 \to t \overline{n}_1) \notin G$ or $(n_2 \to t \overline{n}_2) \notin G$.

The Determinism condition is straightforward, while the Guarded ϵ -productions condition needs more explanation. In Definition 1, we mentioned that the ϵ -production may only be used when no terminal symbol in other productions matches the input string. Consider that the next character to match is c. The case when both the ϵ -production $n_1 \rightarrow \epsilon$ and a production $n_1 \rightarrow c$ can match raises when n_1 's follow-up nonterminal n_2 can also match c, making it possible to use the ϵ -production while $n_1 \rightarrow c$ also matches. Definition 3 captures such cases, requiring that n_1 and n_2 cannot match the same terminal if n_1 has an ϵ -production, and thus rules out example (4) in Section 3.1.

Now we can formally define the important property of DGNF that makes it practically useful.

THEOREM 4.1 (DETERMINISTIC PARSING). If G is a DGNF grammar, then for any expansion $G \vdash n \rightsquigarrow w$, there is a unique derivation for such expansion.

4.3 Well-definedness and correctness

The normalization process serves as the basis for the parsing algorithm, and thus establishing its correctness is crucial for flap. In this section, we prove three key properties of normalization: first, normalization always succeeds for well-typed expressions (Section 4.3.1); then, normalization result will eventually get rid of the internal form $\alpha \overline{n}$ (Section 4.3.2); and finally, the result of normalization is a DGNF grammar (Section 4.3.3).

4.3.1 Normalization is well-defined. As we have briefly discussed in Section 4.1, correctness of normalization depends on types. For example, when normalizing sequencing $g_1 \cdot g_2$, rule (*seq*) returns $N_1 n_2$ with $n_1 \rightarrow N_1$ from g_1 , and n_2 from g_2 . In order for $N_1 n_2$ to be well-formed, we must ensure that N_1 is not ϵ , or otherwise ϵn_2 is ill-formed.

To this end, we make use of the typing information during normalization. In the case of sequencing, since $g_1 \cdot g_2$ is well-typed, the separation relation ($\tau_1 \otimes \tau_2$ in Figure 4) ensures g_1 is not nullable. We then prove that if an expression is not nullable, its normalization cannot have an ϵ production.

LEMMA 4.2 (PRODUCTIONS OF NULL). Given $\Gamma; \Delta \vdash g : \tau$ and $\mathcal{N}[\![g]\!]$ returns $n \Rightarrow G, \tau.NULL =$ true if and only if (1) $n \rightarrow \epsilon \in G$; or (2) $n \rightarrow \alpha \in G$ where $(\alpha : \tau') \in \Gamma$ and $\tau'.NULL =$ true.

In other words, if τ .NULL = false, then $n \to \epsilon \notin G$.

With Lemma 4.2 and similar reasoning about typing during normalization, we prove that normalization is well-defined for well-typed expressions.

THEOREM 4.3 (WELL-DEFINEDNESS). If Γ ; $\Delta \vdash g : \tau$, then $\mathcal{N}[\![g]\!]$ returns $n \Rightarrow G$ for some G and n.

4.3.2 Normalizing closed expressions produces no $\alpha \overline{n}$ form. Theorem 4.3 says that if an expression is well-typed, then normalization returns a grammar successfully. However, this grammar may not be in DGNF. For example, the grammar may include $n \rightarrow \alpha \overline{n}$, which is not a valid DGNF form. Therefore, in order for the normalization result to be in DGNF, we need to prove that all $n \to \alpha \overline{n}$ productions are removed from the normalization result.

In this part, we prove that the normalization result cannot contain any $\alpha \overline{n}$ production for a *closed* well-typed expression. To do so, we need to reason about the occurrences of α . The following lemma says that all α returned as the head of a production must be in the typing context.

LEMMA 4.4 (INTERNAL NORMAL FORM). Given $\Gamma; \Delta \vdash g : \tau$ and $\mathcal{N}[\![g]\!]$ returns $n \Rightarrow G$,

- $if(n \to \alpha \overline{n}) \in G$, then we have $\alpha \in dom(\Gamma)$;
- 645 646

643

644

655

656 657

658

659

660

661

662

663

664 665

666 667

668

669

670

671

672

673

674

675

676

677

678 679

680

• if $(n' \to \alpha \overline{n}) \in G$ for any n', then we have $\alpha \in fv(g)$, and thus $\alpha \in dom(\Gamma, \Delta)$.

Note that the first result applies only to the start symbol *n*, and its proof replies on the typing rule 647 where α is well-typed only if $\alpha \in \Gamma$ (Figure 4). The second result applies to any n', and is proved 648 making use of the first result. Specifically, the proof goes by induction on $\Gamma; \Delta \vdash q : \tau$, and the most 649 650 interesting case is when normalizing $\mu\alpha$. q, where we need to prove that the productions of the start symbol of q cannot start with α , or otherwise normalizing $\mu\alpha$. q would copy all productions 651 from q for α which would result to, e.g., $\alpha \to \alpha$ that fails the lemma. Fortunately, that is exactly 652 what the first result tells us: when typing $\mu \alpha$. *q*, we add α in Δ (Figure 4), and thus normalizing *q* 653 cannot have α at the head of a production for its start symbol. 654

Our goal then follows as a corollary of Lemma 4.4, which says that normalizing any closed well-typed expression only produces the desired normal form.

COROLLARY 4.5 (NORMAL FORM). Given $\bullet; \bullet \vdash g : \tau$ and N[[g]] returns $_ \Rightarrow G$, then for all $(n \rightarrow N) \in G, N$ is ϵ or $t \overline{n}$ for some t and \overline{n} .

4.3.3 *Normalization returns DGNF grammars.* Finally, we prove that normalization returns DGNF grammars. That requires productions to satisfy the conditions given in Definition 3.

We start with Determinism: all $n \to t \overline{n}$ for the same *n* to start with different *t*. To prove that, we again make use of the typing information. The following lemma says that a production can start with *t* if and only if it belongs to the FIRST set of the type.

LEMMA 4.6 (TERMINALS IN FIRST). Given $\Gamma; \Delta \vdash g : \tau$ and $\mathcal{N}[\![g]\!]$ returns $n \Rightarrow G$, we have $t \in \tau$. FIRST if and only if (1) $(n \rightarrow t \overline{n}) \in G$; or (2) $(n \rightarrow \alpha \overline{n}) \in G$ where $(\alpha : \tau') \in \Gamma$ and $t \in \tau$. FIRST.

This lemma is particularly important when proving the case for normalizing $g_1 \vee g_2$, where the typing condition $\tau_1 \# \tau_2$ ensures that g_1 and g_2 have disjoint FIRST, which in turn ensures that rule (*alt*) only copies distinct head terminals from g_1 and g_2 .

The proof for guarded ϵ -productions is more involved, which essentially requires us to show that during expansion $G \vdash n \rightsquigarrow \overline{t} n_1 n_2 \overline{n}$, the FIRST of n_1 is disjoint with the FIRST of n_2 (if n_1 is nullable). The proof replies on showing that "expansion preserves typing". Think about that from the well-typed context free expressions' point of view: if $(g_1 \lor g_2) \cdot g_3$ is well-typed, then $g_1 \cdot g_3$ (and $g_2 \cdot g_3$) must also be well-typed, and going from $(g_1 \lor g_2) \cdot g_3$ to $g_1 \cdot g_3$ is one step of branching, similar to one step of expansion. We refer the reader to the appendix for more details.

With all the conditions proved, we conclude our goal.

THEOREM 4.7 ($\mathcal{N}[\![g]\!]$ produces DGNF). If $\bullet; \bullet \vdash g : \tau$ and $\mathcal{N}[\![g]\!]$ returns $_ \Rightarrow G$, then G is DGNF.

4.4 Normalization Soundness

Normalization transforms a context-free expression into a DGNF grammar. As the final piece of metatheory for normalization, we would like to establish that normalization is sound with respect to the denotational semantics of well-typed context-free expressions defined in Krishnaswami and Yallop [2019]. The denotational semantics $[[g]]_{\gamma}$ interprets *g* as a language (i.e., the set of all strings w matched by *g*), where γ gives interpretation of free variables in *g*.

686

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

Fusing Lexing and Parsing

694 695

696

697 698

699

700

701

702

703

704

705

706

707

708 709 710

711

712

716

717

718

719

720

721

722 723

724

725

726

727 728

729

687	$\llbracket g \rrbracket_v$ inte	erpr	ets <i>g</i> as a language gi	ven the inter	pret	ation of the free variables in γ
688	$\llbracket \epsilon \rrbracket_{v}$	=	$\{\epsilon\}$	$\llbracket g_1 \cdot g_2 \rrbracket_V$	=	$\{w \cdot w' \mid w \in \llbracket g_1 \rrbracket_{Y} \land w' \in \llbracket g_2 \rrbracket_{Y}\}$
689	$\llbracket t \rrbracket_Y$	=	$\{t\}$	$\llbracket \alpha \rrbracket_{Y}$	=	$\gamma(a)$
690	$\llbracket \bot \mathring{\rrbracket}_{\gamma}$	=	Ø	$\llbracket \mu \alpha . g \rrbracket_{Y}$	=	$fix(\lambda L.[[g]]_{(\gamma, L/\alpha)})$
692	$\llbracket g_1 \lor g_2 \rrbracket_{\gamma}$	=	$\llbracket g_1 \rrbracket_\gamma \cup \llbracket g_2 \rrbracket_\gamma$	fix(f) =	=	$L_0 = \emptyset$
693				пх(j) -	i∈ĭ	$\bigcup_{i=1}^{n} \bigcup_{j=1}^{n} L_{n+1} = f(L_n)$

Most definitions are straightforward. The semantics interpret ϵ as the singleton set containing the empty set, t as the singleton set that matches a one-character string t, \perp as an empty language, and $q_1 \lor q_2$ as unions of sets. The interpretation of $q_1 \lor q_2$ appends a string from q_1 to a string from q_2 . Variables α have their interpretation from the environment γ , and $\mu\alpha$. q is interpreted as the least fixed point of *q* with respect to α .

To prove that our normalization is sound, we show that the normalized DGNF denotes exactly the same language as the denotation semantics of an expression. Recall that we have defined the expansion relation in Definition 2, where $G \vdash n \rightsquigarrow w$ denotes that n expands to a complete string w, where all non-terminals have been expanded. We prove the normalized grammar can expand to a string if and only if the string is included in the denotational semantics of the expression. The proof is done by first induction on the length of w and then on the structure of q.

THEOREM 4.8 (SOUNDNESS). Given $\bullet; \bullet \vdash q : \tau$ and N[[q]] returns $n \Rightarrow G$, we have $w \in [[q]]_{\bullet}$ if and only if $G \vdash n \rightsquigarrow w$ for any w.

4.5 Implementation

The compositionality of the normalization algorithm simplifies the implementation of normalization in flap. Since the normalization of each term is defined in terms of the normal forms of its subterms, 713 flap can represent terms in normal form. For example, if q and q' are flap parsers in normal form, 714 then q >>> q' is also a parser in normal form, built from q and q' using the rules in Figure 6. 715

The simplicity of the normalization algorithm is reflected in the implementation of flap. Additionally, the most intricate part of the algorithm – dealing with fixed points – is also the subtlest part of the implementation. The implementation follows the formal algorithm closely, inserting placeholders (α s) that are tracked using an environment and resolved later. This kind of "backpatching" mirrors the way in which recursion is commonly implemented in eager functional languages such as OCaml [Reynaud et al. 2021]; if flap were instead implemented in a lazy language then it would be possible to implement fixed point normalization with significantly less fuss.

5 FUSION

This section shows how flap, starting with a separately-defined lexer and normalized parser, makes use of the syntactic restrictions of DGNF to implement lexer-parser fusion, eliminating tokens from generated code altogether.

Canonicalizing lexer 5.1

We assume a more restrictive definition of lexers than the interface in Section 2 provides. In 730 particular, we assume that rules are disjoint on the left (i.e. there is no string that is matched by 731 more than one regular expression in a set of rules), and disjoint on the right (i.e. there is exactly one 732 Skip rule, and no token appears in more than one Return rule). Using negation and intersection, it 733 is easy to transform a lexer that does not obey these constraints into an equivalent lexer that does, 734

fused grammar $F ::= \{n \to r \overline{n}\} \cup \{n \to ?r\}$

$$\begin{array}{ll} \mathcal{F}[\![L,G]\!] = \mathcal{F}_1 \cup \mathcal{F}_2 \cup \mathcal{F}_3 \\ \text{where } \mathcal{F}_1 = \{n \rightarrow r\bar{n} \mid r \Rightarrow \text{Return } t \in L \land n \rightarrow t \, \bar{n} \in G\} \\ \mathcal{F}_2 = \{n \rightarrow rn \mid r \Rightarrow \text{Skip} \in L \land n \in G\} \\ \mathcal{F}_3 = \{n \rightarrow ? \neg r \mid n \rightarrow \epsilon \in G \land r = \bigvee\{r \mid n \rightarrow r\bar{n} \in \mathcal{F}_1 \cup \mathcal{F}_2\}\} (epsilon \ productions) \\ \mathcal{F}_4 \\ \mathcal{F}_3 = \{n \rightarrow ? \neg r \mid n \rightarrow \epsilon \in G \land r = \bigvee\{r \mid n \rightarrow r\bar{n} \in \mathcal{F}_1 \cup \mathcal{F}_2\}\} (epsilon \ productions) \\ \mathcal{F}_4 \\ \mathcal{F}_4 \\ \mathcal{F}_4 \\ \mathcal{F}_5 \\ \text{so there is no need to restrict the interface exposed to the user. With the lexer thus canonicalized and the parser translated into DGNF, it is straightforward to define fusion. \\ \mathcal{F}_4 \\ \mathcal{F}_4 \\ \mathcal{F}_4 \\ \mathcal{F}_5 \\ \mathcal{F}_2 \\ \mathcal{F}_4 \\ \mathcal{F}_5 \\ \mathcal{F}_2 \\ \mathcal{F}_5 \\ \mathcal{F}_5$$

Figure 8 formally defines the fusion algorithm.

The fusing function $\mathcal{F}[\![L, G]\!]$ operates on a canonicalized lexer L and a normalized grammar G, yielding a fused grammar F. The fused result consists of three parts.

First, we replace each production $n \to t \overline{n}$ with a new production $n \to r \overline{n}$, retrieving the regex r that is associated with the token t in the lexer $L(\mathcal{F}_1)$. This is where the fusion function implicitly specializes the lexer to each nonterminal in the normalized grammar, and discards lexing rules that return tokens not in productions for the nonterminal. Canonicalizing the lexer to enforce disjointness simplifies this discarding of rules.

Then, we add an additional production $n \to r n$ for the **skip** regex r (which may be \perp) for each nonterminal, allowing each nonterminal to match an arbitrary number of the skip regex (\mathcal{F}_2).

Finally, for nonterminals with an epsilon production, the discarded regexes, along with the skip regex, are incorporated into a lookahead regex (\mathcal{F}_3). That is, we add a lookahead production $n \rightarrow ?\neg r$ for the regex that is the complement of the regexes that appear in other productions for *n*.

Fusion with normalized grammars is strikingly simple; it would be significantly more involved to directly fuse the context-free expressions with the lexing rules. Furthermore, as with normalized grammars, an expansion relation for fused grammars would carry the guarantee that every expansion has a unique derivation.

6 IMPLEMENTATION OF PARSING

This section describes the lexing and parsing algorithms, shows how to stage the parsing algorithm to improve performance, and explains details of the implementation of the algorithms in flap.

6.1 The lexing algorithm

Figure 9 presents the lexing algorithm. The algorithm has the *longest-match* semantics conventional for lexers: each token returned corresponds to the lexing rule that matches the longest possible prefix of the input string. This behaviour is implemented by repeatedly updating the best match seen *so far* until none of the lexing rules matches the input string.

 $\mathcal{L}ex$ is the top-level lexing algorithm that takes the lexing rules L and an input string s, with two key utility functions \mathcal{L} and \mathcal{M} . For simplicity, we assume utility functions can freely access the L argument to $\mathcal{L}ex$. At a high level, \mathcal{L} reads a single token from a prefix of a string, pairs the token action with the remainder of the string, and passes it to \mathcal{M} . \mathcal{M} constructs a sequence of tokens, updating the sequence according to the action passed from \mathcal{L} .

The \mathcal{L} function takes four arguments: L' is a set of lexing rules; k is a token action representing the best match so far; rs is the remainder string for the best match; s is the input string. For empty input strings the best match information k and rs is passed to \mathcal{M} . For non-empty input strings

 $\mathcal{L}ex(L,s) = \mathcal{L}(L, \mathrm{NO}, [], s)$ 785 786 $\mathcal{L}(L', k, rs, []) = \mathcal{M}(k, rs)$ $\mathcal{M}(NO, rs)$ = FAIL 787 $\mathcal{L}(L', k, rs, c::cs) = \text{if } L'_c \stackrel{?}{=} \emptyset \text{ then } \mathcal{M}(k, rs)$ $\mathcal{M}(Skip, [])$ = [] 788 $= \mathcal{L}(L, \text{NO}, [], c::cs)$ else case K of $\emptyset \mapsto \mathcal{L}(L'_c, k, rs, cs)$ $\mathcal{M}(\text{Skip}, c::cs)$ 789 $\mathcal{M}(\text{Return } t, [])$ = [t] $\{k'\} \mapsto \mathcal{L}(L'_c, k', cs, cs)$ 790 where $L'_c = \{\partial_c(r) \Rightarrow k \mid r \Rightarrow k \in L' \land \partial_c(r) \neq \bot\}$ $K = \{k \mid r \Rightarrow k \in L'_c \land v(r)\}$ $\mathcal{M}(\text{Return } t, c::cs) = t :: \mathcal{L}(L, \text{NO}, [], c :: cs)$ 791 792 793 Fig. 9. Lexing algorithm 794 795 $\mathcal{P}arse(n \Rightarrow G, s) = \mathcal{P}(n, s)$ 796 $\mathcal{P}(n, []) = \text{if } n \to \epsilon \in G \text{ then } [] \text{ else FAIL}$ Q([], ts)= ts797 $Q(n::ns,ts) = Q(ns,\mathcal{P}(n,ts))$ $\mathcal{P}(n, t::ts) = \text{if } n \to t\overline{n} \in G \text{ then } Q(\overline{n}, ts)$ 798 else if $n \to \epsilon \in G$ then *t*::*ts* else FAIL 799 800 Fig. 10. Parsing algorithm for DGNF grammars 801 802 803 $\mathcal{FP}arse(n \Rightarrow F, s) = \mathcal{G}([n], s)$ 804 805 $\mathcal{F}(F_n, k, rs, s) =$ $\mathcal{G}([],s) = s$ 806 case s of $[] \mapsto Step(k, rs)$ $\mathcal{G}(n::ns,s) = \mathcal{G}(ns,\mathcal{F}(F_n,k,s,s))$ 807 $c::cs \mapsto \text{if } F'_n \stackrel{?}{=} \emptyset \text{ then } Step(k, rs)$ where $F_n = \{ \langle r, \overline{n} \rangle \mid n \to r\overline{n} \in F \}$ 808 $k = \text{if } n \rightarrow ?r \in F \text{ then BACK else NO}$ else case K of $\emptyset \mapsto \mathcal{F}(F'_n, k, rs, cs)$ 809 $\{ns\} \mapsto \mathcal{F}(F'_n, \text{ on } ns, cs, cs)$ Step(BACK, s) = s810 where $F'_n = \{ \langle \partial_c(r), k \rangle \mid \langle r, k \rangle \in F_n \land \partial_c(r) \neq \bot \}$ $Step(ON ns, s) = \mathcal{G}(ns, s)$ 811 $\widetilde{K} = \{k \mid \langle r, k \rangle \in F'_n \land v(r)\}$ Step(NO, s)= FAIL 812 813 814 Fig. 11. Parsing algorithm for fused grammars 815 $SParse_{n \Rightarrow F}(s) = \mathcal{T}([n], s)$ 816 817 818 $S_{F_{n,k}}(rs,s) =$ $\mathcal{T}([], s) = s$ 819 $\mathcal{T}(n::ns, s) = \mathcal{T}(ns, S_{F_n,k}(s, s))$ case s of $[] \mapsto Step(k, rs)$ 820 where $F_n = \{ \langle r, \overline{n} \rangle \mid \overline{n \to r\overline{n} \in F} \}$ $c_i :: c_s \mapsto \text{if } F'_{n_i} \stackrel{?}{=} \emptyset \text{ then } Step(k, r_s)$ 821 $k = \text{if } n \rightarrow ?r \in F \text{ then BACK else NO}$ 822 else case K_i of $\emptyset \mapsto S_{F'_n,k}(rs, cs)$ 823 Step(BACK, s) = s $\{ns\}\mapsto S_{F'_{n,i},\text{ON }ns}(cs,cs)$ 824 $Step(ON ns, s) = \mathcal{T}(ns, s)$ $c_j :: c_s \mapsto \ldots$ 825 Step(NO, s) = FAILwhere $F'_{n,i} = \{ \langle \partial_{c_i}(r), k \rangle \mid \langle r, k \rangle \in F_n \land \partial_{c_i}(r) \neq \bot \}$ $K_i = \{ k \mid \langle r, k \rangle \in F'_{n,i} \land v(r) \}$ 826 827 828 829 Fig. 12. Staged parsing algorithm 830 831 832 833 Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018. c::cs, the result depends on L'_c , the set of lexing rules updated to use the non-empty *derivatives* with respect to c (Figure 2) of the string. If L' is empty, lexing cannot proceed any further, and so \mathcal{L} transfers control to \mathcal{M} , passing the best match information. Otherwise, the result depends on the rule $r \Rightarrow a$ that matches the string up to this point including c (i.e. the rule that accepts ϵ after consuming c). If there is no such rule, then lexing continues with k. If there is such a rule, it is unique (since lexing rules are disjoint (Section 2.2)), and it represents a new longest-match k', and lexing continues with k' and the remainder for the best match cs.

The \mathcal{M} function takes two arguments k, a token action and rs, a remainder string. There are five cases, one for the sentinel NO action, two for Skip actions, and two for Return actions. The sentinel NO indicates that lexing has failed. For Skip, lexing continues if the remainder rs is non-empty. For Return t, t is added to the output sequence, and lexing continues if the remainder rs is non-empty. In the cases where lexing continues, it commences by supplying NO for the best-match-so-far, so that reading the next token only succeeds if \mathcal{L} matches a non-empty prefix of the remaining input.

848 6.2 The DGNF parsing algorithm

Figure 10 presents the parsing algorithm for grammars in Deterministic Greibach Normal Form.
Deterministic parsing makes the algorithm simple, since there is no need for backtracking.

 \mathcal{P} arse is the top-level parsing algorithm which takes the parsing grammar $n \Rightarrow G$ and a sequence of tokens *s*. There are two key functions: \mathcal{P} parses using a single nonterminal *n*, and *Q* parses using a sequence of nonterminals *ns*. Again, we assume \mathcal{P} and *Q* can freely access *G*.

P takes the nonterminal *n* and a sequence of tokens and returns the remainder of the sequence after parsing. For empty sequences parsing succeeds only if the grammar has a rule $n \to \epsilon$. For non-empty sequences *t*::*ts*, if the grammar has a rule $n \to t\bar{n}$, \mathcal{P} consumes *t* and parses *ts* with *Q*. Otherwise, parsing succeeds (consuming nothing) only if the grammar has a rule $n \to \epsilon$.

Q takes a sequence of nonterminals *ns* and a sequence of tokens *ts* and parses successive prefixes of *s* with each nonterminal in *ns*.

6.3 The parsing algorithm for fused grammars

In practice, flap does not need separately-defined lexing and DGNF parsing algorithms, since it
 fuses lexing and parsing. We presented those algorithms to allow a direct comparison with the
 parsing algorithm for fused grammars.

Figure 11 shows an algorithm for parsing with fused grammars. The algorithm combines the features of the lexing algorithm (Figure 9) and the parsing algorithm (Figure 10): like the lexing algorithm it maintains a set of derivatives and an action and remainder string for the current *best match*; like the parsing algorithm, it keeps track of the current non-terminal.

FP arse takes the fused grammar $n \Rightarrow F$ and an input string *s*, with two key functions: *F* parses using a single nonterminal *n*, and *G* parses using a sequence of nonterminals *ns* using *F*.

 \mathcal{F} takes four arguments: F_n , a set of pairs representing non-epsilon productions for *n*; *k*, an action; 871 rs, a remainder string; and s, an input string. For empty input strings the best match information k872 and rs is passed to \mathcal{G} (via the auxiliary function Step). For non-empty input strings c::cs, the result 873 depends on F'_n , the production pairs for *n* updated to use the non-empty *derivatives with respect to* 874 *c* (Figure 2) of the string. If F'_n is empty, parsing cannot proceed any further, and so \mathcal{F} transfers 875 control to \mathcal{G} (via *Step*), passing the best match information. Otherwise, the result depends on the 876 production pair $\langle r, \overline{n} \rangle$ for which r matches the string up to this point including c (i.e. the rule that 877 accepts ϵ after consuming c). If there is no such rule, then parsing continues with k. If there is such 878 a rule, it is unique (since the regular expressions for a particular nonterminal are disjoint), and it 879 represents a new longest-match ns, and parsing continues, updating the best match information to 880 ON \overline{ns} . Here ON \overline{ns} represents one of three continuation types, and indicates that parsing should 881

847

851

852

853

860

861

continue using the nonterminal sequence \overline{ns} ; the others are back, indicating that parsing with *n* should succeed, consuming no input, and no, indicating that parsing with *n* should fail. The *Step* function matches these three cases, and takes an action appropriate to each continuation.

The G function takes a sequence of nonterminals ns and a sequence of characters s and parses successive prefixes of s with each nonterminal in ns by calling \mathcal{F} . The value of \mathcal{F} 's k argument depends on whether there is an epsilon rule for n in the fused grammar: if so, then a parsing failure with n should backtrack, consuming no input; if not, then parsing returns FAIL.

We draw attention to two salient features of the fused parsing algorithm: first, it consists of elements from the lexing and parsing algorithms of Sections 6.1 and 6.2; second, it does not materialize the tokens produced by the lexing algorithm, instead operating directly on the character string. The final algorithm in the next section makes this even more apparent.

6.4 The staged parsing algorithm

The parsing algorithm for fused grammars described in Section 6.3 is impractically inefficient. For each character of the input, the algorithm computes derivatives and checks emptiness and nullability for sets of regular expressions. However, since the regular expressions and other information about the grammar are known in advance of parsing, the inefficient algorithm can be *staged* [Taha 1999] to produce an efficient algorithm. The idea of staging is to identify those parts of the algorithm that do depend only on static information – i.e. on the grammar – and execute them first, leaving only the parts that depend on dynamic information – i.e. on the input string – for later. The result of staging, as illustrated below, is to transform the unstaged parser into a parser generator that produces as output a parser specialized to the input grammar.



Figure 12 shows a staged version of the fused parsing algorithm. The structure of the algorithm is very close to the fused grammar parsing algorithm of Section 6.3: S corresponds to \mathcal{F} and \mathcal{T} corresponds to \mathcal{G} . However, there are three key differences.

First, those parts of the algorithm that depend on the input string are marked as *dynamic*, indicated with red highlighting. These dynamic elements are not executed immediately; instead they become part of the generated specialized parser produced by the first stage of execution.

Second, in the function S, F_n and k have become indexes rather than arguments. Consequently, rather than being passed to the function at run-time, those arguments serve to distinguish generated functions: each instantiation of F_n and k generate a distinct function S in the specialized parser.

Finally, the case match in S is expanded to include a distinct case for each character c_i, c_j , etc. This expansion resolves a tension in the distinction between static and dynamic data: the static computation of derivatives $\partial_c(r)$ in the first stage depends on the value of c, which is only available dynamically. In the expanded case match the value of c_i is known on the right-hand side of the corresponding case, making it possible to compute derivatives valid within that program context. This scrutiny of a statically-unknown expression using a case match over its statically-known set of possible values is known as "The Trick" in the partial evaluation literature [Danvy et al. 1996].

The evaluation of the staged parsing algorithm is largely standard: the unhighlighted (static) expressions are executed first, producing the highlighted (dynamic) expressions as output. Each call to a dynamic indexed function $S_{F_{n,k}}$ triggers the generation of a dynamic function whose body consists of the result of executing the right-hand side of $S_{F_{n,k}}$ in Figure 12. To ensure that the generation process terminates, the generation of these indexed functions is memoized: there is at

most one generated function $S_{F_n,k}$ for any particular F_n and k. The result of the algorithm is a set of mutually recursive functions that operate only on strings, not on components of the grammar:

 1:20

```
S_{n \to r\overline{n}, \dots, \text{BACK}}(r, s) = \text{case } s \text{ of } [] \qquad \mapsto s
 \text{'a'} :: cs \mapsto S_{n \to r_a \overline{n}, \text{BACK}}(r, cs)
 \text{'b'} :: cs \mapsto S_{n \to r_a \overline{n}, \text{on } \overline{ns}}(cs, cs)
 \cdots
S_{n \to r\overline{n}, \dots, \text{on } \overline{ns}}(r, s) = \cdots
```

6.5 Implementing the staged parsing algorithm

The flap library generates code for the fused grammar using MetaOCaml's staging facilities together with Yallop and Kiselyov's [2019] *letrec insertion* library for creating the indexed mutually-recursive functions produced by the staged parsing algorithm (Section 6.4).

There are three key differences between the pseudocode algorithm in Figure 12 and flap's implementation. First, while the pseudocode presents a recognizer that either consumes input or fails, flap supports *semantic actions* – i.e. constructing and returning ASTs or other values when parsing succeeds – as described in Section 2.3.

Second, while the input to the pseudocode is a linked list of characters, flap operates on OCaml's more efficient flat array representation of strings, using indexes to keep track of string positions as parsing proceeds. Relatedly, flap also optimizes the test for the end of input by taking advantage of the fact that OCaml's strings are null-terminated, to ease interoperability with C. This representation allows the check for end of input to be incorporated into the branch on the next character in the generated code: a null character '\000' indicates a *possible* end of input, which can subsequently be confirmed by checking the string length.

Third, while the pseudocode generates a case in each branch for each possible character in the input, flap generates a much smaller number of cases by grouping characters with equivalent behaviour into classes. Branching on these character classes rather than treating characters individually leads to a substantial reduction in code size. Owens et al. [2009] describe the construction of these classes in described in detail.

Here is an excerpt of the code generated by flap for the s-expression parser:

```
and parse<sub>5</sub> r i len s = match s.[i] with

| ' '|' \setminus n' \rightarrow parse_6 r (i + 1) len s

| '(' \rightarrow parse_9 r (i + 1) len s

| 'a'...'z' \rightarrow parse_3 r (i + 1) len s

| '\setminus 000' \rightarrow if i = len then [] else failwith "unexpected"

| \_ \rightarrow []
```

This excerpt shows the code generated for a single indexed function $S_{F_n,k}$. There are four arguments, representing the beginning of the current token r (to support backtracking in the lookahead transition), the current index i, the input length len, and the input string s.

The subscripts 5, 6, etc. attached to the parse functions correspond to the indexes F_n , k in the pseudocode algorithm; the letrec insertion library assigns a fresh subscript to each distinct index.

The character range pattern 'a'...'z' illustrates the character class optimization described above.
Without that optimization, each of the characters from 'a' to 'z' would have a separate case in
the match expression.

The check i = len determines whether the character '\000' represents the end of input or a null character in the input string.

The value [] corresponds to a semantic action: it is the empty list returned when an empty sequence of s-expressions is parsed. It appears twice in the generated code, since (as Figure 12

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

shows), parsing for a particular nonterminal can end in two ways: when it encounters the end of 981 input, and when it encounters a non-matching character. 982

EVALUATION 7 984

983

987

995

996

997 998

1020

1029

985 This section evaluates the performance of flap, and shows that lexer-parser fusion drastically improves performance. Many parser combinator libraries suffer from poor performance, but the 986 experiments described here show that combinator parsing does not need to be slow.

988 Three key techniques account for flap's speed. First, Krishnaswami and Yallop's type system ensures that the time taken for parsing is linear in the length of the input, a substantial advantage over 989 libraries that require backtracking. Second, staging eliminates the overhead arising from parsing ab-990 stractions, generating parsing code that is specialized to a particular grammar. Finally, lexer-parser 991 fusion eliminates the overhead arising from defining lexers and parsers separately: in particular, it 992 993 avoids the materialization of tokens, and eliminates all branching except for approximately one branch on each character in the input. 994

Krishnaswami and Yallop showed that the first two of these techniques can be used to build a parser combinator library that outperforms code generated by ocamlyacc. We focus here on the question of how much additional performance benefit arises from lexer-parser fusion.

Benchmarks 999 7.1

1000 We build on the benchmark suite published by Krishnaswami and Yallop [2019], adding implemen-1001 tations of each benchmark for flap, for the parser generator menhir [Pottier and Régis-Gianas 1002 [n.d.]], and for the ParTS deterministic parsing library [Casinghino and Roux 2020], and extending the suite with an additional benchmark for parsing CSV files. 1003

1004 For each benchmark we compare up to six implementations. We make no comparison with 1005 unstaged parser combinators, which Krishnaswami and Yallop found to be significantly slower 1006 than their staged implementation, and between 4.5 and 125 times slower than ocamlyacc.

- 1007 (a) an implementation generated by ocamllex and ocamlyacc
- 1008 (b) an implementation generated by ocamllex and menhir in table-generation mode
- 1009 (c) an implementation generated by ocamllex and menhir in code-generation (tableless) mode 1010
- (d) an implementation created using flap 1011
 - (e) the staged parser implementation from Krishnaswami and Yallop [2019]
- 1012 (f) an implementation using ParTS, where one is available 1013

(a)-(c) use identically structured grammars and lexers in each benchmark, since menhir accepts 1014 ocamlyacc files as input. (d)-(f) also use identically structured grammars in each benchmark, 1015 since they all use the standard parser combinator interface (Section 2.3). However, (d)-(f) use 1016 differently-structured lexers: (e) and (f), taken respectively from Krishnaswami and Yallop [2019] 1017 and Casinghino and Roux [2020], reuse the deterministic parser combinators for lexing, while flap 1018 uses the more conventional lexing interface described in Section 2.2. 1019

The benchmarks are as follows:

- (1) (pgn) Chess game descriptions in Portable Game Notation format. The semantic actions extract 1021 the result of each game. The input is a corpus of 6759 Grand Master games. 1022
- (2) (ppm) Image files in Netphin format. The semantic actions validate the non-syntactic constraints 1023 of the format, such as colour range and pixel count. 1024
- (3) (sexp) S-expressions with alphanumeric atoms. The semantic actions of the parser count the 1025 number of atoms in each s-expression. 1026
- (4) (csv) A parser for the Comma-Separated Value format. The grammar conforms quite closely to 1027 RFC 4180 [Shafranovich 2005], but makes the terminating CRLF mandatory and does not treat 1028





Fig. 13. Parser throughput: ocamlyacc, menhir, flap, staged combinators and ParTS



headers specially. The semantic actions check that each row contains the same number of fields.
There is no implementation using Krishnaswami and Yallop's combinators for this benchmark,
because more than a single character of lookahead is needed to distinguish escaped (i.e. repeated)
double-quotes "" from unescaped quotes ", and so the lexer cannot be implemented with typed
context-free expressions without substantial changes to its structure. The lexer interface used
in flap (Section 2.2) does not suffer from this limitation.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1096 1097

1102

1103 1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1120

	Input		Norm	nalized	Fused	Output	
Grammar	Lexer rules	Context-free exps	Nonterms Productions		Productions	Functions	
pgn	13	95	38	53	91	206	
ppm	6	10	5	6	16	55	
sexp	4	11	3	6	9	11	
csv	3	14	5	7	7	20	
json	12	42	9	33	42	97	
arith	14	143	28	55	83	209	
1		<i>.</i>			·		

Table 1. Sizes of inputs, intermediate representations, and generated code

Benchmark	Compilation time (ms)
sexp	0.331
pgn	212
ppm	3.60
json	28.5
csv	0.499
arith	460

Table 2. Compilation time (type-checking, normalization, fusion, code generation)

- (5) (json) A parser for JavaScript Object Notation (JSON). The semantic actions count the number of objects represented in the input. Following Krishnaswami and Yallop, we use the simple JSON grammar given by Jonnalagedda et al. [2014].
 - (6) (arith) A miniature programming language with constructs for arithmetic, comparison, letbinding and branching. The semantic actions evaluate the parsed expression.

7.2 Running time

Figure 13 shows the absolute and relative throughput of the six implementations using the six benchmark grammars. For the benchmarks that are taken from [Krishnaswami and Yallop 2019] we use the test corpora from the same source. For the CSV benchmark we have generated a set of files of various sizes and dimensions, using a random variety of textual and numeric data.

The benchmarks were compiled with BER MetaOCaml N111 with flambda optimizations enabled and run on a single Intel i9-12900K core with 8GB memory running Ubuntu Linux, using the Core_bench micro-benchmarking library [Hardin and James 2013].

As the graph shows, our experiments confirm the results reported by Krishnaswami and Yallop: the staged implementation of typed context-free expressions generally outperforms ocamlyacc. The addition of lexer-parser fusion makes flap considerably faster than both typed CFEs and ocamlyacc, reaching around 1.3GB/s (a little under 2.5 cycles per byte) on the json benchmark.

Linear-time parsing. Finally, as Figure 14 illustrates, on Krishnaswami and Yallop's benchmark
suite, parsers built with flap, like parsers built with Krishnaswami and Yallop's combinators,
execute in time linear in the length of their input.

1121 7.3 Code size

A second important measure of usefulness for parsing: if parsing tools are to be usable in practice,it is essential that they do not generate unreasonably large code.

There are several reasons to be apprehensive about the size of code generated by flap. First, conversion to Greibach Normal Form is well known to substantially increase the size of grammars; for example, in the procedure given by Blum and Koch [1999] the result of converting a grammar

1:23

G has size $O(|G|^3)$. Second, the fusion process is inherently duplicative, repeatedly copying the lexer rules into the various grammar productions. Finally, experience in the multi-stage programming community shows that it is easy to inadvertently generate extremely large programs, since antiquotation makes it easy to duplicate terms.

However, measurements largely dispel these concerns. Table 1 gives the size of the representations 1132 of the benchmark parsers at various stages in flap's pipeline. The leftmost pair of columns of 1133 figures shows the size of the input parsers, measured as the number of lexer rules (including 1134 1135 both **Return** and **Skip** rules) and the number of context-free expression nodes, as described in Section 2. The pair of columns to the right shows the number of nonterminals and productions after 1136 the grammar is converted to Deterministic GNF using the procedure in Section 4. As the figures 1137 show, the normalization algorithm for typed context-free expressions does not produce the drastic 1138 increases in size that occurs in the more general form of conversion to GNF. The next column to 1139 the right shows the size of the grammar after fusion (Section 5). Fusion does not alter the number 1140 of nonterminals, but it can add productions: for example, the Skip rules in the s-expression lexer 1141 add additional productions to each nonterminal. Finally, the rightmost column shows the number 1142 of function bindings in the code generated by flap. Comparing this generated function count with 1143 the number of context-free expressions in the input reveals a fairly unalarming relationship: with 1144 one exception (pgn), the ratio between the two barely exceeds 2. 1145

Sharing. The entries for pgn and arith hint at opportunities for further improvement. In both
 cases, the number of context-free expressions that make up the grammar (95 and 143) is surprisingly
 high, since both languages are fairly simple. Inspecting the implementations of the grammars reveals
 the cause: in several places, the combinators that construct the grammar duplicate subexpressions.
 For example, here is the implementation of a Kleene plus operator used in pgn:

1152 **let** oneormore e = (e >>> star e) ...

Normalization turns these two occurrences of e into multiple entries in the normalized form, and ultimately to multiple functions in the generated code.

The core problem is that the parser combinator interface (Section 2.3) provides no way to express sharing of subgrammars. Since duplication of this sort is common, it is likely that extending flap with facilities to express and maintain sharing could substantially reduce generated code size.

7.4 Compilation time

A final measure of practicality is the time taken to perform the fusion transformation. Slow compilation times can have a significant effect on usability; as Nielsen [1993] notes, software that takes more than one second to respond can cause a user to lose concentration, harming interactivity.

Table 2 shows the compilation time for the six benchmark grammars. In each case, the total time taken to type-check and normalize the grammar, fuse the grammar and lexer and generate optimized code is less than half a second.

8 RELATED WORK

Deterministic Greibach Normal Form. There are several longstanding results related to determin istic variants of Greibach Normal Form. For example, Geller et al. [1976] show that every strict
 deterministic language can be given a strict deterministic grammar in Greibach Normal Form, and
 Nijholt [1979] gives a translation into Greibach Normal Form that preserves strict deterministicness.
 The distinctive contributions of this paper are the new normal form that is well suited to fusion,
 and the compositional normalization procedure from typed context-free expressions, allowing
 deterministic GNF to be used in the implementation of parser combinators.

1176

1146

1153

1154

1159

1164

1165

1166

Combining lexers and parsers. The work most closely related to ours, by Casinghino and Roux 1177 [2020] investigates the application of traditional stream fusion techniques to parser combinators in 1178 1179 the ParTS system. We have included their two published benchmarks in the evaluation of Section 7 and found that, as they report, when the flambda suite of compiler optimizations is applied to 1180 their code, its performance is similar to the results achieved by Krishnaswami and Yallop [2019]. A 1181 significant difference between their work and ours is that they approach fusion as a traditional 1182 optimization problem, in which transformations are applied to code that satisfies certain heuristics, 1183 1184 and are not applied in more complex cases. In contrast, we treat lexer-parser fusion as a sequence of total transformations that is guaranteed to convert every input (i.e. every parser) into a form that 1185 enjoys pleasant performance properties. More concretely, Figure 13 shows significant performance 1186 differences between ParTS and flap, with ParTS achieving one half and a tenth of the throughputs 1187 of flap on the sexp and json benchmarks. 1188

Another line of work, on *Scannerless GLR parsing* [Economopoulos et al. 2009; van den Brand et al.
2002], also aims to eliminate the boundary between lexers and parsers, both in the interface and
the implementation. The principal aim is to provide a principled way to handle lexical ambiguity,
in contrast to our focus on performance.

Context-aware scanning, introduced by Van Wyk and Schwerdfeger [2007] is another variant on the parser-scanner interface focused on disambiguation; it passes contextual information from the parser to the scanner about the set of valid tokens at a particular point, in a similar way to the lexer specialization in Section 3.3 of this paper. However, Van Wyk and Schwerdfeger's framework goes further, and allows the automatic selection of a lexer (not just a subset of lexing rules) based on the parsing context.

Fusion. The notion of fusion, in the sense of merging computations to eliminate intermediate structures, has been applied in several domains, including query engines [Shaikhha et al. 2018], GPU kernels [Filipovic et al. 2015] and tree traversals [Sakka et al. 2019].

Perhaps the most widespread is stream fusion, which appears to have originated with Wadler's deforestation [Wadler 1990], and has since been successfully applied as both a traditional compiler optimization [Coutts et al. 2007] and as a staged library [Kiselyov et al. 2017] that provides guarantees similar to those we give here for parsers.

Parser optimization. Finally, in contrast to the constant-time speedups resulting from lexerparser fusion, we note an intriguing piece of work by Klyuchnikov [2010] that applies two-levelsupercompilation to parser optimization, leading to asymptotic improvements.

9 FUTURE WORK

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209 1210

1211

1212

1213

1214

1215

There are a number of promising avenues for future work. First, extending flap's rather minimal lexer and parser interfaces to support common needs such as left-recursive grammars, lexers and parsers with multiple entry points, mechanisms for maintaining state during parsing, and more expressive lexer semantic action could make the library substantially more usable in practice.

Second, building on the proofs of normalization correctness in Section 4 to cover the whole of
 flap, we plan to formally establish that the code generated by Section 5 faithfully represents the
 semantics of the combinators in Section 2.

Third, applying the ideas in this paper to more powerful parsing algorithms such as LR(1) and LALR(1), and incorporating them into traditional standalone parser generator (rather than a staged library) could make lexer-parser fusion available to many more software developers.

Finally, it may be that fusion can be extended to longer pipelines than the lexer-parser interface
that we investigate here. Might it be possible to fuse together (e.g.) decompression, unicode decoding,
lexing and parsing into a single computation that does not materialize intermediate values?

1:26

- Norbert Blum and Robert Koch. 1999. Greibach Normal Form Transformation Revisited. Information and Computation 150, 1 (1999), 112–118. https://doi.org/10.1006/inco.1998.2772
- Anders Bondorf. 1992. Improving Binding Times Without Explicit CPS-conversion. In *Proceedings of the 1992 ACM Conference on LISP and Functional Programming* (San Francisco, California, USA) (*LFP '92*). ACM, New York, NY, USA, 1–10. https://doi.org/10.1145/141471.141483
- 1231
 Janusz A. Brzozowski. 1964. Derivatives of Regular Expressions. J. ACM 11, 4 (Oct. 1964), 481–494. https://doi.org/10.1145/

 1232
 321239.321249
- 1233 Chris Casinghino and Cody Roux. 2020. ParTS: Final Report. HR001120C0016 Final Report.
- Duncan Coutts, Roman Leshchinskiy, and Don Stewart. 2007. Stream fusion: from lists to streams to nothing at all. In
 Proceedings of the 12th ACM SIGPLAN International Conference on Functional Programming, ICFP 2007, Freiburg, Germany,
 October 1-3, 2007, Ralf Hinze and Norman Ramsey (Eds.). ACM, 315–326. https://doi.org/10.1145/1291151.1291199
- Olivier Danvy, Karoline Malmkjær, and Jens Palsberg. 1996. Eta-Expansion Does The Trick. ACM Trans. Program. Lang.
 Syst. 18, 6 (1996), 730–751. https://doi.org/10.1145/236114.236119
- Rowan Davies and Frank Pfenning. 1996. A Modal Analysis of Staged Computation. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (St. Petersburg Beach, Florida, USA) (*POPL '96*).
 Association for Computing Machinery, New York, NY, USA, 258–270. https://doi.org/10.1145/237721.237788
- Giorgios Economopoulos, Paul Klint, and Jurgen J. Vinju. 2009. Faster Scannerless GLR Parsing. In *Compiler Construction*, 18th International Conference, CC 2009, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, York, UK, March 22-29, 2009. Proceedings (Lecture Notes in Computer Science, Vol. 5501), Oege de Moor and Michael I. Schwartzbach (Eds.). Springer, 126–141. https://doi.org/10.1007/978-3-642-00722-4 10
- Jiri Filipovic, Matus Madzin, Jan Fousek, and Ludek Matyska. 2015. Optimizing CUDA code by kernel fusion: application on BLAS. *J. Supercomput.* 71, 10 (2015), 3934–3957. https://doi.org/10.1007/s11227-015-1483-z
- Matthew M. Geller, Michael A. Harrison, and Ivan M. Havel. 1976. Normal forms of deterministic grammars. *Discret. Math.* 16, 4 (1976), 313–321. https://doi.org/10.1016/S0012-365X(76)80004-0
- Sheila A. Greibach. 1965. A New Normal-Form Theorem for Context-Free Phrase Structure Grammars. J. ACM 12, 1 (Jan. 1965), 42–52. https://doi.org/10.1145/321250.321254
- 1249 Christopher S. Hardin and Roshan P. James. 2013. Core_bench: Micro-Benchmarking for OCaml. OCaml Workshop.
- Manohar Jonnalagedda, Thierry Coppey, Sandro Stucki, Tiark Rompf, and Martin Odersky. 2014. Staged Parser Combinators for Efficient Data Processing. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications* (Portland, Oregon, USA) (*OOPSLA '14*). ACM, New York, NY, USA, 637–653. https: //doi.org/10.1145/2660193.2660241
- 1253Oleg Kiselyov. 2014. The Design and Implementation of BER MetaOCaml System Description. In *FLOPS 2014 (LNCS,*
Vol. 8475), Michael Codish and Eijiro Sumii (Eds.). Springer, 86–102. https://doi.org/10.1007/978-3-319-07151-0_6
- Oleg Kiselyov, Aggelos Biboudis, Nick Palladinos, and Yannis Smaragdakis. 2017. Stream fusion, to completeness. In Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017, Giuseppe Castagna and Andrew D. Gordon (Eds.). ACM, 285–299. https://doi.org/10.1145/3009837
 Ilya Klyuchnikov. 2010. Towards effective two-level supercompilation. Preprint 81. Keldysh Institute of Applied Mathematics,
- 1258 Moscow.
- Neelakantan R. Krishnaswami and Jeremy Yallop. 2019. A typed, algebraic approach to parsing, See [McKinley and Fisher
 2019], 379–393. https://doi.org/10.1145/3314221.3314625
- Kathryn S. McKinley and Kathleen Fisher (Eds.). 2019. Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019, Phoenix, AZ, USA, June 22-26, 2019. ACM. https://doi.org/10.1145/3314221
 Jakob Nielsen. 1993. Usability Engineering. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- Kristian Nielsen and Morten Heine Sørensen. 1995. Call-By-Name CPS-Translation As a Binding-Time Improvement. In
 Proceedings of the Second International Symposium on Static Analysis (SAS '95). Springer-Verlag, London, UK, UK, 296–313.
 http://dl.acm.org/citation.cfm?id=647163.717677
- Anton Nijholt. 1979. Strict Deterministic Grammars and Greibach Normal Form. J. Inf. Process. Cybern. 15, 8/9 (1979), 395–401.
- Scott Owens, John H. Reppy, and Aaron Turon. 2009. Regular-expression derivatives re-examined. J. Funct. Program. 19, 2
 (2009), 173–190. https://doi.org/10.1017/S0956796808007090
- 1269 François Pottier and Yann Régis-Gianas. [n.d.]. The Menhir parser generator. http://gallium.inria.fr/~fpottier/menhir/.
- 1270
 Alban Reynaud, Gabriel Scherer, and Jeremy Yallop. 2021. A practical mode system for recursive definitions. Proc. ACM Program. Lang. 5, POPL (2021), 1–29. https://doi.org/10.1145/343326
- Laith Sakka, Kirshanthan Sundararajah, Ryan R. Newton, and Milind Kulkarni. 2019. Sound, fine-grained traversal fusion
 for heterogeneous trees, See [McKinley and Fisher 2019], 830–844. https://doi.org/10.1145/3314221.3314626
- 1273 1274

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

//doi.org/10.17487/RFC4180

1276	//doi.org/10.1/48//KFC4180
1277	Amir Shaikhha, Mohammad Dashti, and Christoph Koch. 2018. Push versus pull-based loop fusion in query engines. <i>J.</i>
1278	Funct. Program. 28 (2018), e10. https://doi.org/10.1017/S0956796818000102
1279	Walid Tana. 1999. Multi-Stage Programming: Its Theory and Applications. Technical Report. Mark van den Brand Jeroen Scheerder Jurgen I. Vinju, and Felco Visser. 2002. Disambiguation Filters for Scannerless.
1280	Generalized LR Parsers. In Compiler Construction. 11th International Conference. CC 2002. Held as Part of the Foint European
1281	Conferences on Theory and Practice of Software, ETAPS 2002, Grenoble, France, April 8-12, 2002, Proceedings (Lecture Notes
1282	in Computer Science, Vol. 2304), R. Nigel Horspool (Ed.). Springer, 143–158. https://doi.org/10.1007/3-540-45937-5_12
1283	Eric R. Van Wyk and August C. Schwerdfeger. 2007. Context-aware Scanning for Parsing Extensible Languages. In Proceedings
1284	of the 6th International Conference on Generative Programming and Component Engineering (Salzburg, Austria) (GPCE '07).
1285	ACM, New York, NY, USA, 63-72. https://doi.org/10.1145/12899/1.1289983 Philin Wadler, 1985. How to Replace Failure by a List of Successes. In Proc. of a Conference on Functional Programming.
1286	Languages and Computer Architecture (Nancy, France). Springer-Verlag, Berlin, Heidelberg, 113–128.
1287	Philip Wadler. 1990. Deforestation: transforming programs to eliminate trees. Theoretical Computer Science 73, 2 (1990),
1288	231-248. https://doi.org/10.1016/0304-3975(90)90147-A
1289	Jeremy Yallop and Oleg Kiselyov. 2019. Generating Mutually Recursive Definitions. In <i>Proceedings of the 2019 ACM SIGPLAN</i>
1200	Workshop on Partial Evaluation and Program Manipulation (Cascais, Portugal) (PEPM 2019). ACM, New York, NY, USA, 75–81. https://doi.org/10.1145/3204022.3204078
1201	7 5-61. https://doi.org/10.1145/5254052.5274076
1202	
1292	
1204	
1205	
1296	
1297	
1298	
1299	
1300	
1301	
1302	
1303	
1304	
1305	
1306	
1307	
1308	
1309	
1310	
1311	
1312	
1313	
1314	
1315	
1316	
1317	
1318	
1319	
1320	
1321	
1322	

Yakov Shafranovich. 2005. Common Format and MIME Type for Comma-Separated Values (CSV) Files. RFC 4180. https://

1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347 1348 1349 1350 1351 1352 1353 1354 1355
A COMPLETE DERIVATION
This section presents the complete derivation for normalizing \vec{r}
$g = \mu \operatorname{sexp} .(\operatorname{lpar} \cdot (\mu \operatorname{sexps} \cdot \epsilon \vee \operatorname{sexps}) \cdot \operatorname{rpar}) \vee \operatorname{atom}$
We automatically remove unreachable productions in the result.
$M[\operatorname{sexp}] = n_2 \Longrightarrow \{ n_2 \to \operatorname{sexp} \} \qquad M[\operatorname{sexps}] = n_3 \Longrightarrow \{ n_3 \to \operatorname{sexps} \}$
$\mathcal{N}[\![\mathfrak{e}]\!] = n_1 \Rightarrow \{ n_1 \to \epsilon \} \qquad \qquad \mathcal{N}[\![\operatorname{sexp} \cdot \operatorname{sexps}]\!] = n_4 \Rightarrow \{ n_4 \to \operatorname{sexp} n_3, n_3 \to \operatorname{sexps} \} \qquad $
$N[\![\epsilon \lor \text{ sexp } \cdot \text{ sexps }]\!] = n_5 \Rightarrow \{ n_5 \to \epsilon, n_5 \to \text{ sexp } n_3, n_3 \to \text{ sexps } \}$
$N[\mu \text{ sexps } \cdot \epsilon \lor \text{ sexps } + \text{ sexps } \Rightarrow \{\text{ sexps } \rightarrow \epsilon, \text{ sexps } \rightarrow \epsilon, n_3 \rightarrow \epsilon, n_$
$\mathcal{N}[[ext{LPAR}]] = n_6 \Rightarrow \{ n_6 o ext{LPAR} \}$
$\overline{\mathcal{M}[\text{LPAR} \cdot (\mu \text{ sexps } \cdot \varepsilon \lor \text{ sexp} \cdot \text{ sexps })]} = n_7 \Rightarrow \{n_7 \to \text{LPAR sexps, sexps} \to \epsilon, \text{sexp } n_3, n_3 \to \epsilon, n_3 \to \epsilon, n_3 \to \epsilon \text{ sexp } n_3 \}$
$N[[\mathbf{RPAR}]] = rpar \Rightarrow \{ rpar \rightarrow RPAR \}$
$N[\text{LPAR} \cdot (\mu \text{ sexps. } \epsilon \vee \text{ sexps}) \cdot \text{RPAR}] = n_8 \Rightarrow \{ n_8 \rightarrow \text{LPAR sexps rpar, sexps} \rightarrow \epsilon, \text{sexps} \rightarrow \text{sexp} n_3, n_3 \rightarrow \epsilon, n_3 \rightarrow \text{sexp} n_3, \text{rpar} \rightarrow \text{RPAR} \} $ $N[\text{ATOM}] = n_9 \Rightarrow \{ n_9 \rightarrow \text{ATOM} \}$
$\mathcal{N}[(\texttt{LPAR} \cdot (\mu \texttt{ sexps} \cdot \epsilon \lor \texttt{ sexps}) \cdot \texttt{RPAR}) \lor \texttt{ATOM}] = n_{10} \Rightarrow \{n_{10} \rightarrow \texttt{LPAR} \texttt{ sexps} \texttt{rpar}, n_{10} \rightarrow \texttt{ATOM}, \texttt{sexps} \rightarrow \epsilon, \texttt{sexps} \rightarrow \epsilon, \texttt{sexps} \rightarrow \epsilon, n_3 \rightarrow \epsilon, n_3 \rightarrow \texttt{sexp} n_3, \texttt{rpar} \rightarrow \texttt{RPAR} \}$
$N[g] = \operatorname{sexp} \Rightarrow \{\operatorname{sexp} \Rightarrow \operatorname{LPAR} \operatorname{sexps} \operatorname{rpar}, \operatorname{sexp} \Rightarrow \operatorname{ATOM}, \operatorname{sexps} \Rightarrow \operatorname{LPAR} \operatorname{sexps} \operatorname{rpar} n_3, \operatorname{sexps} \Rightarrow \operatorname{ATOM} n_3, n_3 \Rightarrow \varepsilon, n_3 \Rightarrow \operatorname{LPAR} \operatorname{sexps} \operatorname{rpar} n_3, \operatorname{rpar} \Rightarrow \operatorname{RPAR} \}$
Comparing the simplified derivation in Section 4.1 with the complete derivation, we note the following simplification: first, we omit the derivation of tokens: second. when normalizing sexus we produce a nonterminal n_3 with a production $n_3 \rightarrow$ sexus. That means n_3 is conivalent
to sexps. However, this <i>n</i> ₃ is retained in the final result, making the final grammar a big larger. It's easy to check that the grammar is equivalent
to the one given in the paper.
It is easy to consider an optimization process that gets rid of n_3 in the middle of the derivation. For example, for the result for n_5 , instead of
$n_5 \Rightarrow n_5 \to \epsilon, n_5 \to \operatorname{sexp} n_3, n_3 \to \operatorname{sexps}$
We can have
$n_5 \Rightarrow n_5 \to \epsilon, n_5 \to \operatorname{sexp} \operatorname{sexps}$
Then the normalization result would be exactly the same as the one in the paper.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:29

1358 1359 1360

1361

1362

1363 1364

1365

1366

1367

1372

1373

1391

1392

1393

1394

1395

1396

1397

1398 1399

1400

1401

1402

1404

1405 1406

B DETERMINISTIC PARSING

THEOREM 4.1 (DETERMINISTIC PARSING). If G is a DGNF grammar, then for any expansion $G \vdash n \rightsquigarrow w$, there is a unique derivation for such expansion.

PROOF. By straightforward induction on $G \vdash n \rightsquigarrow w$.

C WELL-TYPED NORMALIZATION

This section presents well-typed normalization, which shows how normalization captures the type information, and then proves its properties that are important for later proofs.

First, we note that during normalization (Figure 6), we create one fresh nonterminal exactly for one context-free expression. Therefore, we can attach to each nonterminal its type information. That is, instead of *n*, we have n_{τ} , where τ indicates the type of *n*. We also write α_{τ} where τ is the type of α as in $\mu\alpha : \tau$. *g*.

Refining the normalization, we have:

 $\mathcal{N}[\![g]\!]$ returns $n_{\tau} \Rightarrow G$, with a grammar *G* and the start nonterminal *n* of type τ (with *n* fresh)

 $\Gamma; \Delta \vdash \mathbf{n}_{\tau} : \tau$

With that, we can type-check any *N* according to the typing rules, by treating *t* as constants, n_{τ} as nonterminal of type τ , and lists as sequences (e.g., $n_{1\tau_1} n_{2\tau_2}$ as $n_{1\tau_1} \cdot n_{2\tau_2}$).

Now we can prove properties about the well-typed normalization. While those lemmas are proved in the typed normalization, they naturally hold for the untyped normalization as the two are the same process.

LEMMA C.1. Given $\Gamma; \Delta \vdash g : \tau$, and $\mathcal{N}[\![g]\!]$ returns $n_{\tau'} \Rightarrow G$, then $\tau = \tau'$.

PROOF. By a straightforward induction on Γ ; $\Delta \vdash g : \tau$.

LEMMA C.2. Given $\Gamma; \Delta \vdash g : \tau'$, and $\mathcal{N}[\![g]\!]$ returns $_ \Rightarrow G$, then for any $n_{\tau} \in G$, if $N_1, ..., N_i$ are all productions of n. we have

1403 • $\tau = \tau_1 \lor \tau_2 \lor \cdots \lor \tau_i$, where

• $(n_{\tau} \to N_1 \in G \land \Gamma; \Delta \vdash N_1 : \tau_1)$ and $(n_{\tau} \to N_2 \in G \land \Gamma; \Delta \vdash N_2 : \tau_2)$ and \cdots and $(n_{\tau} \to N_i \in G \land \Gamma; \Delta \vdash N_i : \tau_i);$ and

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

• $\tau_1 \# \tau_2 \cdots \# \tau_i$, *i.e.*, all $\tau_1, \tau_2, \cdots, \tau_i$ are apart from each other. 1407 1408 **PROOF.** By induction on Γ ; $\Delta \vdash q : \tau$. 1409 • The cases for ϵ , t, \perp and α follow trivially. 1410 • The case for $q_1 \cdot q_2$. 1411 $\mathcal{N}[\![g_1 \cdot g_2]\!] = n_{\tau_1 \cdot \tau_2} \Longrightarrow \{n_{\tau_1 \cdot \tau_2} \to N_1 \ n_{2\tau_2} \mid n_{1\tau_1} \to N_1 \in G_1\} \cup G_1 \cup G_2$ 1412 $\mathcal{N}\llbracket g_1 \rrbracket = n_{1\tau_1} \Longrightarrow G_1 \land \mathcal{N}\llbracket g_2 \rrbracket = n_{2\tau_2} \Longrightarrow G_2$ 1413 By I.H., we know that for each N_1 , we have Γ ; $\Delta \vdash N_1 : \tau'_1$ for some τ'_1 , and τ_1 is the \lor of all 1414 τ'_1 , and all τ'_1 is apart (#) from each other. 1415 According to well-typedness, we know that $\tau_1 \otimes \tau_2$, which says that τ_1 .FLAST $\cap \tau_2$.FIRST = \emptyset , 1416 and $\neg \tau_1$.NULL. 1417 Since τ_1 is the \lor of all τ'_1 , we know τ'_1 .NULL = false, and τ'_1 .FLAST $\cap \tau_2$.FIRST = \emptyset , and thus 1418 $\tau_1' \circledast \tau_2.$ 1419 So Γ ; $\Delta \vdash N_1 n_{2\tau_2} : \tau'_1 \cdot \tau_2$ 1420 Moreover, since τ'_1 .Null = false, we have $\tau'_1 \cdot \tau_2$.FIRST = τ'_1 .FIRST and $\tau'_1 \cdot \tau_2$.Null = false. 1421 Given that all τ'_1 apart from each other, we can derive that all $\tau'_1 \cdot \tau_2$ apart from each other. 1422 • The case for 1423 $\mathcal{N}\llbracket g_1 \lor g_2 \rrbracket = n_{\tau_1 \lor \tau_2} \Longrightarrow \{n_{\tau_1 \lor \tau_2} \to N_1 \mid n_{1\tau_1} \to N_1 \in G_1\} \cup \{n_{\tau_1 \lor \tau_2} \to N_2 \mid n_{2\tau_2} \to N_2 \in G_2\}$ 1424 \cup $G_1 \cup G_2$ 1425 $\mathcal{N}\llbracket g_1 \rrbracket = n_{1\tau_1} \Longrightarrow G_1 \land \mathcal{N}\llbracket g_2 \rrbracket = n_{2\tau_2} \Longrightarrow G_2$ 1426 By I.H., we know that for each N_1 , we have Γ ; $\Delta \vdash N_1 : \tau'_1$ for some τ'_1 , and τ_1 is the \lor of all 1427 τ'_1 , and all τ'_1 is apart (#) from each other. Moreover, for each N_2 , we have $\Gamma; \Delta \vdash N_2 : \tau'_2$ for 1428 some τ'_2 , and τ_2 is the \lor of all τ'_2 , and all τ'_2 is apart (#) from each other. 1429 It's easy to see that $\tau_1 \lor \tau_2$ is the \lor of all τ'_1 and τ'_2 . 1430 According to well-typedness, we know that $\tau_1 \# \tau_2$. That is τ_1 .FIRST $\cap \tau_2$.FIRST = \emptyset , and 1431 $\neg(\tau_1.\text{Null} \land \tau_2.\text{Null})$. From the former, we can derive that $\tau'_1.\text{FIRST} \cap \tau'_2.\text{FIRST} = \emptyset$. From 1432 the latter, we know that at least one of τ_1 and τ_2 has NULL = false, so at least one of τ'_1 and 1433 τ'_2 has NULL = false. With that, we have $\tau'_1 \# \tau'_2$. Thus, all τ'_1 and τ'_2 apart from each other. 1434 • The case for $\mu\alpha$: τ . q 1435 $\mathcal{N}\llbracket \mu\alpha: \tau, g \rrbracket = \alpha_{\tau} \to \{\alpha_{\tau} \to N \mid n_{\tau} \to N \in G\} \cup \{n'_{\tau'} \to N \overline{n}' \mid n'_{\tau'} \to \alpha_{\tau} \overline{n}' \in G \land n_{\tau} \to N \in G\}$ 1436 $\cup G \setminus_{n' \prec \to \alpha_{\tau}} \overline{n'}$ 1437 $\mathcal{N}[\![q]\!] = \mathbf{n}_{\tau} \Rightarrow \mathbf{G}$ 1438 $G \setminus_{n' \neq \prime} \rightarrow \alpha_{\tau} \overline{n'}$ is G with all $n'_{\tau'} \rightarrow \alpha_{\tau} \overline{n'}$ removed for any n', τ' and $\overline{n'}$ 1439 By I.H., we know that for each *N*, we have Γ ; $\Delta \vdash N : \tau''$ for some τ'' , and τ is the \lor of all 1440 $\tau^{\prime\prime}$, and all $\tau^{\prime\prime}$ is apart (#) from each other. 1441 The goal for α_{τ} follows from n_{τ} . The remaining is to show that the goal holds for each $n'_{\tau'}$ 1442 that has a production that starts with α_{τ} . Essentially what happens is that one production 1443 $n'_{\tau'} \to \alpha_{\tau} \,\overline{n}'$ is replaced by multiple productions $n'_{\tau'} \to N \,\overline{n}'$ for each $n_{\tau} \to N \in G$ where 1444 $\Gamma; \Delta \vdash \mathbf{N} : \tau''.$ 1445 First, we need to show that $N \overline{n}'$ is well-typed. We already know each individual terminal or 1446 nonterminal in $N \overline{n}'$ is well-typed, so the only requirement is the \circledast condition during type-1447 checking. Given that $\alpha_{\tau} \overline{n}'$ is well-typed, we know that τ .NULL = false, so τ'' .NULL = false. 1448 Moreover, τ'' .FLAST $\subseteq \tau$.FLAST. With that, and the fact that $\alpha_{\tau} \overline{n}'$ is well-typed, we can 1449 derive that the \circledast condition is always satisfied when type-checking N \overline{n}' . Therefore, N \overline{n}' is 1450 well-typed. 1451 Because τ is the \vee of all τ'' , it's easy to show that the type of $\alpha_{\tau} \overline{n}'$ is the \vee of the types of 1452 all $N \overline{n}'$. Therefore, the type of n' is the same as before. Also, all types of the productions of 1453 *n*' are still apart with each other. 1454 1455

1456 1457 D NORMALIZATION IS WELL-DEFINED (PROOF) 1458 LEMMA 4.2 (PRODUCTIONS OF NULL). Given Γ ; $\Delta \vdash q : \tau$ and $\mathcal{N}[\![q]\!]$ returns $n \Rightarrow G, \tau.NULL =$ true 1459 if and only if (1) $n \to \epsilon \in G$; or (2) $n \to \alpha \in G$ where $(\alpha : \tau') \in \Gamma$ and τ' .NULL = true. 1460 1461 **PROOF.** Left to right According to Lemma C.2, we must have one $n_{\tau} \to N \in G$, where $\Gamma; \Delta \vdash$ 1462 N : τ , and τ .NULL = true. We case analyze the shape of N: 1463 • If $N = \epsilon$, then we have proved (1). 1464 • If $N = t \overline{n}$, then it's impossible that τ .NULL = true. 1465 • If $N = \alpha \overline{n}$. Since $\alpha \overline{n}$ is well-typed, if \overline{n} is not empty, then the type must have NULL = false. 1466 Therefore \overline{n} must be empty, and α has its type NULL = true. So we have proved (2). 1467 **Right to left** Following Lemma C.2, the type τ is the \lor of all types. If either $n \to \epsilon$ of α has type 1468 NULL = true, we know that τ .NULL = true. 1469 1470 1471 THEOREM 4.3 (WELL-DEFINEDNESS). If $\Gamma; \Delta \vdash q : \tau$, then $\mathcal{N}[\![q]\!]$ returns $n \Rightarrow G$ for some G and n. 1472 **PROOF.** By induction on *q*. Most cases are straightforward. The only interesting cases are when 1473 $q = q_1 \cdot q_2$ or $q = \mu \alpha$. q'. 1474 • $g = g_1 \cdot g_2$. We have: 1475 $\mathcal{N}[\![q_1 \cdot q_2]\!] = n \Longrightarrow \{n \to N_1 \ n_2 \mid n_1 \to N_1 \in G_1\} \cup G_1 \cup G_2$ 1476 $\mathcal{N}\llbracket q_1 \rrbracket = n_1 \Longrightarrow G_1 \land \mathcal{N}\llbracket q_2 \rrbracket = n_2 \Longrightarrow G_2$ 1477 As $g_1 \cdot g_2$ is well-typed, we know that the type of g_1 has NULL = false. By Lemma 4.2, N_1 is 1478 not ϵ , ensuring that $N_1 n_2$ is a valid form. 1479 • $g = \mu \alpha$. g'. We have: 1480 $\Gamma; \Delta, \alpha : \tau \vdash q : \tau$ 1481 $\mathcal{N}\llbracket \mu \alpha, q \rrbracket = \alpha \implies \{\alpha \to N \mid n \to N \in G\} \cup \{n' \to N \ \overline{n'} \mid n' \to \alpha \ \overline{n'} \in G \land n \to N \in G\} \cup G \setminus_{n' \to \alpha \ \overline{n'}}$ 1482 $\mathcal{N}[\![q]\!] = n \Rightarrow G$ 1483 $G \setminus_{n' \to \alpha \ \overline{n}'}$ is *G* with all $n' \to \alpha \ \overline{n}'$ removed for any *n'* and \overline{n}' 1484 We need to show that $N \overline{n}'$ is valid, requiring either N to not be ϵ , or \overline{n}' to be empty. 1485 Since $\alpha \overline{n}'$ is well-typed (Lemma C.2), we know that either \overline{n}' is empty, or α must have 1486 NULL = false. In the first case we are done. In the second case, following Lemma 4.2, we 1487 know N cannot be ϵ . 1488 1489 1490 E NORMALIZATION RETURNS DGNF GRAMMARS (PROOF) 1491 E.1 Normalizing closed expressions produces no $\alpha \overline{n}$ form 1492 1493 LEMMA 4.4 (INTERNAL NORMAL FORM). Given Γ ; $\Delta \vdash q : \tau$ and $\mathcal{N}[\![q]\!]$ returns $n \Rightarrow G$, 1494 • $if(n \to \alpha \overline{n}) \in G$, then we have $\alpha \in dom(\Gamma)$; 1495 • if $(n' \to \alpha \overline{n}) \in G$ for any n', then we have $\alpha \in fv(q)$, and thus $\alpha \in dom(\Gamma, \Delta)$. 1496 **PROOF. Part 1** By induction on Γ ; $\Delta \vdash q : \tau$, most cases are straightforward. We discuss the 1497 following three cases: 1498 • $q = \alpha$. As q is well-typed, it must be $\alpha \in \text{dom}(\Gamma)$. The goal follows directly. 1499 • $q = q_1 \cdot q_2$. The goal follows by the I.H. on q_1 . 1500 • $g = \mu \alpha$. g'. As the well-typedness of g' adds α to Δ , the goal follows directly by the I.H. on 1501 q'. 1502 **Part 2** By induction on Γ ; $\Delta \vdash q : \tau$. The only interesting case is when $q = \mu \alpha$. q'. We have 1503 1504

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

 $\Gamma; \Delta, \alpha : \tau \vdash q : \tau$ 1505 $\mathcal{N}[\![\mu\alpha. g]\!] = \alpha \Longrightarrow \{\alpha \to N \mid n \to N \in G\} \cup \{n' \to N \ \overline{n}' \mid n' \to \alpha \ \overline{n}' \in G \land n \to N \in G\} \cup G \backslash_{n' \to \alpha \ \overline{n}'}$ 1506 1507 $\mathcal{N}[\![q]\!] = n \Rightarrow G$ 1508 $G \setminus_{n' \to \alpha \ \overline{n}'}$ is *G* with all $n' \to \alpha \ \overline{n}'$ removed for any n' and \overline{n}' 1509 By I.H., we know that for all $(n'' \to beta \ \overline{n}) \in G \setminus_{n' \to \alpha} \overline{n'}, \beta \in fv(q)$. 1510 For *N*, if it is *beta* \overline{n} , then either $\beta \in \text{fv}(\mu \alpha, q)$, or $\beta = \alpha$. By Part 1, we know that $\beta \in \text{dom}(\Gamma)$, 1511 so $\beta \neq \alpha$. So it can only be $\beta \in fv(\mu\alpha, q)$. And the goal follows. 1512 1513 1514 COROLLARY 4.5 (NORMAL FORM). Given $\bullet; \bullet \vdash g : \tau$ and $N[\![g]\!]$ returns $_ \Rightarrow G$, then for all 1515 $(n \to N) \in G, N \text{ is } \epsilon \text{ or } t \overline{n} \text{ for some } t \text{ and } \overline{n}.$ 1516 1517 PROOF. Follows directly from Lemma 4.4. 1518 1519 1520 A nonterminal's non- ϵ productions start with distinct terminals. E.2 1521 LEMMA 4.6 (TERMINALS IN FIRST). Given Γ ; $\Delta \vdash g : \tau$ and $\mathcal{N}[\![g]\!]$ returns $n \Rightarrow G$, we have $t \in \tau$. First 1522 if and only if (1) $(n \to t \overline{n}) \in G$; or (2) $(n \to \alpha \overline{n}) \in G$ where $(\alpha : \tau') \in \Gamma$ and $t \in \tau$. First. 1523 1524 **PROOF.** Left to right According to Lemma C.2, we must have one $n_{\tau} \rightarrow N \in G$, where $\Gamma; \Delta \vdash$ 1525 $N : \tau$, and $t \in \tau$. FIRST. We case analyze the shape of N: 1526 1527 • If $N = \epsilon$, then it's impossible. 1528 • If $N = t \overline{n}$, then we have proved (1). 1529 • If $N = \alpha \overline{n}$. Since $\alpha \overline{n}$ is well-typed, the FIRST of the type of $\alpha \overline{n}$ is equivalent to the FIRST of 1530 the type of α . So we have proved (2). 1531 1532 **Right to left** Following Lemma C.2, the type τ is the \lor of all types. If either $n \to t \overline{n}$ of α has 1533 type $t \in$ FIRST, we know that $t \in \tau$.FIRST. 1534 1535 1536 LEMMA E.1 (PRODUCTIONS WITH DISTINCT TERMINALS). If Γ ; $\Delta \vdash q : \tau$, and $\mathcal{N}[\![q]\!]$ returns $_ \Rightarrow G$, 1537 then for any two productions $(n \to t_1 \overline{n}_1) \in G$ and $(n \to t_2 \overline{n}_2) \in G$, we have $t_1 \neq t_2$. 1538 1539 1540 **PROOF.** Suppose there are $n \to t \overline{n}_1$ and $n \to t \overline{n}_2$. 1541 By Lemma C.2, we know that the types of $t \overline{n}_1$ and $t \overline{n}_2$ must be apart. Therefore they have 1542 disjoint FIRST. 1543 By Lemma 4.6, we know that both $t \overline{n}_1$ and $t \overline{n}_2$ have $t \in$ FIRST. However, since their types have 1544 disjoint FIRST, this is impossible. So contradiction. 1545 1546 1547 The ϵ -production may only be used when other productions do not apply. E.3 1548

We defined the notion of containment of types as follows. The key of the definition is rule ST-BASE, which says that a grammar g_1 is a subtype grammar of g_2 , if g_1 is of type τ_1 , g_2 is of type τ_2 , and $\tau_1 = \tau_2 \lor \tau$ for some τ . Notably, we have $\Gamma; \Delta \vdash g \leq g$ for any well-typed grammar $\Gamma; \Delta \vdash g : \tau$, as we have $\tau = \tau \lor \{\text{NULL} = \text{false}; \text{FIRST} = \emptyset\}$.

Anon.

 $\Gamma; \Delta \vdash q_1 \leq q_2$ (containment of types) ST-BASE ST-TRANS $\frac{\Gamma; \Delta \vdash g_1 : \tau_1 \qquad \Gamma; \Delta \vdash g_2 : \tau_2 \qquad \tau_1 = \tau_2 \lor \tau}{\Gamma; \Delta \vdash g_1 \lesssim g_2} \qquad \qquad \frac{\Gamma; \Delta \vdash g_1 \lesssim g_2 \qquad \Gamma; \Delta \vdash g_2 \lesssim g_3}{\Gamma; \Delta \vdash g_1 \lesssim g_3}$ $\frac{\prod_{i=1}^{\text{ST-CON}} \Gamma_{i} \otimes g_{1} \otimes g_{1}' \qquad \Gamma_{i} \otimes f_{2} \otimes g_{2}'}{\prod_{i=1}^{\text{ST-CON}} \Gamma_{i} \otimes f_{2} \otimes g_{1}' \otimes g_{2}'} \qquad \qquad \frac{\prod_{i=1}^{\text{ST-UNION}} \Gamma_{i} \otimes f_{1} \otimes g_{1}' \qquad \Gamma_{i} \otimes f_{2} \otimes g_{2}'}{\prod_{i=1}^{\text{ST-UNION}} \Gamma_{i} \otimes f_{2} \otimes g_{1}' \otimes g_{2}'}$ LEMMA E.2. If Γ ; $\Delta \vdash q_1 : \tau_1$, and Γ ; $\Delta \vdash q_1 \leq q_2$, then Γ ; $\Delta \vdash q_2 : \tau_2$, and $\tau_1 = \tau_2 \lor \tau$ for some τ . **PROOF.** By induction on Γ ; $\Delta \vdash q_1 \leq q_2$. rule st-base follows trivially. • Case st-trans $\frac{\Gamma; \Delta \vdash g_1 \lesssim g_2 \qquad \Gamma; \Delta \vdash g_2 \lesssim g_3}{\Gamma; \Delta \vdash g_1 \lesssim g_3}$ We have q_1 of type τ_1 . By I.H., q_2 of type τ_2 , and $\tau_1 = \tau_2 \vee \tau$. By the second I.H., q_3 of type τ_3 , and $\tau_2 = \tau_3 \vee \tau'$. Therefore, $\tau_1 = \tau_3 \lor (\tau \lor \tau')$. • Case $\frac{\Gamma; \Delta \vdash g_1 \leq g'_1 \qquad \Gamma; \Delta \vdash g_2 \leq g'_2}{\Gamma; \Delta \vdash q_1 \cdot q_2 \leq g'_1 \cdot q'_2}$ I; $\Delta \vdash g_1 \cdot g_2 \lesssim g_1 \cdot g_2$ We have $g_1 \cdot g_2$ of type $\tau_1 \cdot \tau_2$ with g_1 of type τ_1 and g_2 of type τ_2 and $\tau_1 \circledast \tau_2$. By I.H., g'_1 of type τ'_1 , and $\tau_1 = \tau'_1 \lor \tau$. By the second I.H., g'_2 of type τ'_2 , and $\tau_2 = \tau'_2 \lor \tau'$. Now we want to show $g'_1 \cdot g'_2$ is of type $\tau'_1 \cdot \tau'_2$. For that, we need to prove $\tau'_1 \circledast \tau'_2$. That means we need to prove τ'_1 .FLAST $\cap \tau'_2$.FIRST = $\emptyset \land \neg \tau'_1$.NULL We already know $\tau_1 \circledast \tau_2$, which means τ_1 .FLAST $\cap \tau_2$.FIRST = $\emptyset \land \neg \tau_1$.NULL Since $\tau_1 = \tau'_1 \lor \tau$ and $\tau_2 = \tau'_2 \lor \tau'$, we can derive τ'_1 .FLAST $\cap \tau'_2$.FIRST $= \emptyset \land \neg \tau'_1$.NULL Therefore, $\tau'_1 \circledast \tau'_2$, and $g'_1 \cdot g'_2$ is of type $\tau'_1 \cdot \tau'_2$. Now the goal is to relate $\tau_1 \cdot \tau_2$ with $\tau'_1 \cdot \tau'_2$. $\tau_1 \cdot \tau_2 = \begin{cases} \text{Null} = \tau_1.\text{Null} \land \tau_2.\text{Null} \\ \text{First} = \tau_1.\text{First} \cup \tau_1.\text{Null}? \tau_2.\text{First} \\ \text{FLast} = \tau_2.\text{FLast} \cup \tau_2.\text{Null}? (\tau_2.\text{First} \cup \tau_1.\text{FLast}) \end{cases}$ given $\neg \tau_1$.NULL $\tau_1 \cdot \tau_2 = \begin{cases} \text{NULL} = \text{false} \\ \text{FIRST} = \tau_1.\text{FIRST} \\ \text{FLAST} = \tau_2.\text{FLAST} \cup \tau_2.\text{NULL}?(\tau_2.\text{FIRST} \cup \tau_1.\text{FLAST}) \end{cases}$ Similarly, $\tau_1' \cdot \tau_2' = \begin{cases} \text{NULL} = \text{false} \\ \text{FIRST} = \tau_1'.\text{FIRST} \\ \text{FLAST} = \tau_2'.\text{FLAST} \cup \tau_2'.\text{NULL}?(\tau_2'.\text{FIRST} \cup \tau_1'.\text{FLAST}) \end{cases}$ We have $\tau_1 = \tau_1' \vee \tau$ and $\tau_2 = \tau_2' \vee \tau'$. Therefore, with $\neg \tau_2.\text{NULL}$ implying $\neg \tau_2'.\text{NULL}$, $\tau_1 \cdot \tau_2 = (\tau_1' \cdot \tau_2') \vee \begin{cases} \text{NULL} = \text{false} \\ \text{FIRST} = \tau.\text{FIRST} \\ \text{FLAST} = \tau.\text{FIRST} \end{cases}$ FLAST $= \tau.\text{FIRST}$ given $\neg \tau_1$.NULL

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

1:34

Fusing Lexing and Parsing

1603	• Case	
1604	ST-UNION $\Gamma \cdot \Lambda \vdash a_1 \leq a' \qquad \Gamma \cdot \Lambda \vdash a_2 \leq a'$	
1605	$\frac{1, \Delta + g_1 \gtrsim g_1}{-}$	
1606	$\Gamma; \Delta \vdash g_1 \lor g_2 \leq g'_1 \lor g'_2$	1
1607	We have $g_1 \lor g_2$ of type $\tau_1 \lor \tau_2$ with g_1 of type τ_1 and g_2 of type τ_2	and $\tau_1 \# \tau_2$.
1608	By I.H., g'_1 of type τ'_1 , and $\tau_1 = \tau'_1 \lor \tau$.	
1609	By the second I.H., g'_2 of type τ'_2 , and $\tau_2 = \tau'_2 \lor \tau'$.	
1610	Now we want to show $g'_1 \lor g'_2$ is of type $\tau'_1 \lor \tau'_2$. For that, we need	to prove $\tau'_1 \# \tau'_2$.
1611	That means we need to prove $(\tau'_1.FIRST \cap \tau'_2.FIRST = \emptyset) \land \neg(\tau'_1.NU$	LL $\wedge \tau'_2$.NULL)
1612	We already know $\tau_1 \# \tau_2$, which means $(\tau_1.FIRST \cap \tau_2.FIRST = \emptyset) \land$	$\neg(\tau_1.\text{NULL} \land \tau_2.\text{NULL})$
1613	Since $\tau_1 = \tau_1' \lor \tau$ and $\tau_2 = \tau_2' \lor \tau'$, we can derive $(\tau_1'.FIRST \cap \tau_2'.FII$	$RST = \emptyset) \land \neg(\tau_1).NULL \land$
1614	τ_2 .NULL)	
1615	Finally we have $\tau_1 \# \tau_2$, and $g_1 \lor g_2$ is of type $\tau_1 \lor \tau_2$.	
1616	Finally, we have $\tau_1 \lor \tau_2 = (\tau_1 \lor \tau_2) \lor (\tau \lor \tau')$.	
1617		
1618	Lemma E.3 (Expansion preserves typing). Given Γ ; $\Delta \vdash q : \tau, \mathcal{N}[\![q]\!]$	returns \Rightarrow G, if G +
1619	$n_{\tau} \rightsquigarrow \overline{t} n' \overline{n}$, then $\Gamma; \Delta \vdash \overline{t} n' \overline{n} : \tau_1$, and $\tau = \tau_1 \lor \tau'$ for some τ' .	_ , ,
1620		
1621	PROOF. By induction on $G \vdash n_{\tau} \rightsquigarrow t n' n$.	
1622	• In the base case, $G \vdash n_{\tau} \rightsquigarrow n_{\tau}$. The goal follows trivially.	
1623	• In the inductive case, we have $G \vdash n_{\tau} \rightsquigarrow t n' \overline{n}, n' \rightarrow N \in G$ and s	so $G \vdash n \rightsquigarrow t N \overline{n}$,
1624	By I.H., we have Γ ; $\Delta \vdash t n' \overline{n} : \tau_1$, and $\tau = \tau_1 \lor \tau'$.	
1625	According to Lemma C.2, we know that $\Gamma; \Delta \vdash n' \leq N$ by rule ST-E	BASE.
1626	Therefore, Γ ; $\Delta \vdash tn' n \leq tN n$ by rule st-con.	
1627	By Lemma E.2, 1; $\Delta \vdash t N n : \tau_2$, and $\tau_1 = \tau_2 \lor \tau''$.	
1628	Therefore, $\tau = \tau_2 \vee (\tau' \vee \tau'')$.	
1629		
1630	LEMMA E.4 (GUARDED ϵ -production). Given $\Gamma; \Delta \vdash q : \tau, \mathcal{N}[\![q]\!]$ return	$s n \Rightarrow G$, and $G \vdash n \rightsquigarrow^*$
1631	$\cdots n_1 n_2 \cdots$, if $(n_1 \to \epsilon) \in G$, then either $(n_1 \to t \overline{n_1}) \notin G$ or $(n_2 \to t \overline{n_2})$	$\notin G$ for any $t, \overline{n}_1, \overline{n}_2$.
1622		
1634	PROOF. We have: $\dots n_n n_n \dots n_n$ well-typed	By Lemma F 3
1635	The type of $\cdots n_1$ is $\tau \cdot \tau_1$ the type of n_1 is τ_2 and the type of n_2 is τ_2	Suppose
1636	$\tau \cdot \tau_1 \otimes \tau_2$	By typing
1637	$\tau \cdot \tau_1 \in \tau_2$ $\tau \cdot \tau_1$ FLAST $\cap \tau_2$ FIRST = Ø	By ®
1638	$\tau \cdot \tau_1$. FLAST = τ_1 . FLAST $\cup \tau_1$. NULL ? (τ_1 . FIRST $\cup \tau$. FLAST)	By definition
1639	$n_1 \to \epsilon \in G$	Given
1640	n_1 .Null = true	Lemma 4.2
1641	$\tau \cdot \tau_1$.FLAST = τ_1 .FLAST $\cup (\tau_1$.FIRST $\cup \tau$.FLAST)	Follows
1642	τ_1 .First \cap τ_2 .First = \emptyset	Follows
1643	$n_1 \to t \ \overline{n}_1 \in G \land n_2 \to t \ \overline{n}_2 \in G$	Assume
1644	$t \in \tau_1$.First $\land t \in \tau_2$.First	Lemma 4.6
1645		I
1646		
1647		
1648		
1649		
1650		
1651		

Contradiction with τ_1 .FIRST $\cap \tau_2$.FIRST = \emptyset

1655 E.4 Final result

THEOREM 4.7 ($\mathcal{N}[\![g]\!]$ produces DGNF). If $\bullet; \bullet \vdash g : \tau$ and $\mathcal{N}[\![g]\!]$ returns $_ \Rightarrow G$, then G is DGNF.

PROOF. Follows from Corollary 4.5, Lemma E.1, and Lemma E.4.

F SOUNDNESS (PROOF)

F.1 An alternative normalization

To make proofs easier, we consider the definition N, which has the same definition as N except for the case of $\mu\alpha$. g, where we do not substitute α :

$$\boldsymbol{N}(\mu\alpha:\tau,g) = \alpha \Rightarrow \{\alpha \to N \mid n \to N \in G\} \cup G$$

where $\boldsymbol{N}(g) = n \Rightarrow G$

While N does not return a DGNF grammar, it is easy to see that N and N defines the same language:

LEMMA F.1. If $\mathcal{N}[\![g]\!]$ return $n_1 \Rightarrow G_1$, and $\mathcal{N}(g)$ return $n_2 \Rightarrow G_2$, then for all $w, G_1 \vdash n_1 \rightsquigarrow^* w$ if and only if $G_2 \vdash n_2 \rightsquigarrow^* w$.

PROOF. By straightforward induction on g.

F.2 Subexpression

1677 The subexpression relation essentially defines a subset relation between the grammars denoted by1678 context-free expressions.

$$\begin{array}{c} \hline g_1 \sqsubseteq g_2 \\ \hline g \sqsubseteq g \\ \hline g \sqsubseteq g \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-TRANS} \\ g_1 \sqsubseteq g_2 \\ g_2 \sqsubseteq g_3 \\ \hline g_1 \sqsubseteq g_3 \\ \hline g_1 \sqsubseteq g_1 \cdot g_2 \\ \hline g_1 \sqsubseteq g_1 \cdot g_2 \\ \hline g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g \sqsubseteq \mu \alpha \cdot g \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g \sqsubseteq \mu \alpha \cdot g \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \sqsubseteq g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \qquad \begin{array}{c} \begin{array}{c} \text{SG-UNION-L} \\ g_2 \vdash g_1 \vee g_2 \\ \hline \end{array} \end{array}$$

We can show that what subexpression means in terms of the alternative normalization.

LEMMA F.2. If $g_1 \subseteq g_2$, and $N(g_1)$ returns $n_1 \Rightarrow G_1$, and $N(g_2)$ returns $n_2 \Rightarrow G_2$, then for all $n \in \text{dom}(G_1), (n \rightarrow N) \in G_1$ if and only if $(n \rightarrow N) \in G_2$.

PROOF. By straightforward induction on $g_1 \sqsubseteq g_2$.

F.3 Proof of soundness

In the following lemma statement, we denote a natural number as \square , and the length of a word w as |w|. The relations $\gamma \models \Gamma$ and $\delta \models \Delta$ mean that γ and δ give interpretations (i.e., languages L) of variables in Γ and Δ respectively.

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

Anon.

Fusing Lexing and Parsing

1701					
1702		$\delta \models \Delta$ $L \models \tau$			
1703		$\overline{\delta \mid / \alpha \models \Lambda \ \alpha \cdot \tau}$			
1704		0, L/ u P <u>A</u> , u · <i>i</i>			
1705 1706	$L \models \tau \triangleq Null(LL) \Rightarrow \tau.Null \land Firs$	$\tau(L) \subseteq \tau.\mathrm{First} \land \mathrm{FLas}$	$\tau(L) \subseteq \tau.FLast$		
1707	LEMMA F.3. Given Γ ; $\Delta \vdash q : \tau$, and $\gamma \models \Gamma$, and	$\delta \models \Delta$, and $\mathbf{N}(q)$ return	$n \approx n \Rightarrow G, if$		
1708	(1) $a \sqsubset a'$, where $\bullet : \bullet \vdash a' : \tau'$ and $N(a')$ retu	1) $a \sqsubset a'$, where $\bullet: \bullet \vdash a': \tau'$ and $N(a')$ returns $n' \Rightarrow G'$: and			
1709	(2) $\forall \alpha \in \text{dom}(\gamma), \forall w_1 \leq n, w_1 \in \gamma(\alpha) \text{ if and only if } G' \vdash \alpha \rightsquigarrow^* w_1; and$				
1710	(3) $\forall \alpha \in \text{dom}(\delta), \forall w_2 < n, w_2 \in \delta(\alpha) \text{ if and only if } G' \vdash \alpha \rightsquigarrow^* w_2,$				
1711	$n \forall w \leq \mathbb{n}, w \in \llbracket g \rrbracket_{(\gamma, \delta)} \text{ if and only if and } G' \vdash n \rightsquigarrow^* w.$				
1712	PROOF. By first induction on \mathbb{D} . The base case	of 0 is trivial. In the i	nductive case, we have that		
1714	the lemma holds for $ w < n$, and we want to pr	ove it for $ w < n$.			
1715	Now we perform induction on <i>q</i> .				
1716	• The cases for $a = t$ $a = \epsilon$ and $a = 1$ are	straightforward			
1717	• $a = \alpha$. Then $\mathcal{N}(\alpha) = n \Rightarrow n \rightarrow \alpha$. By Let	$ma E.2.$ we know $(n \cdot 1)$	$\rightarrow \alpha$) $\in G'$ and there is no		
1718	other production for n in G' .		.,		
1719	Since q is well-typed, it must be $\alpha \in \text{dom}(\Gamma)$, and thus $\alpha \in \text{dom}(\gamma)$. Then $[\![q]\!]_{(\gamma,\delta)} = \gamma(\alpha)$.				
1720	As given, we know that $\forall w \leq n, w \in Y$	(α) if and only if G' +	$\alpha \sim^* w.$		
1721	Since we know $(n \rightarrow \alpha) \in G'$, we have	$\forall w \leq \mathbb{n}, w \in \gamma(\alpha)$ if	and only if $G' \vdash n \rightsquigarrow^* w$.		
1722	• $g = g_1 \lor g_2$. Then $[g_1 \lor g_2]_{(\chi, \delta)} = [g_1]_{(\chi, \delta)}$	$\cup [g_2]_{(\gamma,\delta)}$	-		
1723	We have $U = U = U(p,0)$) 20 2(1,0)			
1724	$\{n \to N_1 \mid n_1 \to N_1 \in G_1\} \cup \{n \to N_2$	$\mid n_2 \to N_2 \in G_2 \} \cup G_2$	$\cup G_2$		
1725	$\boldsymbol{N}(g_1) = \boldsymbol{n}_1 \Longrightarrow \boldsymbol{G}_1$				
1726	$\mathcal{N}(g_2) = n_2 \Longrightarrow G_2$				
1727	The goal follows from I.H. on g_1 and g_2 .				
1728	• $g = g_1 \cdot g_2$. Then $[[g_1 \cdot g_2]]_{(\gamma,\delta)} = \{w_1 \cdot w_2\}$	$ w_1 \in [g_1]_{(\gamma,\delta)} \wedge w_2$	$\in [[g_2]]_{(\gamma,\delta)}\}.$		
1729	According to \mathbf{N} , we have				
1730	$\{n \to N_1 \ n_2 \mid n_1 \to N_1 \in G_1\} \cup G_1 \cup G$	2			
1731	$\mathcal{N}(g_1) = n_1 \Longrightarrow G_1$				
1732	$N(g_2) = n_2 \Rightarrow G_2$				
1733	$\Gamma \cdot \Lambda \vdash a_i \cdot \tau_i$				
1734	$\Gamma, \Delta \vdash g_1 : \tau_1$ $\Gamma, \Delta \vdash e \vdash a_2 : \tau_2$				
1735	By I.H. on a_1 , we have				
1727	$\forall w_1 \leq \mathbb{n}, w_1 \in [q_1]_{(v,\delta)}$ if and only if ($a' \vdash n_1 \rightsquigarrow^* w_1.$			
1738	By I.H. on q_2 , we have the following. Here	e we use < instead of	\leq as its typing context Γ, Δ		
1730	includes Δ that only has interpretations	for $ w_2 < \mathbb{n}$.	- ,1 0 ,		
1740	$\forall w_2 < \mathbb{n}, w_2 \in \llbracket q_2 \rrbracket_{(\chi, \delta)}$ if and only if ($n' \vdash n_2 \rightsquigarrow^* w_2.$			
1741	We first prove the conclusion from left to	right. Given $w \leq n$, an	d $w \in [q_1 \cdot q_2]_{(v, \delta)}$, it must		
1742	be $w = w_1 \cdot w_2$ and $w_1 \in \llbracket q_1 \rrbracket_{(u, \delta)}$ and	$w_2 \in \llbracket q_2 \rrbracket_{(1,\delta)}$. As $q_1 \cdot$	q_2 is well-typed, we know		
1743	τ_1 .Null = false, so w_1 cannot be empty,	and thus w_2 must have	ve length $< \mathbb{n}$. So following		
1744	I.H., and that <i>n</i> represents the same lang	uage as $n_1 n_2$, we have	$G' \vdash n \rightsquigarrow^* w_1 \cdot w_2.$		
1745	Now we move to the conclusion from	right to left. Given ($G' \vdash n \rightsquigarrow^* w$, it must be		
1746	$w = w_1 \cdot w_2$, and $G' \vdash n_1 \rightsquigarrow^* w_1$, and G'	$n_2 \rightsquigarrow^* w_2$. As $g_1 \cdot g_2$	is well-typed, we know that		
1747	τ_1 .Null = false, so by Lemma 4.2, w_1 car	not be empty, and thu	is w_2 must have length $< \square$.		
1748	So following I.H., we have $w_1 \in \llbracket g_1 \rrbracket_{(v,\delta)}$, and $w_2 \in [\![g_2]\!]_{(\gamma,\delta)}, =$	and thus $w \in \llbracket g_1 \cdot g_2 \rrbracket_{(v,\delta)}$.		
1749	(),	(,,,,)	(1,0)		

1750	• $q = \mu \alpha$. q_1 . Then $\llbracket \mu \alpha$. $q_1 \rrbracket_{(\nu, \delta)} = \llbracket q_1 \rrbracket_{(\nu, \delta)} \llbracket q_2 \rrbracket_{(\nu, \delta)}$.
1751	We have
1752	$\boldsymbol{N}(\mu\alpha:\tau,g_1) = \alpha \Longrightarrow \{\alpha \to N \mid n \to N \in G\} \cup G$
1753	$\mathcal{N}(g_1) = n \Longrightarrow G$
1754	According to typing, we have Γ ; Δ , α : $\tau \vdash g_1 : \tau$.
1755	According to the I.H. on n , we have
1756	$\forall w' < \mathbb{n}, w' \in \llbracket \mu \alpha. g_1 \rrbracket_{(\gamma, \delta)}$ if and only if and $G' \vdash \alpha \rightsquigarrow^* w'$.
1757	We have $(\gamma, \delta, \llbracket \mu \alpha. g_1 \rrbracket_{(\gamma, \delta)} / \alpha)(\alpha) = \llbracket \mu \alpha. g_1 \rrbracket_{(\gamma, \delta)}.$
1758	That means we have
1759	$\forall \beta \in \operatorname{dom}\left(\delta, \left[\!\left[\mu \alpha. g_1\right]\!\right]_{(\gamma, \delta)} / \alpha\right),$
1760	$\forall w' < \mathbb{n}, w' \in (\delta, \llbracket \mu \alpha, g_1 \rrbracket_{(\gamma, \delta)} / \alpha)(\beta)$ if and only if $G' \vdash beta \sim w'$.
1761	Now by I.H. on g_1 ,
1762	$\forall w \leq \mathbb{n}, w \in \llbracket g_1 \rrbracket_{(v, \delta, \llbracket u\alpha, g_1 \rrbracket_{(v, \delta)}/\alpha)}$ if and only if $G' \vdash n \sim w$
1763	equivalent to
1764	$\forall w \leq \mathbb{n}, w \in \llbracket \mu \alpha, q_1 \rrbracket_{(w, \delta)}$ if and only if $G' \vdash n \rightsquigarrow^* w$.
1765	We have $\alpha \to N \in \mathcal{N}(\mu\alpha, q_1)$, where $n \to N \in \mathcal{N}(\mu\alpha, q_1)$. By Lemma F.2, we have
1766	$\alpha \to N \in G'$ and there is no other productions for α .
1767	Therefore,
1768	$\forall w \leq \mathbb{n}, w \in \llbracket \mu \alpha, g_1 \rrbracket_{(\chi \delta)}$ if and only if $G' \vdash \alpha \rightsquigarrow w$.
1769	
1770	
1771	THEOREM 4.8 (SOUNDNESS). Given \bullet ; $\bullet \vdash g : \tau$ and $N [[g]]$ returns $n \Rightarrow G$, we have $w \in [[g]]_{\bullet}$ if
1772	and only if $G \vdash n \rightsquigarrow w$ for any w .
1774	PROOF. Follows by Lemma F.3, making use of Lemma F.1.
1775	
1776	
1777	
1778	
1779	
1780	
1781	
1782	
1783	
1784	
1785	
1786	
1787	
1788	
1789	
1790	
1791	
1792	
1793	
1794	
1795	
1796	
1797	
1798	