# Intrinsic Verification of Parsers and Formal Grammar Theory in Dependent Lambek Calculus

**Steven Schaefer** [1]    Nathan Varner [1]    Pedro H. Azevedo de Amorim [2]
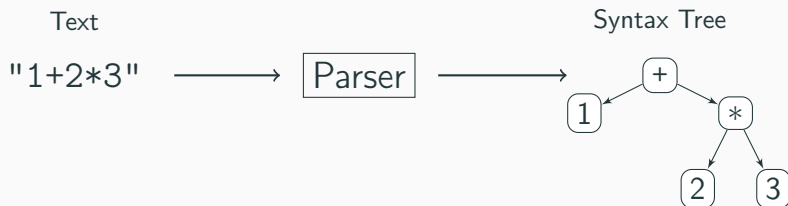Max S. New [1]
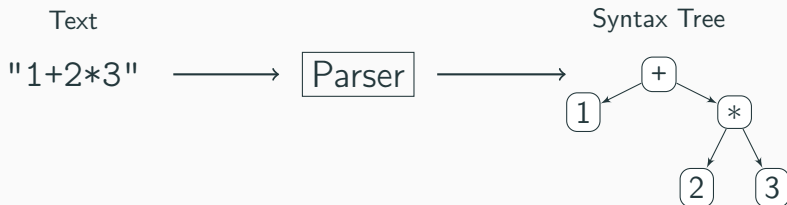
[1] University of Michigan

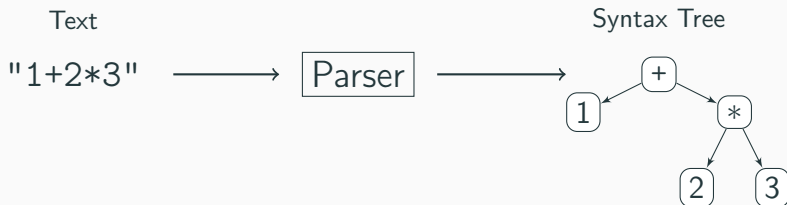[2] University of Oxford

July 24, 2025

**Text**

"1+2*3" ⟶ Parser ⟶

**Syntax Tree**

**The Cost of Incorrect Parsers**

**The Cost of Incorrect Parsers**

- Critical real-world vulnerabilities, including remote code execution **(Equifax, 2017)**

**The Cost of Incorrect Parsers**

- Critical real-world vulnerabilities, including remote code execution **(Equifax, 2017)**
- The guarantees of other software verification can be invalidated

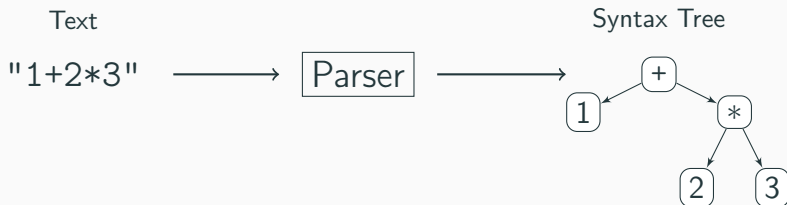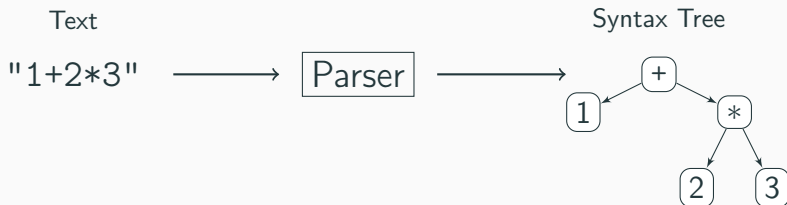**The Cost of Incorrect Parsers**

- Critical real-world vulnerabilities, including remote code execution **(Equifax, 2017)**
- The guarantees of other software verification can be invalidated
- Early versions of CompCert
  - 0 middle or backend bugs, 6 parser bugs **(Yang et al., 2011)**

For a formal grammar $A$, a parser is a function $f :$ **String** $\rightharpoonup$ **Parse**($A$)

For a formal grammar A, a parser is a function f : **String** $\rightharpoonup$ **Parse**(A)

**Soundness** For a string w, f(w) is a valid parse tree for the grammar A of w

For a formal grammar A, a parser is a function f : **String** $\rightharpoonup$ **Parse**(A)

**Soundness** For a string w, f(w) is a valid parse tree for the grammar A of w

**Completeness** Parser returns a parse tree whenever one exists

Existing verified parsers often implement a single formalism/algorithm

Existing verified parsers often implement a single formalism/algorithm

- **LL**, **LR**, **CYK**, parser combinators, Brzozowski derivatives, etc

Existing verified parsers often implement a single formalism/algorithm

- **LL**, **LR**, **CYK**, parser combinators, Brzozowski derivatives, etc
- Each formalism requires its own implementation

**Lambek$^{\mathrm{D}}$: DSL for *intrinsically* verified parsing**

**Lambek$^D$: DSL for *intrinsically* verified parsing**

- Functional programming language for writing verified parsers

### Lambek$^{\mathrm{D}}$: DSL for *intrinsically* verified parsing

- Functional programming language for writing verified parsers
- Ordered, linear type system guarantees parser soundness for free

## Lambek$^{\mathrm{D}}$: DSL for *intrinsically* verified parsing

- Functional programming language for writing verified parsers
- Ordered, linear type system guarantees parser soundness for free
  - Completeness with a little effort

# Dependent Lambek Calculus (Lambek$^{\text{D}}$)

**Lambek$^{\text{D}}$: DSL for *intrinsically* verified parsing**

- Functional programming language for writing verified parsers
- Ordered, linear type system guarantees parser soundness for free
  - Completeness with a little effort
- Extensible framework for implementing grammar formalisms, their parsers, and their parser generators

# Dependent Lambek Calculus (Lambek$^D$)

## Lambek$^D$: DSL for *intrinsically* verified parsing

- Functional programming language for writing verified parsers
- Ordered, linear type system guarantees parser soundness for free
  - Completeness with a little effort
- Extensible framework for implementing grammar formalisms, their parsers, and their parser generators

Universal formalism that expresses grammars as types

# Dependent Lambek Calculus (Lambek$^D$)

## Lambek$^D$: DSL for *intrinsically* verified parsing

- Functional programming language for writing verified parsers
- Ordered, linear type system guarantees parser soundness for free
  - Completeness with a little effort
- Extensible framework for implementing grammar formalisms, their parsers, and their parser generators

Universal formalism that expresses grammars as types

Implemented as an embedded DSL in Cubical Agda ✍

- Verified parsers for selected context-free grammars
- Verified parser generator for regular expressions

Parsing is a deductive problem

$$\frac{\dfrac{\cdots}{\texttt{"a"} \vdash \texttt{'a'} \oplus \texttt{'b'}} \qquad \dfrac{\cdots}{\texttt{"b"} \vdash (\texttt{'a'} \oplus \texttt{'b'})^{*}}}{\texttt{"ab"} \vdash (\texttt{'a'} \oplus \texttt{'b'})^{*}}$$

Parsing is a deductive problem
- User provides a type-level specification

$$\frac{\dfrac{\cdots}{\texttt{"a"} \vdash \texttt{'a'} \oplus \texttt{'b'}} \qquad \dfrac{\cdots}{\texttt{"b"} \vdash (\texttt{'a'} \oplus \texttt{'b'})^*}}{\texttt{"ab"} \vdash (\texttt{'a'} \oplus \texttt{'b'})^*}$$

Parsing is a deductive problem

- User provides a type-level specification
- Parse trees are synthesized via logic programming

$$\frac{\cdots}{\texttt{"a"} \vdash \texttt{'a'} \oplus \texttt{'b'}} \qquad \frac{\cdots}{\texttt{"b"} \vdash (\texttt{'a'} \oplus \texttt{'b'})^{*}}$$
$$\frac{}{\texttt{"ab"} \vdash (\texttt{'a'} \oplus \texttt{'b'})^{*}}$$

# Lambek Calculus: The Basis for Lambek$^D$

Parsing is a deductive problem

- User provides a type-level specification
- Parse trees are synthesized via logic programming

$$\frac{\dfrac{\cdots}{\text{"a"} \vdash \text{'a'} \oplus \text{'b'}} \qquad \dfrac{\cdots}{\text{"b"} \vdash (\text{'a'} \oplus \text{'b'})^*}}{\text{"ab"} \vdash (\text{'a'} \oplus \text{'b'})^*}$$

Lambek$^D$ extends the expressivity of Lambek Calculus by adding dependent types

- Characters `'c'`

# Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`

# Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`
- Concatenation $\otimes$

## Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`
- Concatenation $\otimes$
- Disjunction $(\oplus, 0)$

## Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`
- Concatenation $\otimes$
- Disjunction $(\oplus, 0)$
- Conjunction $(\&, \top)$

## Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`
- Concatenation $\otimes$
- Disjunction ($\oplus$, $0$)
- Conjunction (`&`, $\top$)
- Kleene Star $\cdot^*$

# Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`
- Concatenation $\otimes$
- Disjunction $(\oplus, 0)$
- Conjunction $(\&, \top)$
- Kleene Star $\cdot^*$
  - Inductive linear type of lists

```
data A* : LinTy where
  nil  : I ⊸ A*
  cons : A ⊸ A* ⊸ A*
```

# Types in Lambek$^D$

- Characters `'c'`
- Empty string `I`
- Concatenation $\otimes$
- Disjunction $(\oplus, 0)$
- Conjunction $(\&, \top)$
- Kleene Star $\cdot^*$
  - Inductive linear type of lists
  - Arbitrary inductive types generalize to context-free grammars and beyond

```
data A* : LinTy where
  nil  : I ⊸ A*
  cons : A ⊸ A* ⊸ A*
```

# A Linear, Ordered Theory

> Parse trees of "cat"

$$x : \text{'c'}, y : \text{'a'}, z : \text{'t'} \vdash (x, (y, z)) : \text{'c'} \otimes \text{'a'} \otimes \text{'t'} \quad \checkmark$$

$$\boxed{\text{Parse trees of "cat"}}$$

$$x : \text{'c'}, y : \text{'a'}, z : \text{'t'} \vdash (x, (y, z)) : \text{'c'} \otimes \text{'a'} \otimes \text{'t'} \quad \checkmark$$

Variables used in order (no exchange)

$$x : \text{'c'}, y : \text{'a'}, z : \text{'t'} \nvdash (y, (x, z)) : \text{'a'} \otimes \text{'c'} \otimes \text{'t'} \quad \textcolor{red}{\times}$$

# A Linear, Ordered Theory

> **Parse trees of `"cat"`**

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \vdash (x, (y, z)) : \texttt{'c'} \otimes \texttt{'a'} \otimes \texttt{'t'} \quad ✓$$

Variables used in order (no exchange)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (y, (x, z)) : \texttt{'a'} \otimes \texttt{'c'} \otimes \texttt{'t'} \quad ✗$$

At least once (no weakening)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (y, z) : \texttt{'a'} \otimes \texttt{'t'} \quad ✗$$

# A Linear, Ordered Theory

$$\boxed{\text{Parse trees of "cat"}}$$

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \vdash (x, (y, z)) : \texttt{'c'} \otimes \texttt{'a'} \otimes \texttt{'t'} \quad \checkmark$$

Variables used in order (no exchange)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (y, (x, z)) : \texttt{'a'} \otimes \texttt{'c'} \otimes \texttt{'t'} \quad \textcolor{red}{\boldsymbol{\times}}$$

At least once (no weakening)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (y, z) : \texttt{'a'} \otimes \texttt{'t'} \quad \textcolor{red}{\boldsymbol{\times}}$$

At most once (no contraction)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (x, (y, (z, z))) : \texttt{'c'} \otimes \texttt{'a'} \otimes \texttt{'t'} \otimes \texttt{'t'} \quad \textcolor{red}{\boldsymbol{\times}}$$

# A Linear, Ordered Theory

$$\boxed{\text{Parse trees of "cat"}}$$

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \vdash (x, (y, z)) : \texttt{'c'} \otimes \texttt{'a'} \otimes \texttt{'t'} \quad ✓$$

Variables used in order (no exchange)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (y, (x, z)) : \texttt{'a'} \otimes \texttt{'c'} \otimes \texttt{'t'} \quad ✗$$

At least once (no weakening)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (y, z) : \texttt{'a'} \otimes \texttt{'t'} \quad ✗$$

At most once (no contraction)

$$x : \texttt{'c'}, y : \texttt{'a'}, z : \texttt{'t'} \nvdash (x, (y, (z, z))) : \texttt{'c'} \otimes \texttt{'a'} \otimes \texttt{'t'} \otimes \texttt{'t'} \quad ✗$$

These restrictions ensure parsers are sound-by-construction

All strings parsed by $(A \otimes A)^*$ are also parsed by $A^*$

$$aas : (A \otimes A)^* \vdash \texttt{flatten } aas : A^*$$

All strings parsed by $(A \otimes A)^*$ are also parsed by $A^*$

$$aas : (A \otimes A)^* \vdash \text{flatten } aas : A^*$$

```
flatten : (A ⊗ A)* ⊸ A*
flatten nil = nil
flatten (cons (a , a') aas) =
   cons a (cons a' (flatten aas))
```

Context-free grammar of balanced parentheses

$$D \to \varepsilon \mid (D)D$$

Context-free grammar of balanced parentheses

$$D \rightarrow \varepsilon \mid (D)D$$

```
data Dyck : LinTy where
  nil : I ⊸ Dyck
  bal : '(' ⊸ Dyck ⊸ ')' ⊸ Dyck ⊸ Dyck
```

Context-free grammar of balanced parentheses

$$D \to \varepsilon \mid (D)D$$

```
data Dyck : LinTy where
  nil :  I  ⊸ Dyck
  bal : '(' ⊸ Dyck ⊸ ')' ⊸ Dyck ⊸ Dyck
```

$l : \text{'('}, r : \text{')'}, l' : \text{'('}, r' : \text{')'} \vdash bal\ l\ nil\ r\ (bal\ l'\ nil\ r') : \text{Dyck}$

**Dyck Parser**

A parser for Dyck is a function

$$\text{String} \multimap \text{Dyck} \oplus \top$$

# A Parser for the Dyck Grammar

**Dyck Parser**

A parser for Dyck is a function

$$\text{String} \multimap \text{Dyck} \oplus \top$$

**Soundness** Guaranteed for free by the type system

**Dyck Parser**

A parser for Dyck is a function

$$\texttt{String} \multimap \texttt{Dyck} \oplus \top$$

**Soundness** Guaranteed for free by the type system

**Completeness** Not guaranteed, as the parser can choose to always fail

# A Parser for the Dyck Grammar

**Complete Dyck Parser**

A *complete* parser for `Dyck` is a choice of type `Dyck¬` and function

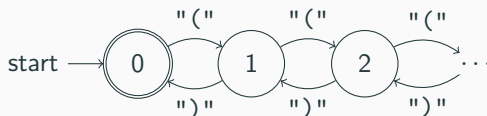$$\texttt{String} \multimap \texttt{Dyck} \oplus \texttt{Dyck}_\neg$$

where `Dyck & Dyck¬` $\cong 0$

**Complete Dyck Parser**

A *complete* parser for `Dyck` is a choice of type `Dyck`$_\neg$ and function

$$\texttt{String} \multimap \texttt{Dyck} \oplus \texttt{Dyck}_\neg$$

where `Dyck & Dyck`$_\neg \cong 0$

To define the parser we will use an automaton
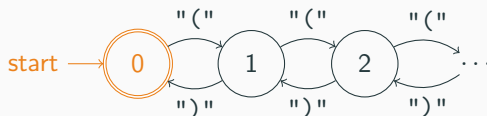
```
data Trace : Bool → ℕ → LinTy where
  eof : Trace true 0
  leftovers : ∀ n → Trace false (suc n)
  push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
  pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
  unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace $b$ $n$ denotes a path through the automaton from state $n$ ending in some state with acceptance criteria $b$
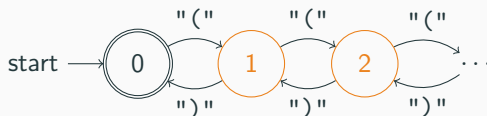
# Deterministic Automaton for the Dyck Grammar



```
data Trace : Bool → ℕ → LinTy where
 eof : Trace true 0
 leftovers : ∀ n → Trace false (suc n)
 push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
 pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
 unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace *b n* denotes a path through the automaton from state *n* ending in some state with acceptance criteria *b*
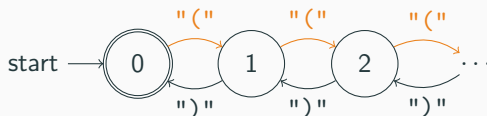
# Deterministic Automaton for the Dyck Grammar



```
data Trace : Bool → ℕ → LinTy where
 eof : Trace true 0
 leftovers : ∀ n → Trace false (suc n)
 push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
 pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
 unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace *b n* denotes a path through the automaton from state *n* ending in
some state with acceptance criteria *b*

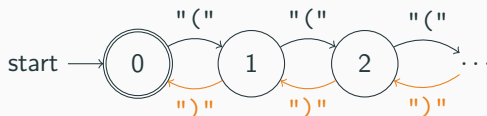# Deterministic Automaton for the Dyck Grammar



```
data Trace : Bool → ℕ → LinTy where
 eof : Trace true 0
 leftovers : ∀ n → Trace false (suc n)
 push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
 pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
 unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace *b n* denotes a path through the automaton from state *n* ending in
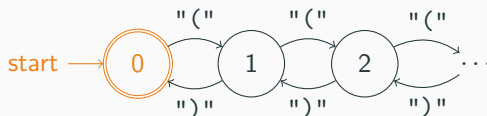some state with acceptance criteria *b*

# Deterministic Automaton for the Dyck Grammar



```
data Trace : Bool → ℕ → LinTy where
 eof : Trace true 0
 leftovers : ∀ n → Trace false (suc n)
 push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
 pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
 unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace *b n* denotes a path through the automaton from state *n* ending in some state with acceptance criteria *b*
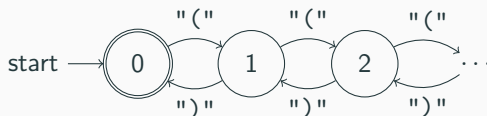
# Deterministic Automaton for the Dyck Grammar



```
data Trace : Bool → ℕ → LinTy where
 eof : Trace true 0
 leftovers : ∀ n → Trace false (suc n)
 push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
 pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
 unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace *b n* denotes a path through the automaton from state *n* ending in some state with acceptance criteria *b*

```
data Trace : Bool → ℕ → LinTy where
 eof : Trace true 0
 leftovers : ∀ n → Trace false (suc n)
 push : ∀ b n → '(' ⊸ Trace b (suc n) ⊸ Trace b n
 pop : ∀ b n → ')' ⊸ Trace b n ⊸ Trace b (suc n)
 unexpected : ')' ⊸ ⊤ ⊸ Trace false 0
```

Trace *b n* denotes a path through the automaton from state *n* ending in some state with acceptance criteria *b*

Looks like an ordinary functional program

```
parse : String ⊸ &[ n ∈ ℕ ] (Trace true n ⊕ Trace false n)
parse nil zero = σ true eof
parse nil (suc n) = σ false leftovers
parse (cons (σ '(' a) w) n =
  let σ b tr = parse w (suc n) in
  σ b (push a tr)
parse (cons (σ ')' a) w) zero =
  σ false (unexpected a _)
parse (cons (σ ')' a) w) (suc n) =
  let σ b tr = parse w n in
  σ b (pop a tr)
```

**Deterministic Automaton**

$$\text{String} \cong \text{Trace true } 0 \oplus \text{Trace false } 0$$

# A Parser for the Dyck Grammar

**Deterministic Automaton**

$$\text{String} \cong \text{Trace true } 0 \oplus \text{Trace false } 0$$

**The Dyck Grammar is Recognized by the Automaton**

$$\text{Dyck} \cong \text{Trace true } 0$$

# A Parser for the Dyck Grammar

**Deterministic Automaton**

$$\text{String} \cong \text{Trace true } 0 \oplus \text{Trace false } 0$$

**The Dyck Grammar is Recognized by the Automaton**

$$\text{Dyck} \cong \text{Trace true } 0$$

$$\text{String} \multimap \text{Trace true } 0 \oplus \text{Trace false } 0$$
$$\multimap \text{Dyck} \oplus \text{Trace false } 0$$

## Strategy for Building a Parser

1. Describe a grammar $A$ and the traces of an automaton each as
   *linear types*
2. In the equational theory of Lambek$^D$, prove $A$ is equivalent to the
   type of traces
3. Use the equivalence to create a parser for $A$

## Strategy for Building a Parser

1. Describe a grammar $A$ and the traces of an automaton each as *linear types*
2. In the equational theory of Lambek$^D$, prove $A$ is equivalent to the type of traces
3. Use the equivalence to create a parser for $A$

Instantiated the same strategy in

- Parser for an **LL**(1) grammar of arithmetic expressions and a lookahead automaton

## Strategy for Building a Parser

1. Describe a grammar $A$ and the traces of an automaton each as *linear types*
2. In the equational theory of Lambek$^D$, prove $A$ is equivalent to the type of traces
3. Use the equivalence to create a parser for $A$

Instantiated the same strategy in

- Parser for an **LL**(1) grammar of arithmetic expressions and a lookahead automaton
- *Parser generator* for regular expressions and finite automata
  - Thompson's construction and Powerset construction

Proving unambiguity of grammars in Lambek$^D$

## In the Paper

Proving unambiguity of grammars in Lambek$^D$

Weak equivalence v. strong equivalence of grammars

## In the Paper

Proving unambiguity of grammars in Lambek$^D$

Weak equivalence v. strong equivalence of grammars

Linear-Non-Linear Dependent Types

$$\bigoplus_{x:X} A\ x \qquad\qquad \bigotimes_{x:X} A\ x$$

Key to describing grammars of unrestricted complexity

## In the Paper

Proving unambiguity of grammars in Lambek$^D$

Weak equivalence v. strong equivalence of grammars

Linear-Non-Linear Dependent Types

$$\bigoplus_{x:X} A\ x \qquad\qquad \underset{x:X}{\&}\ A\ x$$

Key to describing grammars of unrestricted complexity

Denotational Semantics

Interpret grammars in the category **String** $\rightarrow$ **Set**

- Encode richer parsing algorithms (e.g. **LL**, **LR**, and forms of context-sensitivity)

- Encode richer parsing algorithms (e.g. **LL**, **LR**, and forms of context-sensitivity)
- Verified semantic actions to integrate with language backends

## Ongoing and Future Work

- Encode richer parsing algorithms (e.g. **LL**, **LR**, and forms of context-sensitivity)
- Verified semantic actions to integrate with language backends
- Verified frontends (typechecking and semantic analysis) in a modified version of Lambek$^D$

# Dependent Lambek Calculus (Lambek$^{\text{D}}$)

## Lambek$^{\text{D}}$: DSL for *intrinsically* verified parsing

- Functional programming language for writing verified parsers
- Ordered, linear type system guarantees parser soundness for free
  - Completeness with a little effort
- Extensible framework for implementing grammar formalisms, their parsers, and their parser generators

Universal formalism that expresses grammars as types

Implemented as an embedded DSL in Cubical Agda ✍

- Verified parsers for selected context-free grammars
- Verified parser generator for regular expressions