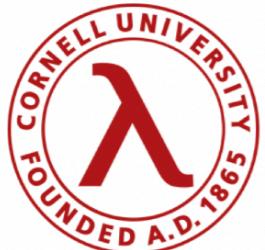


Certified Parsing of Dependent Regular Grammars



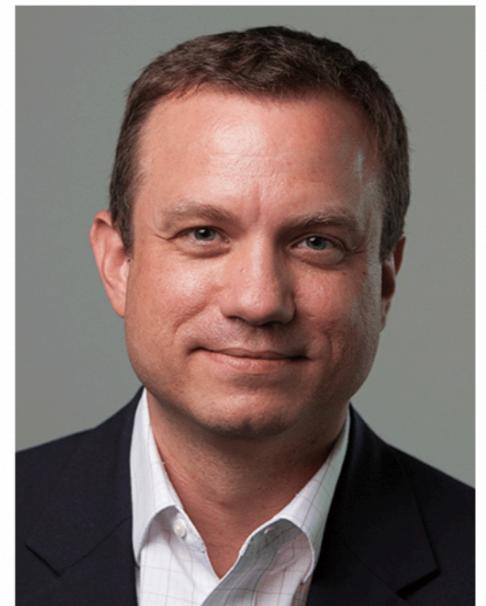
John Sarracino



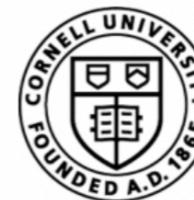
Gang Tan



PennState®



Greg Morrisett



**CORNELL
TECH**

Format Parsing is Hard

Name	Description
CVE-2021-38644	Microsoft MPEG-2 Video Extension Remote Code Execution Vulnerability
CVE-2021-38381	Live555 through 1.08 does not handle MPEG-1 or 2 files properly. Sending two successive RTSP SETUP commands for the same track causes a Use-After-Free and daemon crash.
CVE-2021-36937	Windows Media MPEG-4 Video Decoder Remote Code Execution Vulnerability
CVE-2021-21862	Multiple exploitable integer truncation vulnerabilities exist within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. The implementation of the parser used for the '“Xtra”' FOURCC code is handled. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21861	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. When processing the 'hdlr' FOURCC code, a specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21860	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. The FOURCC code, 'trik', is parsed by the function within the library. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21859	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. The stri_box_read function is used when processing atoms using the 'stri' FOURCC code. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21858	Multiple exploitable integer overflow vulnerabilities exist within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an integer overflow due to unchecked addition arithmetic resulting in a heap-based buffer overflow that causes memory corruption. An attacker can convince a user to open a video to trigger this vulnerability.

MPEG CVEs, -Meridith Patterson

Format Parsing is Hard

Name	Description
CVE-2021-38644	Microsoft MPEG-2 Video Extension Remote Code Execution Vulnerability
CVE-2021-38381	Live555 through 1.08 does not handle MPEG-1 or 2 files properly. Sending two successive RTSP SETUP commands for the same track causes a Use-After-Free and daemon crash.
CVE-2021-38645	Advanced Content library v1.0.1. When processing the 'hdlr' FOURCC code, a specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21860	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an integer overflow due to unchecked addition arithmetic resulting in a heap-based buffer overflow that causes memory corruption. An attacker can convince a user to open a video to trigger this vulnerability.

This is known as Dependent Parsing because parse results depend on prior parsed values

MPEG CVEs, -Meridith Patterson

Dependent Parsing: Netstrings

Netstrings: length-prefixed, colon and
comma delimited encoding of ascii values

Dependent Parsing: Netstrings

Netstrings: length-prefixed, colon and
comma delimited encoding of ascii values

Valid

5:hello,

12:hello world,

Dependent Parsing: Netstrings

Netstrings: length-prefixed, colon and
comma delimited encoding of ascii values

Valid

5:hello,
12:hello world,

Invalid

3:hello,
5:hello

Dependent Parsing: Netstrings

Netstrings: length-prefixed, colon and
comma delimited encoding of ascii values

Valid

5:hello,
12:hello world,

Invalid

3:hello,
5:hello

Netstrings are *dependent* because the
body depends on the value of the prefix

Dependent Regular Grammars

- We extend regular languages with monadic bind.
- Enables reasoning about dependent parsers within a proof assistant.
- Implemented and formalized in Coq.

Outline

- Syntax and Semantics
- Parsing algorithm and correctness
- Netstring case study

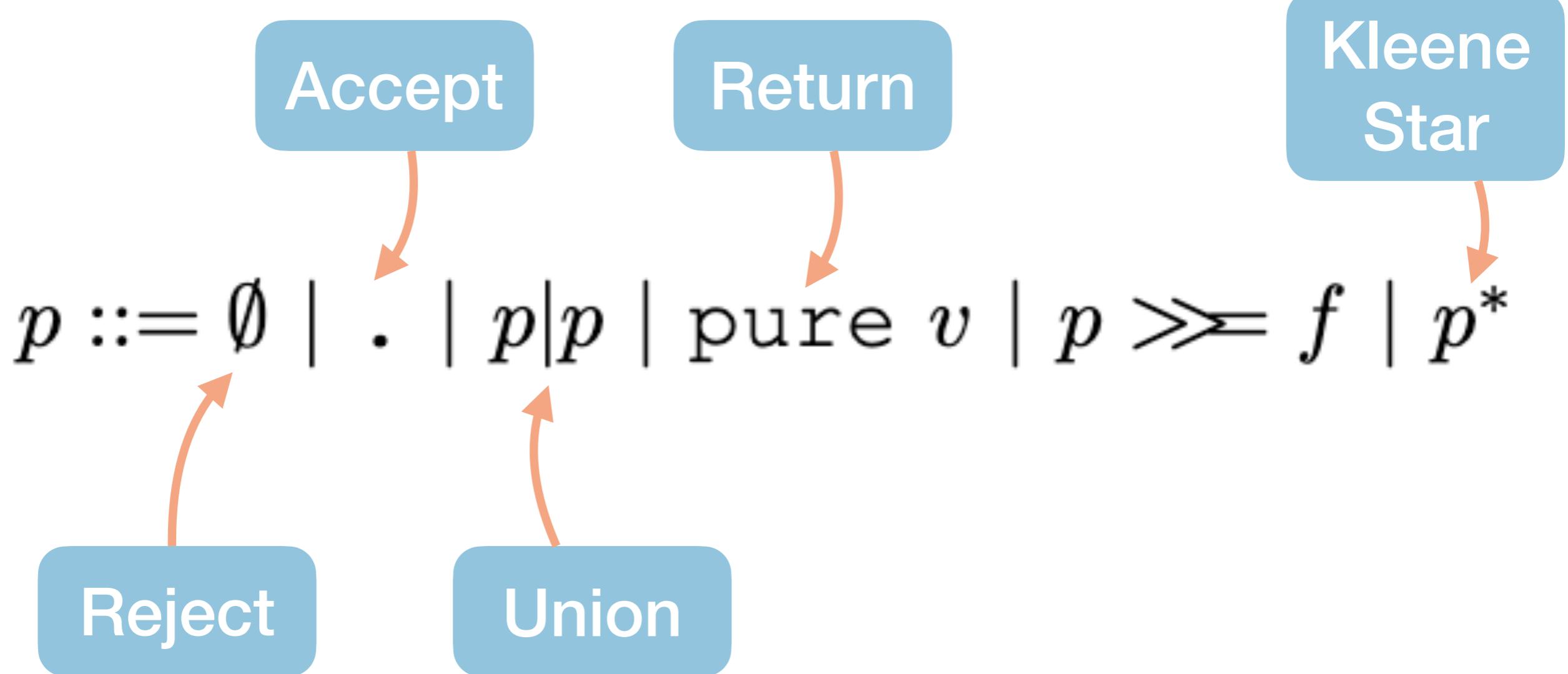
Outline

- **Syntax and Semantics**
- Parsing algorithm and correctness
- Netstring case study

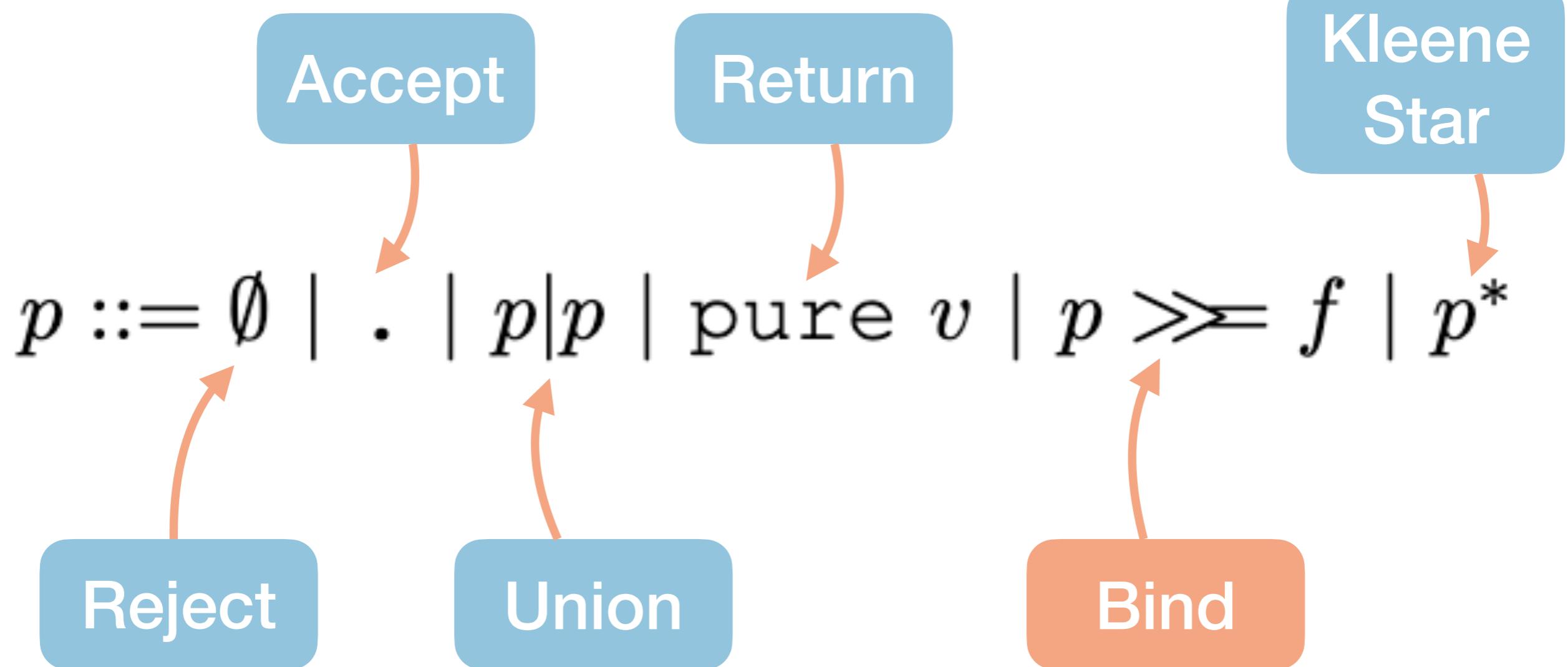
Dependent Parser Syntax

$$p ::= \emptyset \mid . \mid p|p \mid \text{pure } v \mid p \gg f \mid p^*$$

Dependent Parser Syntax



Dependent Parser Syntax



Bind Semantics

p parses s to v: $(s, p) \rightarrow v$

Bind Semantics

p parses s to v: $(s, p) \rightarrow v$

PBIND

$$\frac{(s_l, p) \rightarrow x \quad (s_r, f\ x) \rightarrow y}{(s_l ++ s_r, p \gg f) \rightarrow y}$$

Bind Semantics

p parses s to v : $(s, p) \rightarrow v$

If p parses s_l to x ,

$$\text{PBIND} \quad \frac{(s_l, p) \rightarrow x \quad (s_r, f\ x) \rightarrow y}{(s_l ++ s_r, p \gg f) \rightarrow y}$$

Bind Semantics

p parses s to v : $(s, p) \rightarrow v$

If p parses s_l to x ,

and $f x$ parses s_r to y ,

PBIND

$$\frac{(s_l, p) \rightarrow x \quad (s_r, f x) \rightarrow y}{(s_l ++ s_r, p \gg f) \rightarrow y}$$

Bind Semantics

p parses s to v : $(s, p) \rightarrow v$

If p parses s_l to x ,

and $f x$ parses s_r to y ,

PBIND

$$\frac{(s_l, p) \rightarrow x \quad (s_r, f x) \rightarrow y}{(s_l ++ s_r, p \gg f) \rightarrow y}$$

then $p \gg f$ parses $s_l ++ s_r$ to y

Parsing an Even Digit

digit $\gg= \lambda x \rightarrow \text{if } x \% 2 \text{ then return } x \text{ else fail}$

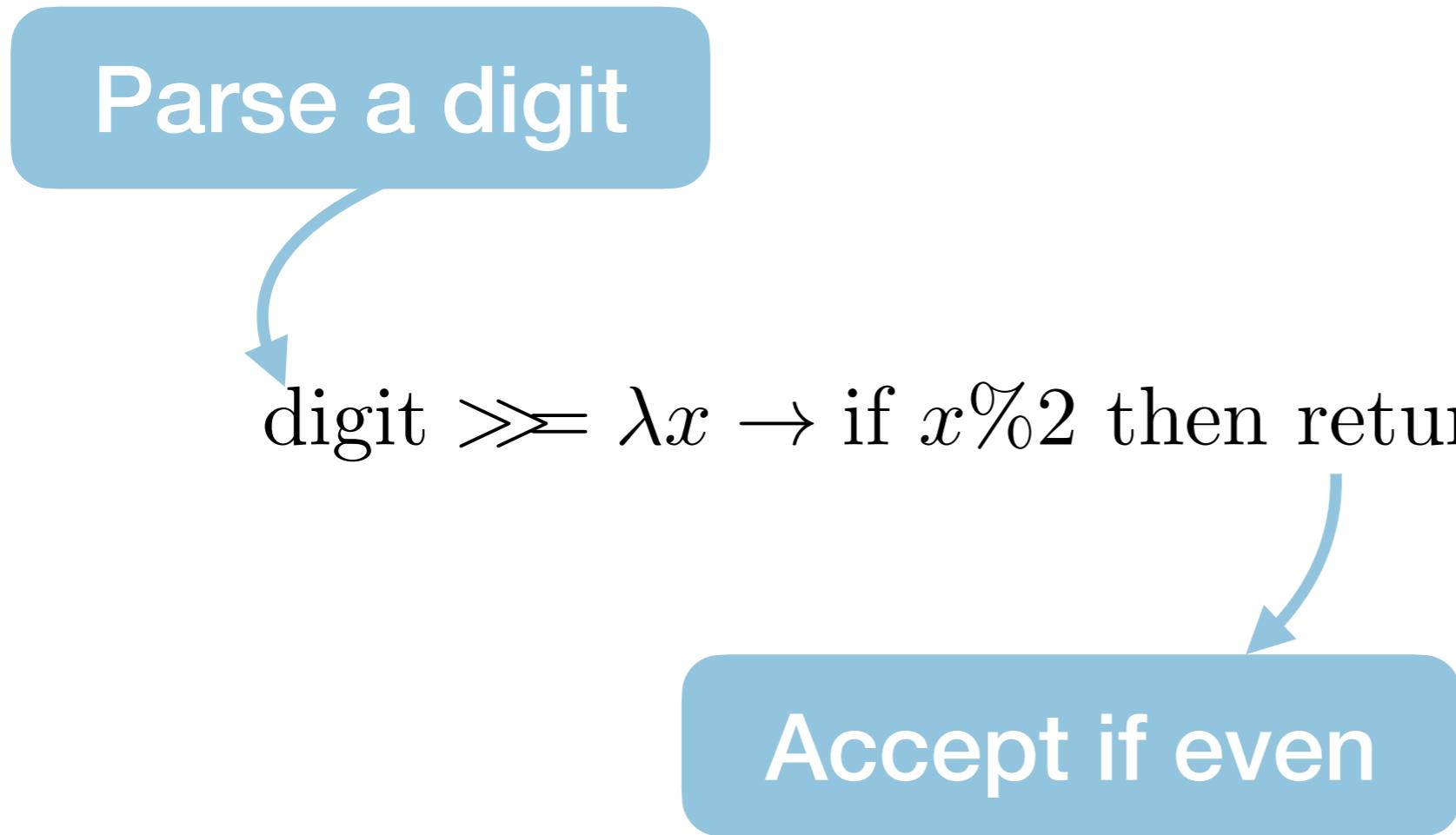
Parsing an Even Digit

Parse a digit

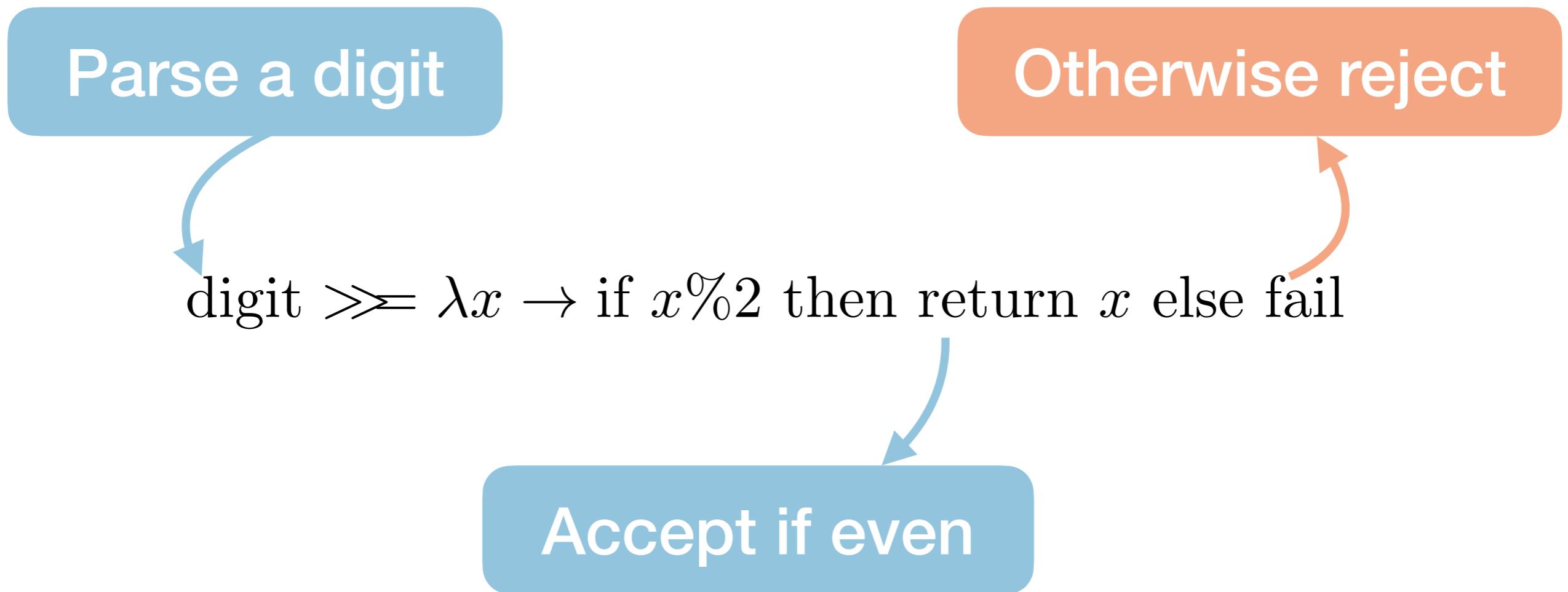


digit $\gg= \lambda x \rightarrow$ if $x \% 2$ then return x else fail

Parsing an Even Digit



Parsing an Even Digit



Outline

- **Syntax and Semantics**
- Parsing algorithm and correctness
- Netstring case study

Outline

- Syntax and Semantics
- **Parsing algorithm and correctness**
- Netstring case study

Parsing Algorithm

- How do we know that the actual parsing algorithm is implemented correctly?
- How do we know it terminates?
- We address both:
 - 1) we define the algorithm in Coq and prove correctness using the parsing semantics;
 - 2) we prove termination in Coq as well.

Parsing with Derivatives

We extend Brzozowski derivatives to Dependent Regular Grammars by defining two functions:

$$\| p \|_{\epsilon}$$

Finite set of return values of p

$$D_c p$$

Derivative of p with respect to c

We parse by iteratively applying derivatives and then building the set of return values.

Epsilon Interpretations

Specification: $v \in \| p \|_{\epsilon} \iff p \text{ parses } \epsilon \text{ to } v$

Parser

$\| p \|_{\epsilon}$

return 5 $\{5\}$

fail $\{\}$

return 2 | return 3 | . $\{2,3\}$

return 5 $\gg= \lambda n \rightarrow \text{return } n + 1$ $\{6\}$

digit* $\{\boxed{\hspace{1em}}\}$

Derivative Intuition

$$D_c p$$

Specification: p parses $c :: s$ to $v \iff D_c p$ parses s to v

Since the derivative preserves parse behavior,
we can parse by iteratively applying the
derivative operator.

Iterated Derivative

Parsing two characters on the string “xy”

. $\gg \lambda x \rightarrow . \gg \lambda y \rightarrow \text{return } (x, y)$

Iterated Derivative

Parsing two characters on the string “xy”

. $\gg \lambda x \rightarrow . \gg \lambda y \rightarrow \text{return } (x, y)$

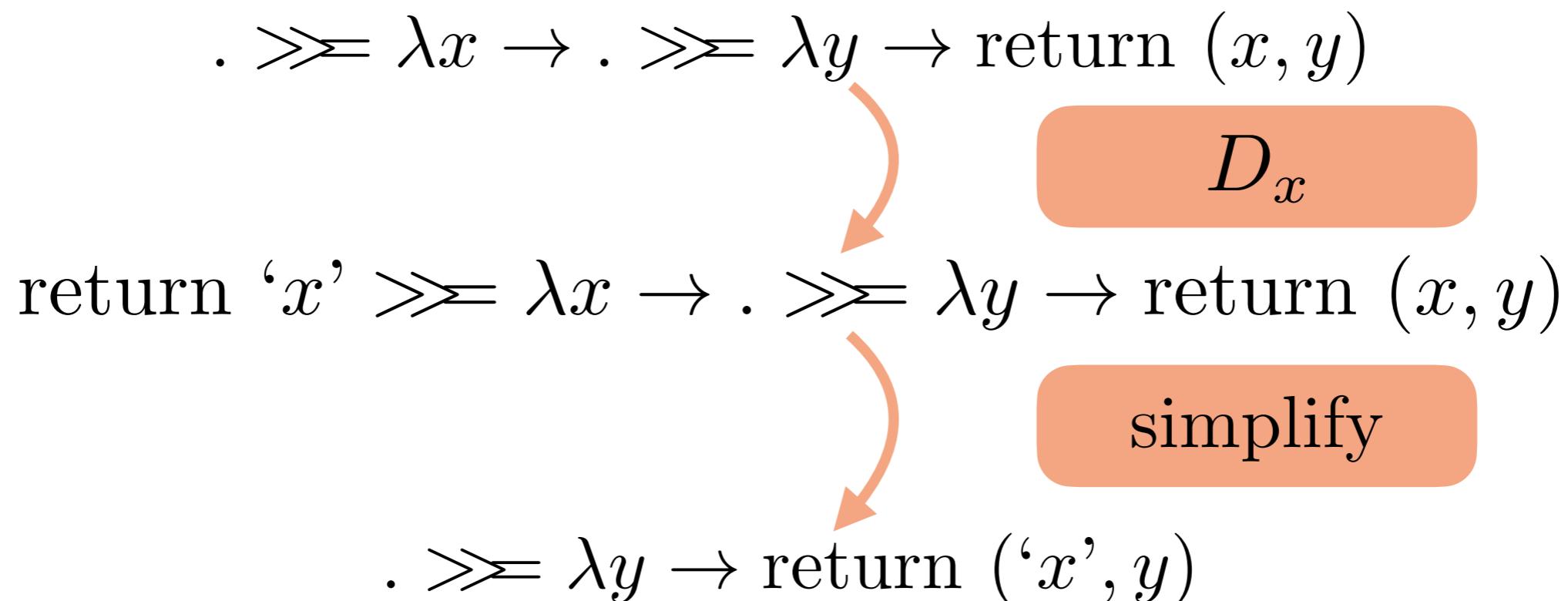
return ‘x’ $\gg \lambda x \rightarrow . \gg \lambda y \rightarrow \text{return } (x, y)$



D_x

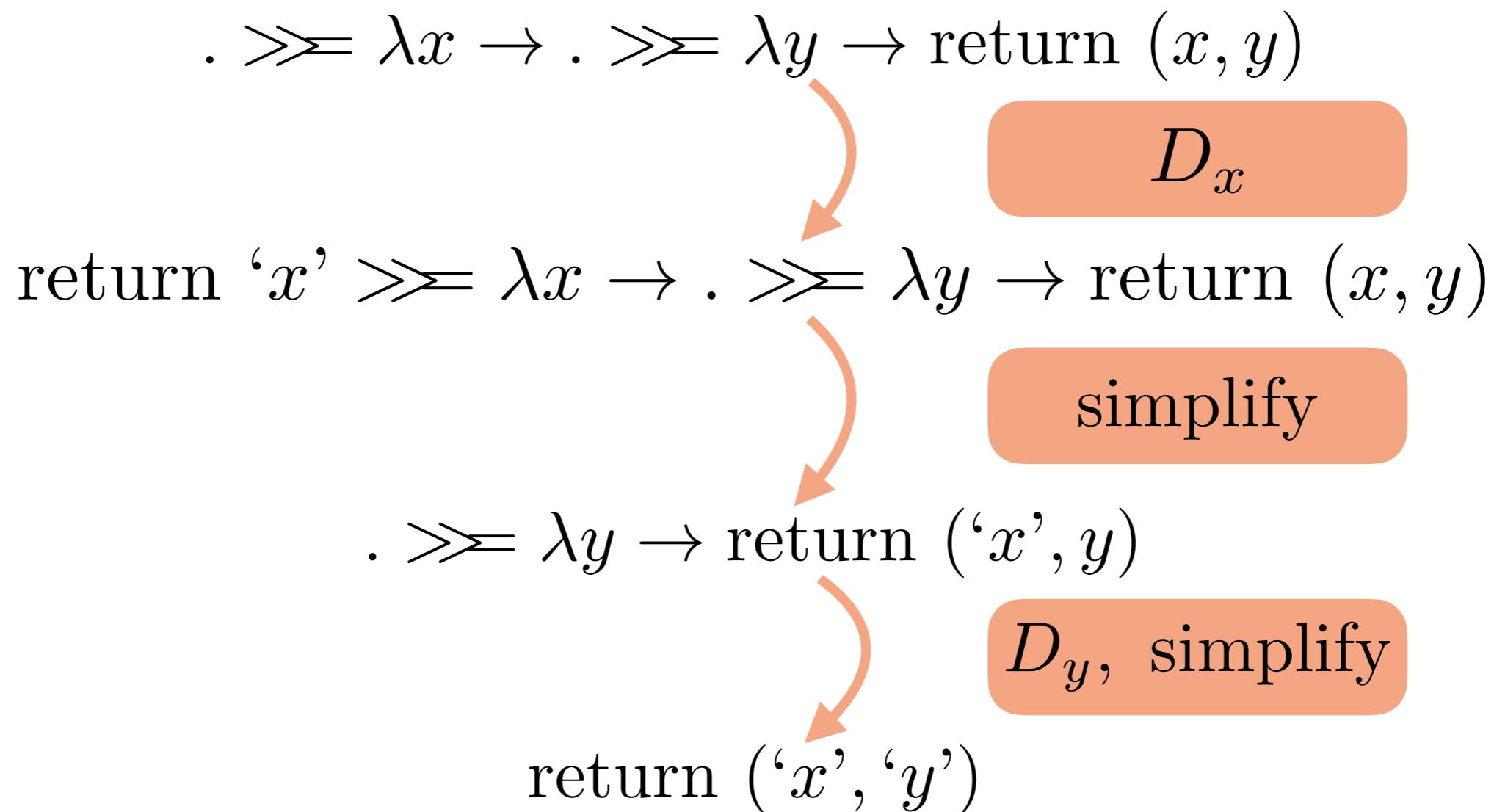
Iterated Derivative

Parsing two characters on the string “xy”



Iterated Derivative

Parsing two characters on the string “xy”



Outline

- Syntax and Semantics
- **Parsing algorithm and correctness**
- Netstring case study

Outline

- Syntax and Semantics
- Parsing algorithm and correctness
- **Netstring case study**

Case study: Netstrings

Proof of concept verified parser for a dependent format (known as calc-regular languages).

```
Definition net_str :=  
  num >>= fun n =>  
    char ":" $> ((repeat n pAny) <$ char ",").
```

Case study: Netstrings

Proof of concept verified parser for a dependent format (known as calc-regular languages).

Clear,
concise!!

```
Definition net_str :=  
  num >>= fun n =>  
    char ":" $> ((repeat n pAny) <$ char ",").
```

Case study: Netstrings

Proof of concept verified parser for a dependent format (known as calc-regular languages).

Clear,
concise!!

```
Definition net_str :=  
  num >>= fun n =>  
  char ":" $> ((repeat n pAny) <$ char ",").
```

Theorem net_str_spec:

```
forall s v,  
  Parses s net_str v <->  
  exists s',  
    s = s' ++ (":" :: v ++ ("," :: nil)) /\  
    Parses s' num (length v).
```

Case study: Netstrings

Proof of concept verified parser for a dependent format (known as calc-regular languages).

Clear,
concise!!

Strong spec!!

```
Definition net_str :=  
  num >>= fun n =>  
  char ":" $> ((repeat n pAny) <$ char ",").
```

Theorem net_str_spec:

```
forall s v,  
  Parses s net_str v <->  
  exists s',  
    s = s' ++ (":" :: v ++ ("," :: nil)) /\br/>    Parses s' num (length v).
```

Future directions

- Interpreting parsers in Coq is slow (too slow for realistic input sizes).
- I want to collaborate with **YOU** on verifying a security property of a real-world parser!

- Interpreting parsers in Coq is slow (too slow for realistic input sizes).
- I want to collaborate with **YOU** on verifying a security property of a real-world parser!

$$p ::= \emptyset \mid . \mid p|p \mid \text{pure } v \mid p \gg f \mid p^*$$

Thanks! Questions?

jsarracino@cornell.edu

<https://github.com/jsarracino/dependent-regular-grammars>

Backup slides

Dependent

Format parsing is hard

Name	Description
CVE-2021-38644	Microsoft MPEG-2 Video Extension Remote Code Execution Vulnerability
CVE-2021-38381	Live555 through 1.08 does not handle MPEG-1 or 2 files properly. Sending two successive RTSP SETUP commands for the same track causes a Use-After-Free and daemon crash.
CVE-2021-36937	Windows Media MPEG-4 Video Decoder Remote Code Execution Vulnerability
CVE-2021-21862	Multiple exploitable integer truncation vulnerabilities exist within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. The implementation of the parser used for the '“Xtra”' FOURCC code is handled. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21861	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. When processing the 'hdlr' FOURCC code, a specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21860	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an improper memory allocation resulting in a heap-based buffer overflow that causes memory corruption. The FOURCC code, 'trik', is parsed by the function in the library. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21859	An exploitable integer truncation vulnerability exists within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. The stri_box_read function is used when processing atoms using the 'stri' FOURCC code. An attacker can convince a user to open a video to trigger this vulnerability.
CVE-2021-21858	Multiple exploitable integer overflow vulnerabilities exist within the MPEG-4 decoding functionality of the GPAC Project on Advanced Content library v1.0.1. A specially crafted MPEG-4 input can cause an integer overflow due to unchecked addition arithmetic resulting in a heap-based buffer overflow that causes memory corruption. An attacker can convince a user to open a video to trigger this vulnerability.

Parser Semantics

$$\text{PANY} \quad \frac{}{(c :: \epsilon, .) \rightarrowtail c}$$

$$\text{PPURE} \quad \frac{(s, p_l) \rightarrowtail v}{(s, p_l | p_r) \rightarrowtail v}$$

$$\text{PALTL} \quad \frac{(s, p_r) \rightarrowtail v}{(s, p_l | p_r) \rightarrowtail v}$$

$$\text{PALTR}$$

$$\frac{}{(\epsilon, \text{pure } v) \rightarrowtail v}$$

$$\text{PSTAR0}$$

$$\frac{}{(\epsilon, p^*) \rightarrowtail []}$$

$$\text{PSTARITER}$$

$$\frac{s \neq \epsilon \quad (s, p) \rightarrowtail x \quad (s', p^*) \rightarrowtail xs}{(s + s', p^*) \rightarrowtail x :: xs}$$

$$\text{PBIND}$$

$$\frac{(s_l, p) \rightarrowtail x \quad (s_r, f\ x) \rightarrowtail y}{(s_l + s_r, p \gg f) \rightarrowtail y}$$

To ensure well-behaved parser semantics, Star interior is productive

Parser Syntax

Traditional

$$r ::= 1 \mid 0 \mid c \mid r_l \oplus r_r \mid r @ F \mid r^* \mid r_l \# r_r$$

Dependent

$$p ::= \emptyset \mid . \mid p|p \mid \text{pure } v \mid p \gg f \mid p^*$$

Interestingly, bind subsumes
map and concat

Epsilon and Derivative Definitions

$\| p \|_\epsilon$

$$\begin{aligned}
 \|\emptyset\|_\epsilon &= [] \\
 \|\cdot\|_\epsilon &= [] \\
 \|p_l | p_r\|_\epsilon &= \|p_l\|_\epsilon + \|p_r\|_\epsilon \\
 \|\text{pure } v\|_\epsilon &= [v] \\
 \|p^*\|_\epsilon &= [[]] \\
 \|p \gg f\|_\epsilon &= \text{concat_map } (\lambda x \Rightarrow \|f x\|_\epsilon) \ \|p\|_\epsilon
 \end{aligned}$$

$D_c p$

$$\begin{aligned}
 D_c \emptyset &= \emptyset \\
 D_c . &= \text{pure } c \\
 D_c p_l | p_r &= D_c p_l | D_c p_r \\
 D_c \text{pure } v &= \emptyset \\
 D_c p^* &= (D_c p) \$ p^* \\
 D_c p \gg f &= (D_c p \gg f) | \bigvee_{x \in \|p\|_\epsilon} D_c (f x)
 \end{aligned}$$