

# Slaying Yacc: Parsing with Derivatives

Matt Might  
University of Utah  
[matt.might.net](http://matt.might.net)



# A little history...



```

(define-struct ø      {})
(define-struct ε      {})
(define-struct token  {value})
(define-struct δ      {lang})
(define-struct ∪      {this that})
(define-struct °      {left right})
(define-struct ★      {lang})

(define (D c L)
  (match L
    [(_)
     (if (eqv? _ c) (ε) (ø))]
    [(ø)
     (ø)]
    [(ε)
     (ø)]
    [(δ _)
     (ø)]
    [(token a)
     (if (eqv? a c) (ε) (ø))]
    [(∪ L1 L2)
     (∪ (D c L1) (D c L2))]
    [(★ L1)
     (° (D c L1) L1)]
    [(° L1 L2)
     (∪ (° (δ L1) (D c L2))
         (° (D c L1) L2)))]))

(define (nullable? L)
  (match L
    [(_)
     #f]
    [(ε)
     #t]
    [(token _)
     #f]
    [(★ _)
     #t]
    [(δ L1)
     (nullable? L1)]
    [(∪ L1 L2)
     (or (nullable? L1)
         (nullable? L2))]
    [(° L1 L2)
     (and (nullable? L1)
          (nullable? L2)))]))

(define (recognizes? w p)
  (cond [(null? w) (nullable? p)]
        [else (recognizes? (cdr w) (D (car w) p))]))

```

```

(define-struct ø      {})
(define-struct ε      {})
(define-struct token  {value})
(define-struct δ      {lang})
(define-struct υ      {this that})
(define-struct °      {left right})
(define-struct ∗      {lang})

```

```

(define (D c L)
  (match L
    [((ø))          ((ø))]
    [(ε)]           ((ø))
    [(δ _)]         ((ø))
    [(token a)]     (if (eqv? a c) (ε) (ø)))
    [(υ L1 L2)]    (υ (D c L1) (D c L2)))
    [(∗ L1)]        (° (D c L1) L)
    [(° L1 L2)]    (υ (° (δ L1) (D c L2))
                    (° (D c L1) L2)))))

```

```

(define (nullable? L)
  (match L
    [((ø))]          #f]
    [(ε)]            #t]
    [(token _)]       #f]
    [(∗ _)]          #t]
    [(δ L1)]         (nullable? L1)]
    [(υ L1 L2)]     (or (nullable? L1)
                        (nullable? L2)))
    [(° L1 L2)]     (and (nullable? L1)
                        (nullable? L2)))))

```

```

(define (recognizes? w p)
  (cond [(null? w) (nullable? p)]
        [else (recognizes? (cdr w) (D (car w) p))]))

```

```

(define-struct ø      {})
(define-struct ε      {tree-set})
(define-struct token  {value?})
(define-lazy-struct δ {lang})
(define-lazy-struct υ {this that})
(define-lazy-struct ° {left right})
(define-lazy-struct ∗ {lang})
(define-lazy-struct → {lang reduce})

(define/memoize (D c p)
  #:order ([p #:eq] [c #:equal])
  (match p
    [((ø))]          ((ø))
    [(ε _)]          ((ø))
    [(δ _)]          ((ø))
    [(token p?)]    (if (p? c) (ε (set c)) (ø)))
    [(υ p1 p2)]    (υ (D c p1) (D c p2)))
    [(∗ p1)]        (° (D c p1) p))
    [(→ p1 f)]      (→ (D c p1) f))
    [(° p1 p2)]    (υ (° (δ p1) (D c p2))
                      (° (D c p1) p2)))))

(define/fix (parse-null p)
  #:bottom (set)
  (match p
    [(ε S)]          S]
    [((ø))]          (set))
    [(δ p)]          (parse-null p)]
    [(token _)]       (set))
    [(∗ _)]          (set '()))
    [(υ p1 p2)]     (set-union (parse-null p1)
                                (parse-null p2)))
    [(° p1 p2)]     (for*/set ([t1 (parse-null p1)]
                               [t2 (parse-null p2)])
                               (cons t1 t2)))
    [(→ p1 f)]       (for/set ([t (parse-null p1)])
                               (f t)))))

(define (parse w p)
  (cond [(null? w) (parse-null p)]
        [else (parse (cdr w) (D (car w) p))]))

```



+ Laziness

+ Memoization

+ Fixed points

# Laziness?

**Wait until you need it.**

# Memoization?

**Cache it once you know it.**

# Fixed points?

**Iterate until stable.**

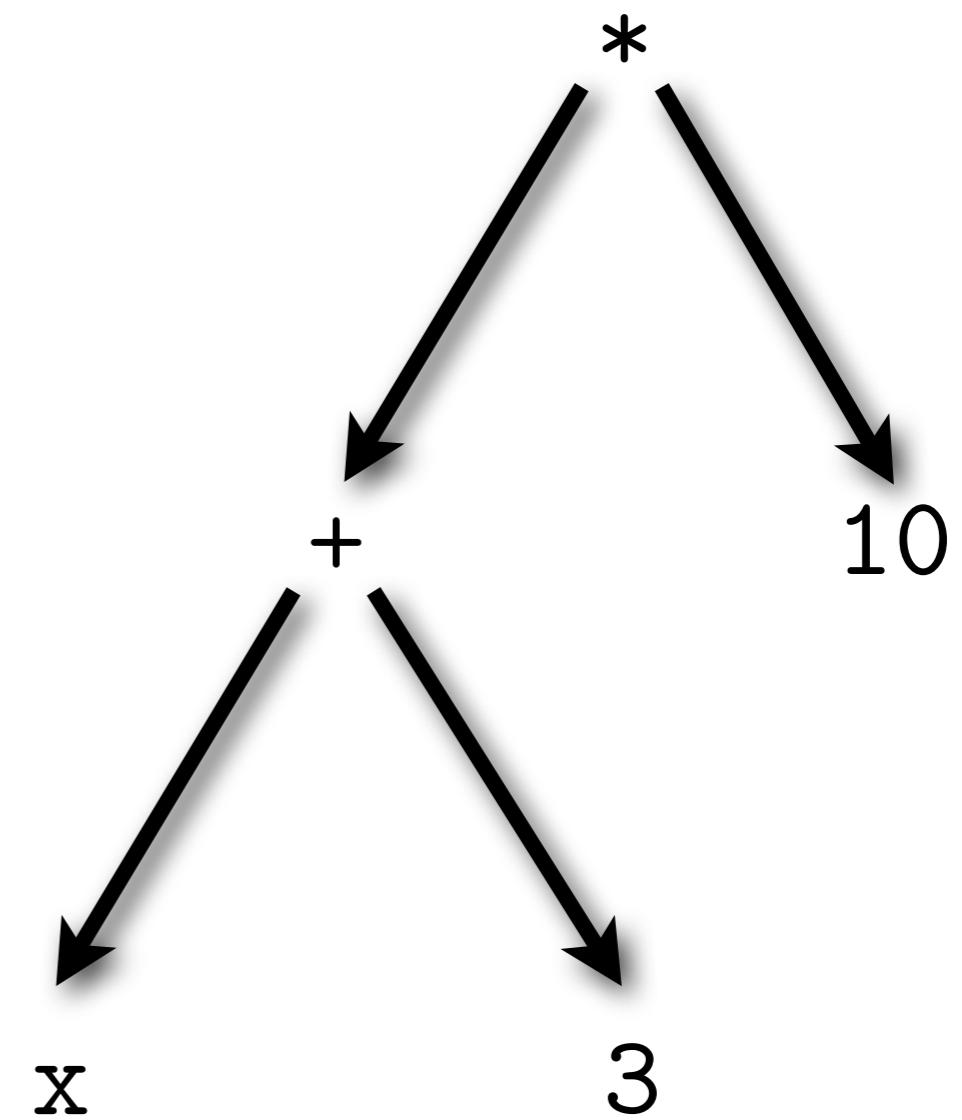
# What is parsing?

**Sequence => Tree**

$$(x + 3) * 10$$

$$(x + 3) * 10$$

(x + 3) \* 10



**“Parsing is ubiquitous.”**

Trevor Jim



# Compilers

## Interpreters

- Syntax highlighting

- Natural language

- Network protocols

- Command lines

- Configuration files

- Serialization

- Query languages

- API design

# Compilers

## Interpreters

Syntax highlighting

Natural language

Network protocols

Command lines

Configuration files

Serialization

Query languages

API design

# RFC 2616 (HTTP 1.1)

## 2.1 Augmented BNF

All of the mechanisms specified in this document are described in both prose and an augmented Backus-Naur Form (BNF) similar to that used by [RFC 822](#) [9]. Implementors will need to be familiar with the notation in order to understand this specification. The augmented BNF includes the following constructs:



```

media-type = type "/" subtype *( ";" parameter )
type      = token
subtype   = token

HTTP-Version = "HTTP" "/" 1*DIGIT "." 1*DIGIT

LWS = [ CRLF ] 1*( SP | HT )

```

```

separators = "(" | ")" | "<" | ">" | "@"
| "," | ";" | ":" | "\"" | "<*>"
| "/" | "[" | "]" | "?" | "="
| "{" | "}" | SP | HT

```

```
http_URL = "http:" "//" host [ ":" port ] [ abs_path [ "?" query ] ]
```

Chunked-Body	= *chunk last-chunk trailer CRLF	HTTP-date = <a href="#">rfc1123-date</a>   <a href="#">rfc850-date</a>   asctime-date <a href="#">rfc1123-date</a> = wkday , SP date1 SP time SP "GMT" <a href="#">rfc850-date</a> = weekday , SP date2 SP time SP "GMT" asctime-date = wkday SP date3 SP time SP 4DIGIT date1 = 2DIGIT SP month SP 4DIGIT ; day month year (e.g., 02 Jun 1982)
chunk	= chunk-size [ chunk-extension ] CRLF chunk-data CRLF	date2 = 2DIGIT "--" month "--" 2DIGIT ; day-month-year (e.g., 02-Jun-82)
chunk-size	= 1*HEX	date3 = month SP ( 2DIGIT   ( SP 1DIGIT )) ; month day (e.g., Jun 2)
last-chunk	= 1*( "0" ) [ chunk-extension ] CRLF	time = 2DIGIT ":" 2DIGIT ":" 2DIGIT ; 00:00:00 - 23:59:59
chunk-extension	= *( ";" chunk-ext-name [ "=" chunk-ext-val ] )	wkday = "Mon"   "Tue"   "Wed"   "Thu"   "Fri"   "Sat"   "Sun"
chunk-ext-name	= token	weekday = "Monday"   "Tuesday"   "Wednesday"   "Thursday"   "Friday"   "Saturday"   "Sunday"
chunk-ext-val	= token   quoted-string	month = "Jan"   "Feb"   "Mar"   "Apr"   "May"   "Jun"   "Jul"   "Aug"   "Sep"   "Oct"   "Nov"   "Dec"
chunk-data	= chunk-size(OCTET)	
trailer	= *(entity-header CRLF)	

# RFC 3501 (IMAPv4)



```

address      = "(" addr-name SP addr-adl SP addr-mailbox SP
              addr-host ")"
addr-adl     = nstring
              ; Holds route from [RFC-2822] route-addr if
              ; non-NIL
addr-host    = nstring
              ; NIL indicates [RFC-2822] group syntax.
              ; Otherwise, holds [RFC-2822] domain name
addr-mailbox = nstring
              ; NIL indicates end of [RFC-2822] group; if
              ; non-NIL and addr-host is NIL, holds
              ; [RFC-2822] group name.
              ; Otherwise, holds [RFC-2822] local-part
              ; after removing [RFC-2822] quoting
addr-name     = nstring
              ; If non-NIL, holds phrase from [RFC-2822]
              ; mailbox after removing [RFC-2822] quoting
append       = "APPEND" SP mailbox [SP flag-list] [SP date-time] SP
              literal
astring      = 1*ASTRING-CHAR / string
ASTRING-CHAR = ATOM-CHAR / resp-specials
atom         = 1*ATOM-CHAR
ATOM-CHAR    = <any CHAR except atom-specials>
atom-specials = "(" / ")" / "[" / SP / CTL / list-wildcards /
              quoted-specials / resp-specials
authenticate  = "AUTHENTICATE" SP auth-type *(CRLF base64)
auth-type    = atom
              ; Defined by [SASL]
base64       = *(base64-char) [base64-terminal]
base64-char  = ALPHA / DIGIT / "+" / "/"
              ; Case-sensitive
base64-terminal = (2base64-char ==") / (3base64-char ==")
body         = "(" body-type-lpart / body-type-mpart ")"
body-extension = nstring / number /
              ("body-extension" *(SP body-extension))
              ; Future expansion. Client implementations
              ; MUST accept body-extension fields. Server
              ; implementations MUST NOT generate
              ; body-extension fields except as defined by
              ; future standard or standards-track
              ; revisions of this specification.
body-ext-lpart = body-fld-md5 [SP body-fld-dsp [SP body-fld-lang
              [SP body-fld-loc *(SP body-extension)]]
              ; MUST NOT be returned on non-extensible
              ; "BODY" fetch
body-ext-mpart = body-fld-param [SP body-fld-dsp [SP body-fld-lang
              [SP body-fld-loc *(SP body-extension)]]
              ; MUST NOT be returned on non-extensible
              ; "BODY" fetch
body-fields   = body-fld-param SP body-fld-id SP body-fld-desc SP
              body-fld-rnc SP body-fld-octets
body-fld-desc = nstring
body-fld-dsp  = "(" string SP body-fld-param ")" / nil
body-fld-enc  = DQUOTE ("7BIT" / "8BIT" / "BINARY" / "BASE64" /
              "QUOTED-PRINTABLE") DQUOTE / string
body-fld-id   = nstring
body-fld-lang  = nstring / "(" string *(SP string))"
body-fld-loc  = nstring
body-fld-lines = number
body-fld-md5  = nstring
body-fld-octets = number
body-fld-param = "(" string SP string *(SP string SP string))" / nil
body-type-lpart = (body-type-basic / body-type-msg / body-type-text)
              [SP body-ext-lpart]
body-type-basic = media-basic SP body-fields
              ; MESSAGE subtype MUST NOT be "RFC822"
body-type-mpart = 1*body SP media-subtype
              [SP body-ext-mpart]
body-type-msg  = media-message SP body-fields SP envelope
              SP body SP body-fld-lines
body-type-text = media-text SP body-fields SP body-fld-lines
capability    = ("AUTH=" auth-type) / atom
              ; New Capabilities MUST begin with "X" or be
              ; registered with IANA as standard or
              ; standards-track
capability-data = "CAPABILITY" *(SP capability) SP "IMAP4rev1"
              *(SP capability)
              ; Servers MUST implement the STARTTLS, AUTH=PLAIN,
              ; and LOGIN/REGISTER capabilities
              ; Servers which offer RFC 1730 compatibility MUST
              ; list "IMAP4" as the first capability.
CHAR8        = %x01-ff
              ; any OCTET except NUL, %x00
command      = tag SP (command-any / command-auth / command-nonauth /
              command-select) CRLF
              ; Modal based on state
command-any   = "CAPABILITY" / "LOGOUT" / "NOOP" / x-command
              ; Valid in all states
command-auth  = append / create / delete / examine / list / lsub /
              rename / select / status / subscribe / unsubscribe
              ; Valid only in Authenticated or Selected state
command-nonauth = login / authenticate / "STARTTLS"
              ; Valid only when in Not Authenticated state
command-select = "CHECK" / "CLOSE" / "EXPUNGE" / copy / fetch / store /
              uid / search
              ; Valid only when in Selected state
continue-req  = "+" SP (resp-text / base64) CRLF
copy         = "COPY" SP sequence-set SP mailbox
create       = "CREATE" SP mailbox
              ; Use of INBOX gives a NO error
date         = date-text / DQUOTE date-text DQUOTE
date-day     = 1*DIGIT
              ; Day of month
date-day-fixed = (SP DIGIT) / 2DIGIT
              ; Fixed-format version of date-day
date-month   = "Jan" / "Feb" / "Mar" / "Apr" / "May" / "Jun" /
              "Jul" / "Aug" / "Sep" / "Oct" / "Nov" / "Dec"
date-text    = date-day -- date-month -- date-year
date-year    = 4DIGIT
date-time    = DQUOTE date-day-fixed -- date-month -- date-year
              SP SP zone DQUOTE
delete       = "DELETE" SP mailbox
              ; Use of INBOX gives a NO error
digit-nz     = %x31-39
              ; 1-9
envelope    = "(" env-date SP env-subject SP env-from SP
              env-sender SP env-reply-to SP env-to SP env-cc SP
              env-bcc SP env-in-reply-to SP env-message-id ")"
env-bcc      = "(" 1*address ")" / nil
env-cc       = "(" 1*address ")" / nil
env-date     = nstring
env-from     = "(" 1*address ")" / nil
env-in-reply-to = nstring
env-message-id = nstring
env-reply-to = "(" 1*address ")" / nil
env-sender   = "(" 1*address ")" / nil
env-subject  = nstring
              ; The Universal Time zone is "+0000".
env-to       = "(" 1*address ")" / nil
examine     = "EXAMINE" SP mailbox
fetch       = "FETCH" SP sequence-set SP ("ALL" / "FULL" / "FAST" /
              fetch-att / "(" fetch-att *(SP fetch-att))")
fetch-att   = "ENVELOPE" / "FLAGS" / "INTERNALDATE" /
              "RFC822" [".HEADER" / ".SIZE" / ".TEXT"] /
              "BODY" [".STRUCTURE"] / "UID" /
              "BODY.PEEK" section [<" number .">] /
              "UID.PEEK" section [<" number .">]
flag        = "\Answered" / "\Flagged" / "\Deleted" /
              "\Seen" / "\Draft" / flag-keyword / flag-extension
              ; Does not include "Recent"
flag-extension = "\ atom
              ; Future expansion. Client implementations
              ; MUST accept flag-extension flags. Server
              ; implementations MUST NOT generate
              ; flag-extension flags except as defined by
              ; future standard or standards-track
              ; revisions of this specification.
flag-fetch   = flag / "\Recent"
flag-keyword = atom
flag-list    = "(" [flag *(SP flag)] ")"
flag-perm   = flag / "\*"
greeting   = "*" SP (resp-cond-auth / resp-cond-bye) CRLF
header-fld-name = astring
header-list  = "(" header-fld-name *(SP header-fld-name)) "
list        = "LIST" SP mailbox SP list-mailbox
list-mailbox = 1*list-char / string
list-wildcards = "* / *"
literal     = "(" number ")"
              ; Number represents the number of CHAR8s
login       = "LOGIN" SP userid SP password
lsub        = "LSUB" SP mailbox SP list-mailbox
mailbox    = "INBOX" / astring
              ; INBOX is case-insensitive. All case variants of
              ; INBOX (e.g., "InBox") MUST be interpreted as INBOX
              ; as an astring. An astring which consists of
              ; the zero-inclusive sequence "I" "N" "B" "O" "X"
              ; is considered to be INBOX and not an astring.
              ; Refer to section 5.1 for further
              ; semantic details of mailbox names.
mailbox-data = "FLAGS" SP flag-list / "LIST" SP mailbox-list /
              "LSUB" SP mailbox-list / "SEARCH" *(SP nz-number) /
              "STATUS" SP mailbox *(status-att-list) /
              number SP "EXISTS" / number SP "RECENT"
mailbox-list = "(" [mbx-list-flags] ")" SP
              (DQUOTE QUOTED-CHAR DQUOTE / nil) SP mailbox
mbx-list-flags = *(mbx-list-oflag SP mbx-list-sflag
              *SP mbx-list-oflag) /
              mbx-list-oflag *(SP mbx-list-oflag)
mbx-list-oflag = "\Noinferiors" / flag-extension
              ; Other flags; multiple possible per LIST response
mbx-list-sflag = "\Noselect" / "\Marked" / "\Unmarked"
              ; Selectability flags; only one per LIST response
media-basic  = (DQUOTE ("APPLICATION" / "AUDIO" / "IMAGE" /
              "MESSAGE" / "VIDEO") DQUOTE) / string) SP
media-subtype = string
              ; Defined in [MIME-INT]
media-message = DQUOTE "MESSAGE" DQUOTE SP DQUOTE "RFC822" DQUOTE
              ; Defined in [MIME-INT]
media-subtype = string
              ; Defined in [MIME-INT]
media-text   = DQUOTE "TEXT" DQUOTE SP media-subtype
              ; Defined in [MIME-INT]
message-data = nz-number SP ("EXPUNGE" / ("FETCH" SP msg-att))
msg-att     = "(" (msg-att-dynamic / msg-att-static)
              *(SP (msg-att-dynamic / msg-att-static)) ")"
msg-att-dynamic = "FLAGS" SP "(" flag-fetch *(SP flag-fetch)) "
              ; MAY change for a message
msg-att-static = "ENVELOPE" SP envelope / "INTERNALDATE" SP date-time /
              "RFC822" [".HEADER" / ".TEXT"] SP nz-string /
              "RFC822.SIZE" SP number /
              "BODY" [".STRUCTURE"] SP body /
              "BODY" section [<" number .">] SP nz-string /
              "UID" SP uniqueid
              ; MUST NOT change for a message
nil         = "NIL"
nstring     = string / nil
number     = 1*DIGIT
              ; Unsigned 32-bit integer
              ; (0 <= n < 4,294,967,296)
nz-number   = digit-nz *DIGIT
              ; Non-zero unsigned 32-bit integer
              ; (0 < n < 4,294,967,296)
password   = astring
quoted     = DQUOTE *QUOTED-CHAR DQUOTE
QUOTED-CHAR = <any TEXT-CHAR except quoted-specials> /
              "\ quoted-specials
quoted-specials = DQUOTE / "\*
rename     = "RENAME" SP mailbox SP mailbox
              ; Use of INBOX as a destination gives a NO error
response   = *(continue-req / response-data) response-done
response-data = ** SP (resp-cond-state / resp-cond-bye /
              mailbox-data / message-data / capability-data) CRLF
response-done = response-tagged / response-fatal
response-fatal = ** SP resp-cond-bye CRLF
              ; Server closes connection immediately
response-tagged = tag SP resp-cond-state CRLF
resp-cond-auth = "OK" / "PREAUTH" SP resp-text
              ; Authentication condition
resp-cond-bye = "BYE" SP resp-text
resp-cond-state = "OK" / "NO" / "BAD" SP resp-text
              ; Status condition
resp-specials = ""
resp-text   = "[!]" resp-text-code "] SP text
resp-text-code = "ALERT" /
              "BADCHARSET" SP "(" astring *(SP astring) ") "
              / capability-data / "PARSE" /
              "PERMANTITLEG" SP "("
              [flag-perm *(SP flag-perm)] ") "
              / "READ-ONLY" / "READ-WRITE" / "TRYCREATE" /
              "UIDNEXT" SP nz-number / "UIDVALIDITY" SP nz-number /
              "UNSEEN" SP nz-number
              atom [SP 1*<any TEXT-CHAR except ">"]
search     = "SEARCH" [SP "CHARSET" SP astring] 1*(SP search-key)
              ; CHARSET argument to MUST be registered with IANA
search-key  = "ALL" / "ANSWERED" / "BCC" SP astring /
              "BEFORE" SP date / "BODY" SP astring /
              "CC" SP astring / "DELETED" / "FLAGGED" /
              "FROM" SP astring / "KEYWORD" SP flag-keyword /
              "NEW" / "OLD" / "ON" SP date / "RECENT" / "SEEN" /
              "SINCE" SP date / "SUBJECT" SP astring /
              "TEXT" SP astring / "TO" SP astring /
              "UNANSWERED" / "UNDELETED" / "UNFLAGGED" /
              "UNKEYWORD" SP flag-keyword / "UNSEEN" /
              ; Above this line were in [IMAP2]
              "DRAFT" / "HEADER" SP header-fld-name SP astring /
              "LARGER" SP number / "NOT" SP search-key /
              "OR" SP search-key SP search-key /
              "SENTBEFORE" SP date / "SENTON" SP date /
              "SENTSINCE" SP date / "SMALLER" SP number /
              "UID" SP sequence-set / "UNDRAFT" / sequence-set /
              (" search-key *(SP search-key) ")
section     = "[" [section-spec] "]"
section-mgtext = "HEADER" / "HEADER.FIELDS" [".NOT"] SP header-list /
              "TEXT"
              ; top-level or MESSAGE/RFC822 part
section-part = nz-number "(" nz-number)
              ; body part nesting
section-spec = section-mgtext / (section-part ["." section-text])
section-text = section-mgtext / "MIME"
              ; text other than actual body part (headers, etc.)
select      = "SELECT" SP mailbox
seq-number  = nz-number / "*"
              ; message sequence number (COPY, FETCH, STORE
              ; commands) or unique identifier (UID COPY,
              ; UID FETCH, UID STORE commands).
              ; * represents the largest number in use. In
              ; the case of message sequence numbers, it is
              ; the size of the message. In the case of
              ; unique identifiers of the last message in the
              ; mailbox or, if the mailbox is empty, the
              ; mailbox's current UIDNEXT value.
              ; The server should respond with a tagged BAD
              ; response to a command that uses a message
              ; sequence number greater than the number of
              ; messages in the selected mailbox. This
              ; includes "*" if the selected mailbox is empty.
seq-range   = seq-number ":" seq-number
              ; two seq-number values and all values between
              ; these two regardless of order.
              ; Example: 2:4 and 4:2 are equivalent and indicate
              ; values 2, 3, and 4.
              ; Example: a unique identifier sequence range of
              ; 3291* includes the UID of the last message in
              ; the mailbox, even if that value is less than 3291.
sequence-set = (seq-number / seq-range) * ","
              ; set of seq-number values, regardless of order.
              ; Service with coalesce overlaps and/or execute the
              ; sequence in any order.
              ; Example: a message sequence number set of
              ; 2,4,7,9,12* for a mailbox with 15 messages is
              ; equivalent to 2,4,5,6,7,9,12,13,14,15
              ; Example: a message sequence number set of *14,5:7
              ; for a mailbox with 10 messages is equivalent to
              ; 10,9,8,7,6,5,4,5,6,7 and MAY be reordered and
              ; overlap coalesced to be 4,5,6,7,8,9,10.
status      = "STATUS" SP mailbox SP
              ; (" status-att *(SP status-att) ")
status-att  = "MESSAGES" / "RECENT" / "UIDNEXT" / "UIDVALIDITY" /
              "UNSEEN"
status-att-list = status-att SP number *(SP status-att SP number)
store       = "STORE" SP sequence-set SP store-att-flags
store-att-flags = ("+" / "-") "FLAGS" ["SILENT"]
              ; flag-list / (flag *(SP flag))
string     = quoted / literal
subscribe  = "SUBSCRIBE" SP mailbox
tag         = 1*any ASTRING-CHAR except "+>"
text        = 1*TEXT-CHAR
TEXT-CHAR   = <any CHAR except CR and LF>
time       = 2DIGIT ":" 2DIGIT
              ; Hours minutes seconds
uid         = "UID" SP (copy / fetch / search / store)
              ; Unique identifiers used instead of message
              ; sequence numbers
uniqueid   = nz-number
              ; Strictly ascending
unsubscribe = "UNSUBSCRIBE" SP mailbox
userid     = astring
x-command   = "X" atom <experimental command arguments>
zone       = ("+" / "-") 4DIGIT
              ; Signed four-digit value of hhmm representing
              ; hours and minutes east of Greenwich (that is,
              ; the amount that the given time differs from
              ; Universal Time). Subtracting the timezone
              ; from the given time will give the UT form.

```

# **RFC 2812 (IRC)**



The Augmented BNF representation for this is:

```
message      = [ ":" prefix SPACE ] command [ params ] crlf
prefix       =servername / ( nickname [ [ !" user ] "@" host ] )
command     = 1*letter / 3digit
params      = *14( SPACE middle ) [ SPACE ":" trailing ]
              =/ 14( SPACE middle ) [ SPACE [ ":" ] trailing ]

nospcrlfcl = %x01-09 / %x0B-0C / %x0E-1F / %x21-39 / %x3B-FF
              ; any octet except NUL, CR, LF, " " and ":" 
middle      = nospcrlfcl *( ":" / nospcrlfcl )
trailing   = *( ":" / " " / nospcrlfcl )

SPACE       = %x20          ; space character
crlf        = %x0D %x0A      ; "carriage return" "linefeed"

target      = nickname / server
msgtarget   = msgto *( "," msgto )
msgto       = channel / ( user [ "%" host ] "@" servername )
msgto       =/ ( user "%" host ) / targetmask
msgto       =/ nickname / ( nickname !" user "@" host )
channel    = ( "#" / "+" / ( "!" channelid ) / "&" ) chanstring
              [ ":" chanstring ]
servername = hostname
host        = hostname / hostaddr
hostname   = shortname *( "." shortname )
shortname  = ( letter / digit ) *( letter / digit / "-" )
              *( letter / digit )
              ; as specified in RFC 1123 [HNAME]
hostaddr   = ip4addr / ip6addr
ip4addr    = 1*3digit "." 1*3digit "." 1*3digit "." 1*3digit
ip6addr    = 1*hexdigit 7( ":" 1*hexdigit )
ip6addr    =/ "0:0:0:0:0:" ( "0" / "FFFF" ) ":" ip4addr
nickname   = ( letter / special ) *8( letter / digit / special / "-" )
targetmask = ( "$" / "#" ) mask
              ; see details on allowed masks in section 3.3.1
chanstring = %x01-07 / %x08-09 / %x0B-0C / %x0E-1F / %x21-2B
chanstring =/ %x2D-39 / %x3B-FF
              ; any octet except NUL, BELL, CR, LF, " ", ",", " and ":" 
channelid  = 5( %x41-5A / digit ) ; 5( A-Z / 0-9 )

user        = 1*( %x01-09 / %x0B-0C / %x0E-1F / %x21-3F / %x41-FF )
              ; any octet except NUL, CR, LF, " " and "@"
key         = 1*23( %x01-05 / %x07-08 / %x0C / %x0E-1F / %x21-7F )
              ; any 7-bit US_ASCII character,
              ; except NUL, CR, LF, FF, h/v TABs, and " "
letter      = %x41-5A / %x61-7A ; A-Z / a-z
digit       = %x30-39 ; 0-9
hexdigit   = digit / "A" / "B" / "C" / "D" / "E" / "F"
special     = %x5B-60 / %x7B-7D
              ; "[", "]", "\", "^", "_", "^{", "|", "}"
```

**Efficient parsing techniques exist.**



LALR( $k$ )

Earley

LL( $k$ )

Operator

GLR

precedence

SLR

CYK

LR( $k$ )

Combinators

PEG

packrat

# Parsing tools abound.



**Yacc**

**Parsec**

**ANTLR**

**NLTK**

**Happy**

**CUPS**

**Flex**

**Bison**

**PLY**

**Ragel**

# State of the art?

\*buff++

# Apache

**2,179 lines of C**

# lighttpd

|,2| | lines of C

# freenode IRCD

> 2000 lines of C

# Courier IMAP

**2,633 lines of C**

# Result?

**CVE-ID****CVE-2004-0786**

Learn more at National Vulnerability Database

- Severity Rating
- Fix Information
- Vulnerable Software Versions

**Description**

The Internet Relay Chat (irc) protocol

**References****Note:****Vulnerability Details : [CVSS](#)**

mod\_access.c in lighttpd 1.4.1

Publish Date : 2007-07-22 Last Update

**IRC buffer overflow (IRC\_Daemon)**[Collapse All](#) [Expand All](#)[Click here if you have problems reading this page](#)**About this signature**

- [Get this signature](#)
- [Update this signature](#)
- [Mark this signature as useful](#)
- [Report this signature as useless](#)

**Apache 1.3.37 has a buffer overflow vulnerability**

**From:** "Matias Soler" <gomezz@users.sourceforge.net>  
**Date:** Tue, 2 Jan 2007 17:35:31 +0100

**Synopsis:** Apache 1.3.37 has a buffer overflow vulnerability.  
**Version:** 1.3.37 (latest)

**Product:** Apache**Component:** htpasswd utility**Version:** 1.3.37 (latest)**Issue:** Buffer overflow vulnerability**Platform:** Local system**Environment:** Suid root environment**Impact:** Denial of service**Target Milestone:** Details**Assigned To:** Details**Assigned:** Details

Incorrect validation on the size of user input allows to copy a string, via `strcpy`, to a fixed size buffer.

File: `htpasswd.c`, Line 421.**Keywords:** FixedInTrunk, PatchAvailable**Depends on:****Blocks:**[Show dependency tree](#)**Details****Vulnerable systems:**

- \* mIRC version 6.1 and prior

**Immune systems:**

- \* mIRC version 6.11

When mIRC is installed, it registers its own handler for URL of the type "irc". Calling "irc://irc.hackme.com" from our web browser causes mIRC to connect to the irc.hackme.com server. By inputting an overly long string to the "irc" protocol, an attacker can control the instruction pointer, thus controls the program's execution.

**Example:**

irc://[buffer]..... where's buffer &gt;998 bytes

An attacker would be able to gain access to the target system if he was able to control the instruction pointer. Hence, he can have his code executed under the current user's privilege.

be modified,

e.)

Mar 11 2004 12:00AM

Jul 12 2009 03:06AM

These issues were disclosed by the vendor.

Inter7 Courier-IMAP 2.2.1

Inter7 Courier-IMAP 2.2.0

Inter7 Courier-IMAP 2.1.2

Inter7 Courier-IMAP 2.1.1

“High-performance  
network daemon”

“High-throughput  
exploit server”

**Yet. . .**

All of them also  
use lex and yacc.

**(For config files.)**

**The tools exist.**

**They know how to use them.**

**So, why don't they?**

# Compilers

## Interpreters

Syntax highlighting

Natural language

Network protocols

Command lines

Configuration files

Serialization

Query languages

API design

# Compilers

## Interpreters

Syntax highlighting

Natural language

Network protocols

Command lines

Configuration files

Serialization

Query languages

API design



```
$ cat /etc/*
```



man.conf

networks

protocols

crontab

asl.conf

hosts

sudoers

fstab

bind zones

passwd

syslog.conf

ftpd.conf

notify.conf

ssh\_config

services

gettytab

# A tribute to false dichotomy.

# Human *versus* Machine



There is no  
tradeoff.

**XML: a compromise  
solution to a problem  
that does not exist.**

**Neither human nor  
machine-readable.**

**The tools exist.**

**So, why don't we use them?**

# Compilers

## Interpreters

Syntax highlighting

Natural language

Network protocols

Command lines

Configuration files

Serialization

Query languages

API design

# Compilers

## Interpreters

Syntax highlighting

Natural language

Network protocols

Command lines

Configuration files

Serialization

Query languages

API design

```
#include <stdio.h>
#include <stdlib.h>
```

```
b. 5M red b
622B Jan 27
519B Jan 27
```

```
typedef unsigned int count ;
t@raptor-mk2:~/talks/parsing-sta
count parens1 ; night2011derivative
unsigned int parens2 ;
```

```
#include <stdio.h>
#include <stdlib.h>
```

```
b. 5M red b
622B Jan 27
519B Jan 27
```

```
typedef unsigned int count ;
count parens1 ;
```

count

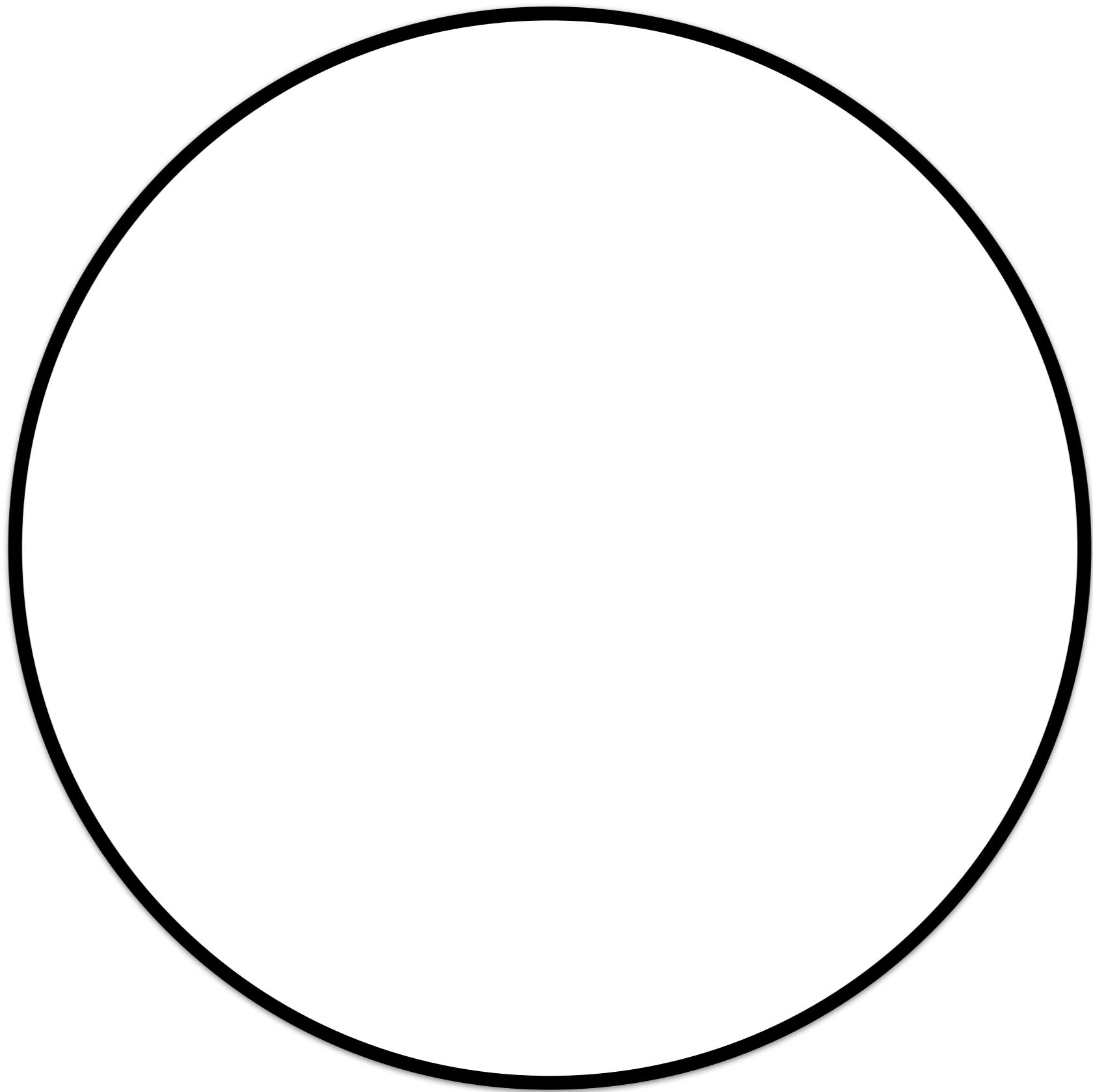
```
unsigned int parens2 ;
```

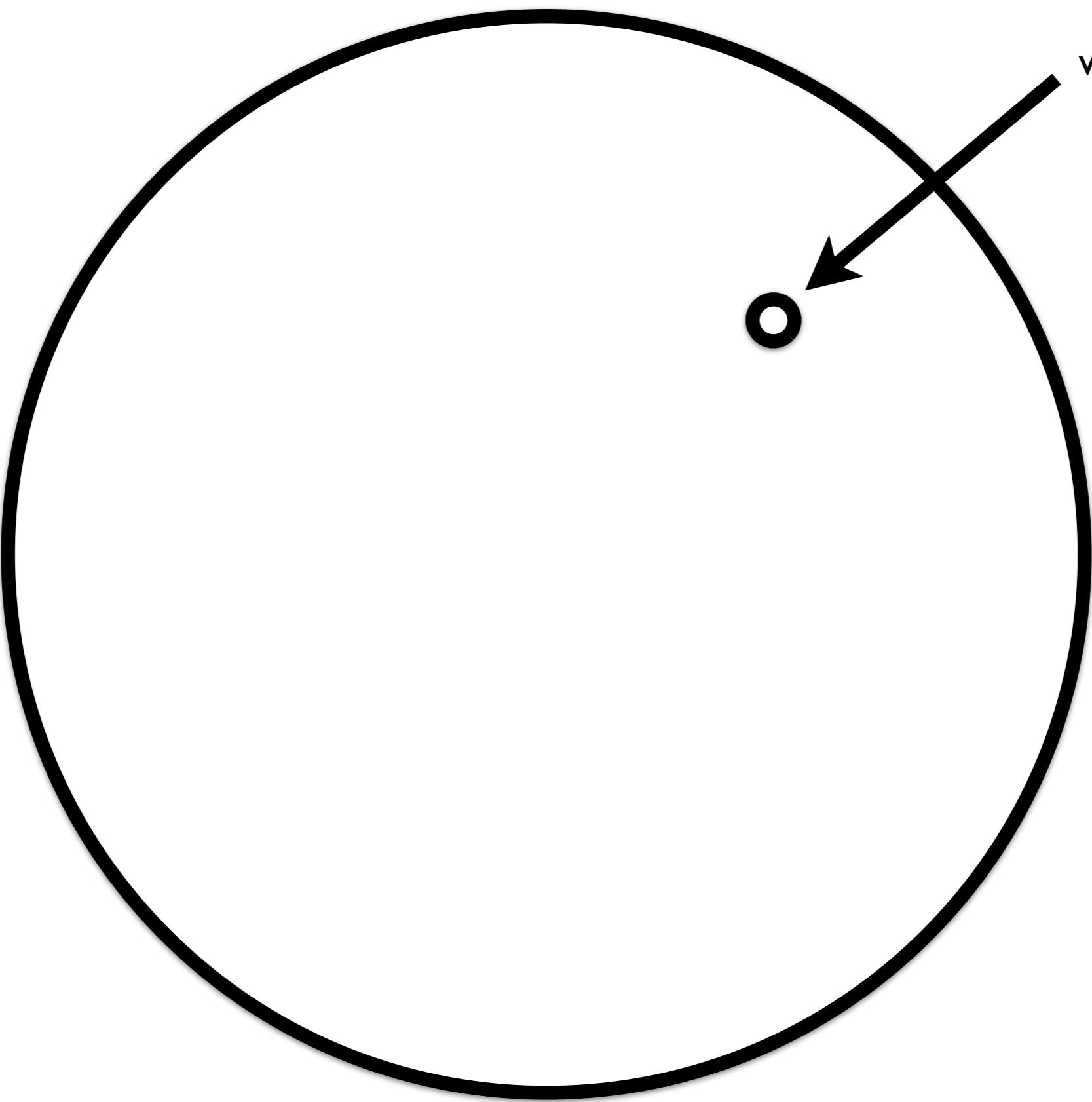
# **emacs, vi: over 30 years**



**WHY!?**







where we use it

# A list of grievances

# **Opportunity cost**

# **Transaction cost**

# **Barriers to entry**

# Opportunity cost

**What do you sacrifice  
with a parser generator?**

# Performance

## Thread-safety

## Non-blocking I/O

~~Performance~~

Thread-safety

Non-blocking I/O

**Performance**

**Thread safety**

**Non-blocking I/O**

# **What if a server blocks on I/O?**

**LOS**

# **Transaction cost**

Time to set up lex & yacc

*versus*

Time to hack it with strtok

# lex/yacc

# strtok()

# lex/yacc

Create .y yacc file

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

# strtok()

# **lex/yacc**

Create .y yacc file

Define %union lval

Define tokens/types

# **strtok()**

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

#include y.tab.h

# strtok()

# **lex/yacc**

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

#include y.tab.h

Define states

# **strtok()**

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

#include y.tab.h

Define states

Define lex rules

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

#include y.tab.h

Define states

Define lex rules

Deal with yywrap()

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

#include y.tab.h

Define states

Define lex rules

Deal with yywrap()

Call yyparse()

# strtok()

# lex/yacc

Create .y yacc file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc spec

Create .l lex file

#include y.tab.h

Define states

Define lex rules

Deal with yywrap()

Call yyparse()

# strtok()

Call strtok()

# Barriers to entry

# How does yacc work?

```
start : exp { printf("ans: %i\n", $1 ) ; }

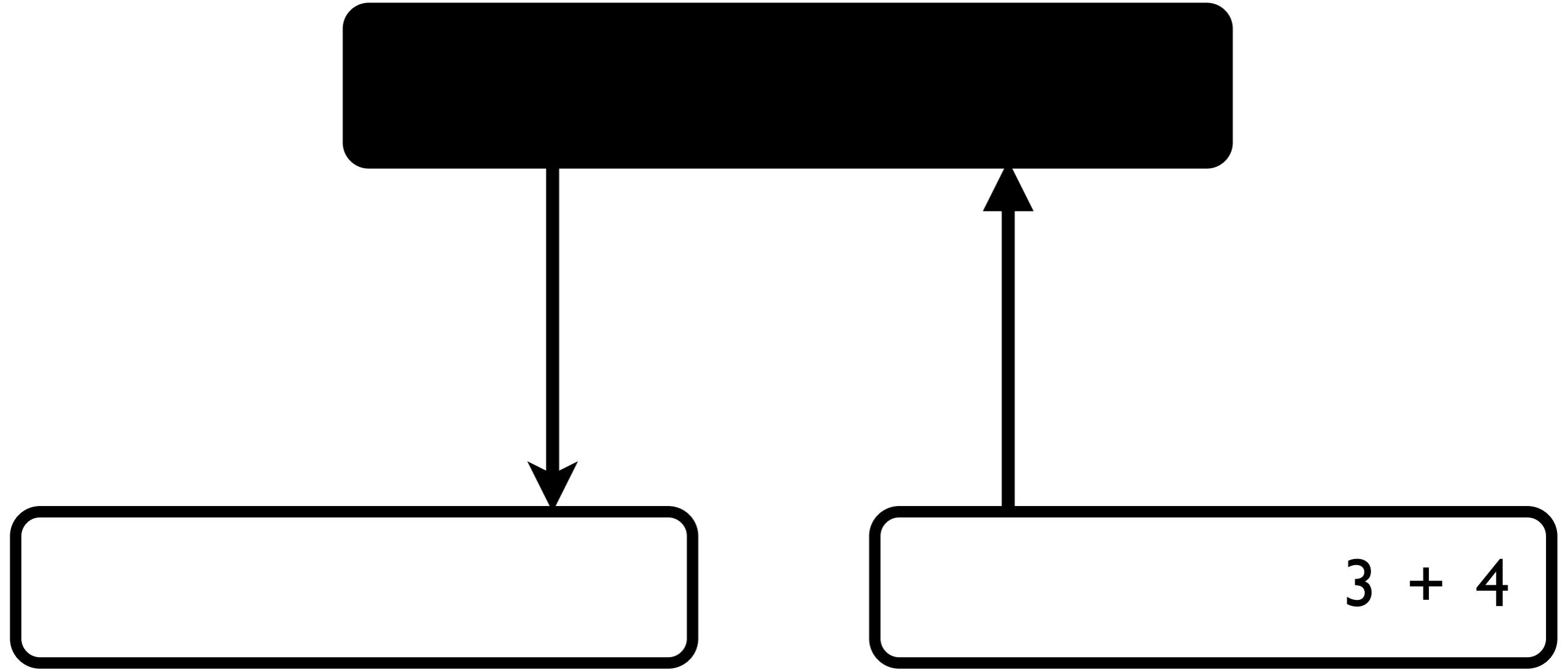
;

exp : exp '+' term { $$ = $1 + $3 ; }
| term           { $$ = $1 ; }
;

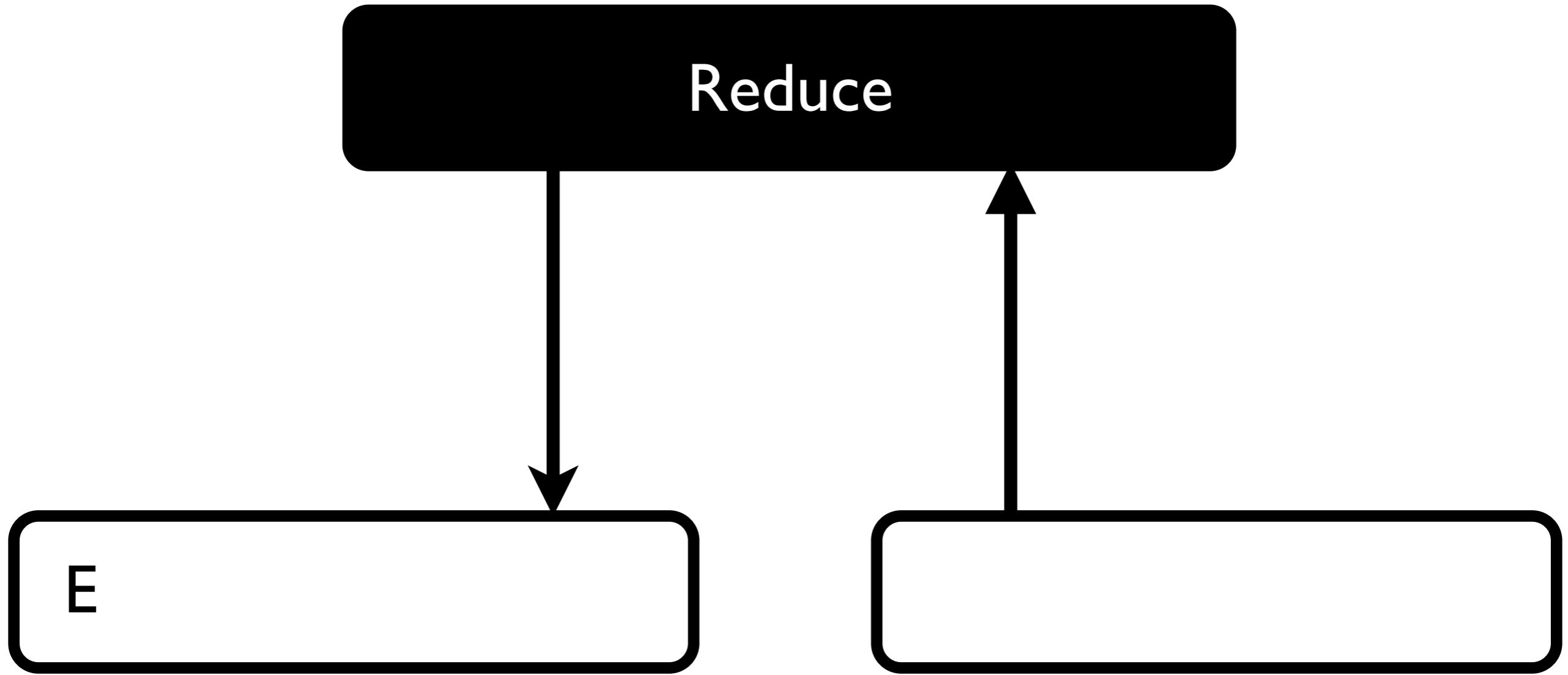
term : term '*' factor { $$ = $1 * $3 ; }
| factor         { $$ = $1 ; }
;

factor : INT { $$ = $1 ; }
| SYM { $$ = lookup($1) ; }
| '(' exp ')' { $$ = $2 ; }
;
```





Reduce



All good until...

```
$ bison -d exp.y  
exp.y: conflicts: 57 shift/  
reduce
```

```
$ bison -d exp.y  
exp.y: conflicts: 57 shift/  
reduce
```

What's in the black box?

```
$ wc -l exp.output  
5109 exp.output
```

state 301:

2 exp: exp . '+' exp   [\$end, '+', '\*', ')']  
2      | exp '+' exp .   [\$end, '+', '\*', ')']  
3      | exp . '\*' exp

'+' shift, and go to state 8  
'\*' shift, and go to state 9

'+'               [reduce using rule 2 (exp)]  
'\*'               [reduce using rule 2 (exp)]  
\$default   reduce using rule 2 (exp)

%left '+'

%left '\*'

```
$ bison -d exp.y  
exp.y: conflicts: 53 shift/  
reduce
```

**Two days later...**

```
%right VAR
%left '@'
%right FOR
%left WHILE
%nonassoc '-'
%left FUNCTION
%nonassoc RETURN
%right '?'
%left '+'
%left ','
%left '*'
%left '['
%right ']'
%left '.'
%left '{'
%right '}'
%nonassoc ';'
```

```
$ bison -d exp.y  
exp.y: conflicts: 1 shift/reduce
```

```
$ bison -d exp.y  
exp.y: conflicts: 1 shift/reduce
```

“Probably just dangling else.”

Barrier to entry:  
One must *learn* each tool.

**Barrier to entry:  
Tool must exist for language.**

How do we change the  
economics of parsing?

# What we need

- Allow non-blocking I/O
- Act as library within the language
- Handle all CFGs: parse forest
- Draw on regular expressions
- Cacheable partial parses
- Be easy to implement

(Jim, Mandelbaum, Walker, POPL 2010)

# Would be nice

- Dynamic reconfigurability
- Beyond context-free languages
- Built into the language
- Support compositionality

(Jim, Mandelbaum, Walker, POPL 2010)

# Parsing with derivatives

# Brzozowski, 1964

## Derivatives of Regular Expressions

JANUSZ A. BRZOZOWSKI

*Princeton University, Princeton, New Jersey†*

*Abstract.* Kleene's regular expressions, which can be used for describing sequential circuits, were defined using three operators (union, concatenation and iterate) on sets of sequences. Word descriptions of problems can be more easily put in the regular expression language if the language is enriched by the inclusion of other logical operations. However, in the problem of converting the regular expression description to a state diagram, the existing methods either cannot handle expressions with additional operators, or are made quite complicated by the presence of such operators. In this paper the notion of a derivative of a regular expression is introduced and the properties of derivatives are discussed. This leads, in a very natural way, to the construction of a state diagram from a regular expression containing any number of logical operators.



D C J

## I. Filter:

Keep every string starting with  $c$ .

## 2. Chop:

Remove  $c$  from the start of each.

*D*  
f

foo frak bar

*D*  
f

foo

frak

*D*<sub>f</sub>

oo

rak

$$D_c L = \{w : cw \in L\}$$

$$cw \in L \text{ iff } w \in D_c(L).$$

# Recognition algorithm

- Derive with respect to each character.
- Does the derived language contain  $\varepsilon$ ?

**foo** ∈ (foo)\*

oo ∈ f(foo)\*

oo ∈  $D_f(\text{oo})^*$

$$oo \in oo(\text{foo})^*$$

$$oo \in oo(foo)^*$$

$$o \in o(\text{foo})^*$$

$$\varepsilon \in (\text{foo})^*$$

$$\varepsilon \in (\text{foo})^*$$

# Deriving atomic languages

$$\epsilon \equiv \{ "", "" \}$$

$$c \equiv \{ c \}$$

$$\emptyset \equiv \{ \}$$

```
(define-struct Ø {})  
(define-struct ε {})  
(define-struct token {value})
```

$$D_c \emptyset =$$

$$D_c \emptyset = \emptyset$$

**(define (D c L)**

```
(define (D c L)
  (match L
```

```
(define (D c L)
  (match L
    [ (ø)           (ø) ]
```

$$D_c(\epsilon) =$$

$$D_c(\epsilon) = \emptyset$$

```
(define (D c L)
  (match L
    [ (ε) (ø)]
```

$$D_c\{c\} = \epsilon$$

$$D_c\{c\} = \epsilon$$

$$D_c\{c'\} = \emptyset \text{ if } c \neq c'$$

```
(define (D c L)
  (match L
    [(token a)           (cond [(eqv? a c) (ε)]
                                 [else          (Ø)] )]
    [else]))
```

# Deriving regular languages

$$L_1\cup L_2$$

$$L_1 \cdot L_2$$

$$L_1^\star$$

```
(define-struct U {this that})  
(define-struct O {left right})  
(define-struct ★ {lang})
```

$$D_c(L_1 \cup L_2)$$

$$\begin{aligned}
D_c(L_1 \cup L_2) &= \{w : cw \in L_1 \cup L_2\} \\
&= \{w : cw \in L_1 \text{ or } cw \in L_2\} \\
&= \{w : w \in D_c L_1 \text{ or } w \in D_c L_2\} \\
&= \{w : w \in D_c L_1\} \cup \{w : w \in D_c L_2\} \\
&= D_c L_1 \cup D_c L_2.
\end{aligned}$$

```
(define (D c L)
  (match L
    [(u L1 L2)      (u (D c L1)
                           (D c L2))]
```

$$D_c(L^\star) =$$

$$D_c(L^\star) = (D_c L) \cdot L^\star$$

```
(define (D c L)
  (match L
    [ (★ L1)      (◦ (D c L1) (★ L1)) ]
```

# Concatenation?

# Needs nullability operator

$$\delta(L) = \epsilon \text{ if } \epsilon \in L$$

$$\delta(L) = \emptyset \text{ if } \epsilon \notin L$$

```
(define-struct δ {lang})
```

$$D_c(\delta(L)) = \emptyset$$

```
(define (D c L)
  (match L
    [ (δ _) (ø) ]
```

$$D_c(L_1 \cdot L_2) =$$

$$D_c(L_1 \cdot L_2) = (D_cL_1 \cdot L_2)$$

$$D_c(L_1 \cdot L_2) = (D_cL_1 \cdot L_2) \cup (\delta(L_1) \cdot D_cL_2)$$

```
(define (D c L)
  (match L
    [ (o L1 L2)      (u (o (δ L1) (D c L2))
                           (o (D c L1) L2))]))
```

```

(define (D c L)
  (match L
    [ (ø)           (ø) ]
    [ (ε)           (ø) ]
    [ (token a)     (cond [(eqv? a c) (ε)]
                           [else             (ø) ] ) ]
    [ (δ _)         (ø) ]

    [ (U L1 L2)    (U (D c L1)
                       (D c L2)) ]
    [ (★ L1)        (◦ (D c L1) L) ]
    [ (◦ L1 L2)    (U (◦ (δ L1) (D c L2))
                       (◦ (D c L1) L2)) ] ) )

```

# To recognize?

# Need nullability

Need

nullability

Need to *compute* nullability

$$\delta(e) = e$$

$$\delta(c) = \emptyset$$

$$\delta(\emptyset) = \emptyset$$

$$\delta(L_1 \cup L_2) = \delta(L_1) \cup \delta(L_2)$$

$$\delta(L_1 \cdot L_2) = \delta(L_1) \cdot \delta(L_2)$$

$$\delta(L_1^\star) = \epsilon$$

```

(define (nullable? L)
  (match L
    [(\emptyset) #f]
    [(ε) #t]
    [(token _) #f]
    [(δ L1) (nullable? L1)]

    [(* _ ) #t]
    [(∪ L1 L2) (or (nullable? L1)
                      (nullable? L2))]

    [(◦ L1 L2) (and (nullable? L1)
                      (nullable? L2))]))
```

```
(define (recognizes? w L)
  (if (null? w)
    (nullable? L)
    (recognizes? (cdr w) (D (car w) L)))) )
```

# How about context-free grammars?

# context-free grammars

# Recursive regular expressions

# Problem

$$L = L \cdot x$$

$$\cup \epsilon$$

# Problem

$$D_{\mathbf{x}} L = D_{\mathbf{x}} L \cdot \mathbf{x}$$

$$\cup \epsilon$$

(D<sup>-1</sup> × L) =

$$\begin{aligned}
 (D^{-1}x L) &= (D^{-1}x (U (\circ L^{-1}x) \\
 &\quad \varepsilon)) \\
 &= (U (U (\circ (D^{-1}x L)^{-1}x) \\
 &\quad (\circ (\delta L) (D^{-1}x^{-1}x))) \\
 &\quad (D^{-1}x \varepsilon))
 \end{aligned}$$



(D<sup>-1</sup> × L) =

$$(D^{-1}x L) = (D^{-1}x (U (\circ^{-1}x L)$$

$\varepsilon))$



$$= (U (U (\circ (D^{-1}x^{-1}x) L) (D^{-1}x L)))$$

$(D^{-1}x \varepsilon))$

# Solution?

```
(define-struct Ø {})
(define-struct ε {})
(define-struct token {value})

(define-struct U {this that})
(define-struct ° {left right})
(define-struct ★ {lang})

(define-struct δ {lang})
```

```
(define-struct Ø {})
(define-struct ε {})
(define-struct token {value})

(define-lazy-struct U {this that})
(define-lazy-struct O {left right})
(define-lazy-struct ★ {lang})

(define-lazy-struct δ {lang})
```

# Problem

$$\delta(L) = \delta(L) \cdot \delta(x)$$

$$\cup \delta(e)$$

# Problem

$$\delta(L) = \delta(L) \cdot \delta(x)$$
$$\cup \delta(e)$$


# Solution?

**Fix it.**

```

(define nullable? L)
  (match L
    [ (empty)          #f]
    [ (ε)              #t]
    [ (token _)        #f]
    [ (δ L1)           (nullable? L1)]

    [ (★ _)            #t]
    [ (∪ L1 L2)         (or (nullable? L1)
                           (nullable? L2)) ]
    [ (◦ L1 L2)         (and (nullable? L1)
                           (nullable? L2)) ] ) )

```

```

(define/fix (nullable? L)
  #:bottom #f
  (match L
    [(\emptyset)           #f]
    [(ε)                  #t]
    [(token _)            #f]
    [(δ L1)               (nullable? L1)]
    [(* _ )               #t]
    [(∪ L1 L2)            (or (nullable? L1)
                                (nullable? L2)) ]
    [(◦ L1 L2)            (and (nullable? L1)
                                (nullable? L2)) ] ) )
  
```

# Final problem

# **Grammar unfolds forever**

# Solution?

# Memoize

```

(define (D c L)
  (match L
    [(∅)           (∅)]
    [(ε)           (∅)]
    [(token a)     (cond [(eqv? a c) (ε)]
                          [else          (∅)]))]
    [(δ _)         (∅)]

    [(∪ L1 L2)     (∪ (D c L1)
                      (D c L2))]

    [(★ L1)        (◦ (D c L1) L))

    [(◦ L1 L2)     (∪ (◦ (δ L1) (D c L2))
                      (◦ (D c L1) L2)))])))

```

```

(define/memoize (D c L)
  #:order [([L #:eq] [c #:equal])])
  (match L
    [(∅)           (∅)]
    [(ε)           (∅)]
    [(token a)     (cond [(eqv? a c) (ε)]
                          [else             (∅)]))]
    [(δ _)         (∅)]

    [(∪ L1 L2)     (∪ (D c L1)
                      (D c L2)))]
    [(★ L1)        (◦ (D c L1) L))]
    [(◦ L1 L2)     (∪ (◦ (δ L1) (D c L2))
                      (◦ (D c L1) L2)))])))

```

**It works!**

**(for recognition)**

# What about parsing?

$$D_c : \mathbb{L} \rightarrow \mathbb{L}$$

$$D_c : \mathbb{P}(A,T) \rightarrow \mathbb{P}(A,T)$$

$$\mathbb{P}(A,T) = A^* \rightarrow \mathcal{P}(T \times A^*)$$

```

(define/memoize (D c L)
  #:order [([L #:eq] [c #:equal])])
  (match L
    [(∅)           (∅)]
    [(ε)           (∅)]
    [(token a)     (cond [(eqv? a c) (ε)]
                          [else             (∅)])])
    [(δ _)         (∅)])

    [(∪ L1 L2)     (∪ (D c L1)
                      (D c L2))])
    [(★ L1)        (◦ (D c L1) L))]
    [(◦ L1 L2)     (∪ (◦ (δ L1) (D c L2))
                      (◦ (D c L1) L2)))]
  )

```

```

(define/memoize (D c L)
  #:order [([L #:eq] [c #:equal])])
  (match L
    [(∅) (∅)]
    [(ε _ ) (∅)]
    [(token a) (cond [(eqv? a c) (ε (set c))]
                        [else (∅))])]
    [(δ _ ) (∅)]
    [(∪ L1 L2) (∪ (D c L1)
                    (D c L2))])
    [(★ L1) (◦ (D c L1) L))]
    [(◦ L1 L2) (∪ (◦ (δ L1) (D c L2))
                    (◦ (D c L1) L2))))]
    [(→ L1 f) (→ (D c L1) f))])

```

# Computing nullability

# Computing null parses

$$\lfloor \emptyset \rfloor(\epsilon) = \{\}$$

$$\lfloor \epsilon \downarrow T \rfloor(\epsilon) = T$$

$$\lfloor \delta(p) \rfloor = \lfloor p \rfloor(\epsilon)$$

$$\lfloor p \cup q \rfloor(\epsilon) = \lfloor p \rfloor(\epsilon) \cup \lfloor q \rfloor(\epsilon)$$

$$\lfloor p \circ q \rfloor(\epsilon) = \lfloor p \rfloor(\epsilon) \times \lfloor q \rfloor(\epsilon)$$

$$\lfloor p \rightarrow f \rfloor(\epsilon) = \{f(t_1), \dots, f(t_n)\}$$

where  $\{t_1, \dots, t_n\} = \lfloor p \rfloor(\epsilon)$

$$\lfloor p^* \rfloor(\epsilon) = (\lfloor p \rfloor(\epsilon))^*$$

```

(define/fix (parse-ε p)
  #:bottom (set)
  (match p
    [(_ε S)           S]
    [(_∅)              (set)]
    [(_δ p)            (parse-ε p)]]
    [(_token _)        (set)]
    [(_★ _)            (set '())]
    [(_∪ p1 p2)        (set-union (parse-ε p1)
                                    (parse-ε p2))])
    [(_◦ p1 p2)        (for*/set ([t1 (parse-ε p1)]
                                    [t2 (parse-ε p2)])
                                (cons t1 t2)))]
    [(_→ p1 f)          (for/set ([t (parse-ε p1)])
                                (f t)))])))

```

```
(define (recognizes? w L)
  (if (null? w)
    (nullable? L)
    (recognizes? (cdr w) (D (car w) L)))))
```

```
(define (parse w L)
  (if (null? w)
      (parse-ε L)
      (parse (cdr w) (D (car w) L)))) )
```



$$\epsilon \equiv \lambda w. \{(\epsilon, w)\} \quad \mathbb{P}(A, T) = A^* \rightarrow \mathcal{P}(T \times A^*) \quad \begin{cases} D_c(c) = \epsilon \rightarrow \lambda \epsilon. c \\ D_c(c') = \emptyset \text{ if } c \neq c' \end{cases}$$

$$p \in \mathbb{P}(A, T) \quad \emptyset \equiv \lambda w. \{\} \quad [\mathbb{P}](A, T) = A^* \rightarrow \mathcal{P}(T)$$

$$\lfloor p \rfloor(w) = \{t : (t, \epsilon) \in p(w)\}$$

$$f \in X \rightarrow Y \quad w \equiv \lambda w'. \begin{cases} \{(w, w'')\} & w' = ww'' \\ \emptyset & \text{otherwise.} \end{cases}$$

$$p \in \mathbb{P}(A, X)$$

$$p \rightarrow f \in \mathbb{P}(A, Y) \quad D_c : \mathbb{L} \rightarrow \mathbb{L} \quad D_c : [\mathbb{P}](A, T) \rightarrow [\mathbb{P}](A, T) \quad \begin{matrix} & p \in \mathbb{P}(A, X) \\ p \rightarrow f = \lambda w. \{((f(x), w') : (x, w') \in p(w)\} & q \in \mathbb{P}(A, X) \end{matrix}$$

$$D_c : \mathbb{P}(A, T) \rightarrow \mathbb{P}(A, T) \quad p \cup q \in \mathbb{P}(A, X)$$

$$D_c(p) = \lambda w. p(cw) - (\lfloor p \rfloor(\epsilon) \times \{cw\}) \quad p \cup q = \lambda w. p(w) \cup q(w)$$

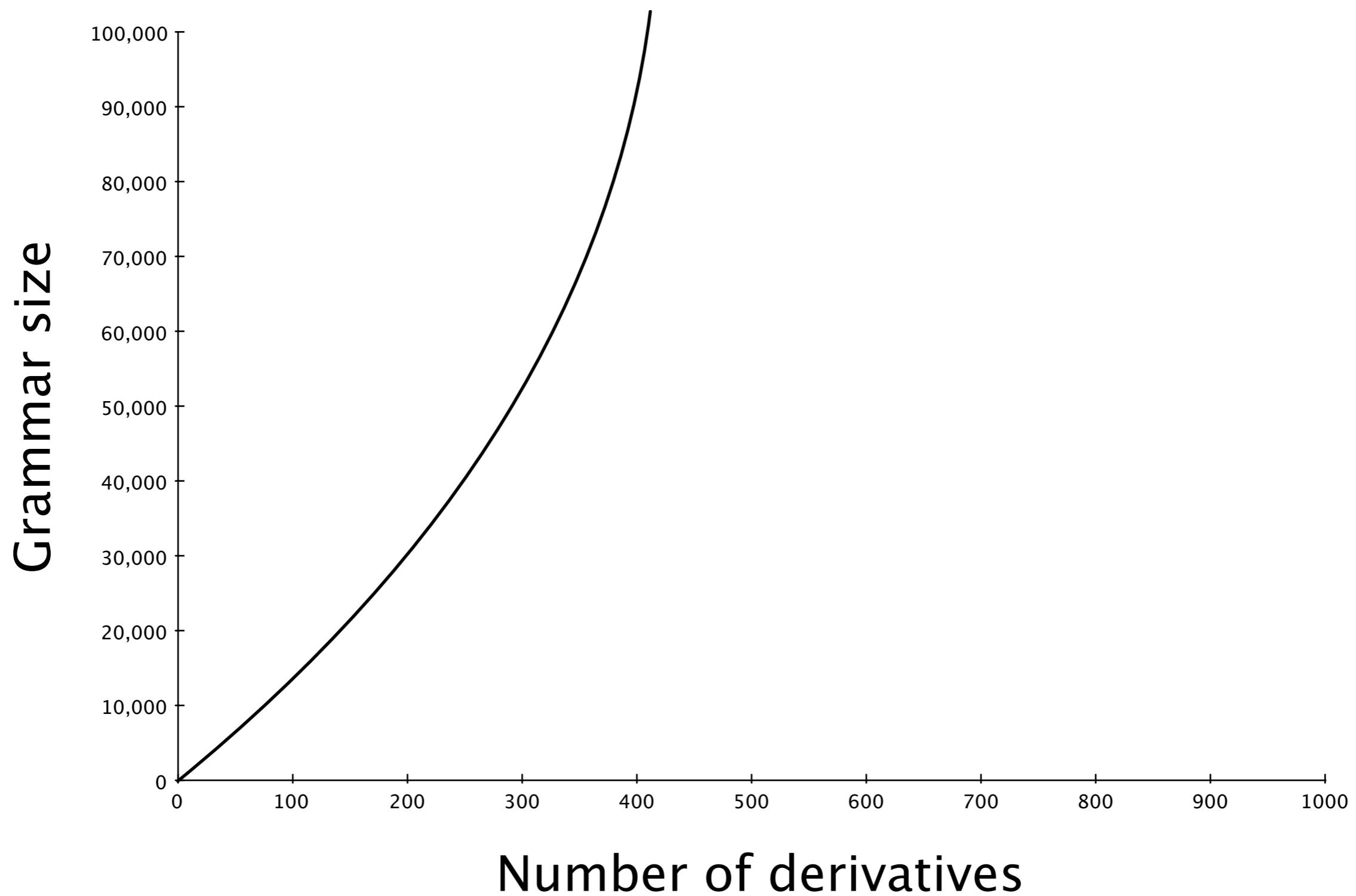
$$p(cw) = D_c(p)(w) \cup (\lfloor p \rfloor(\epsilon) \times \{cw\}) \quad D_c(p \cup q) = D_c(p) \cup D_c(q)$$

$$D_c(p \cdot q) = \begin{cases} D_c(p) \cdot q & \epsilon \notin \mathcal{L}(p) \\ D_c(p) \cdot q \cup (\epsilon \rightarrow \lambda \epsilon. \lfloor p \rfloor(\epsilon)) \cdot D_c(q) & \text{otherwise.} \end{cases} \quad D_c(p \rightarrow f) = D_c(p) \rightarrow f$$

$$p \cdot q = \lambda w. \{((x, y), w'') : (x, w') \in p(w), (y, w'') \in q(w')\}$$

# Performance?

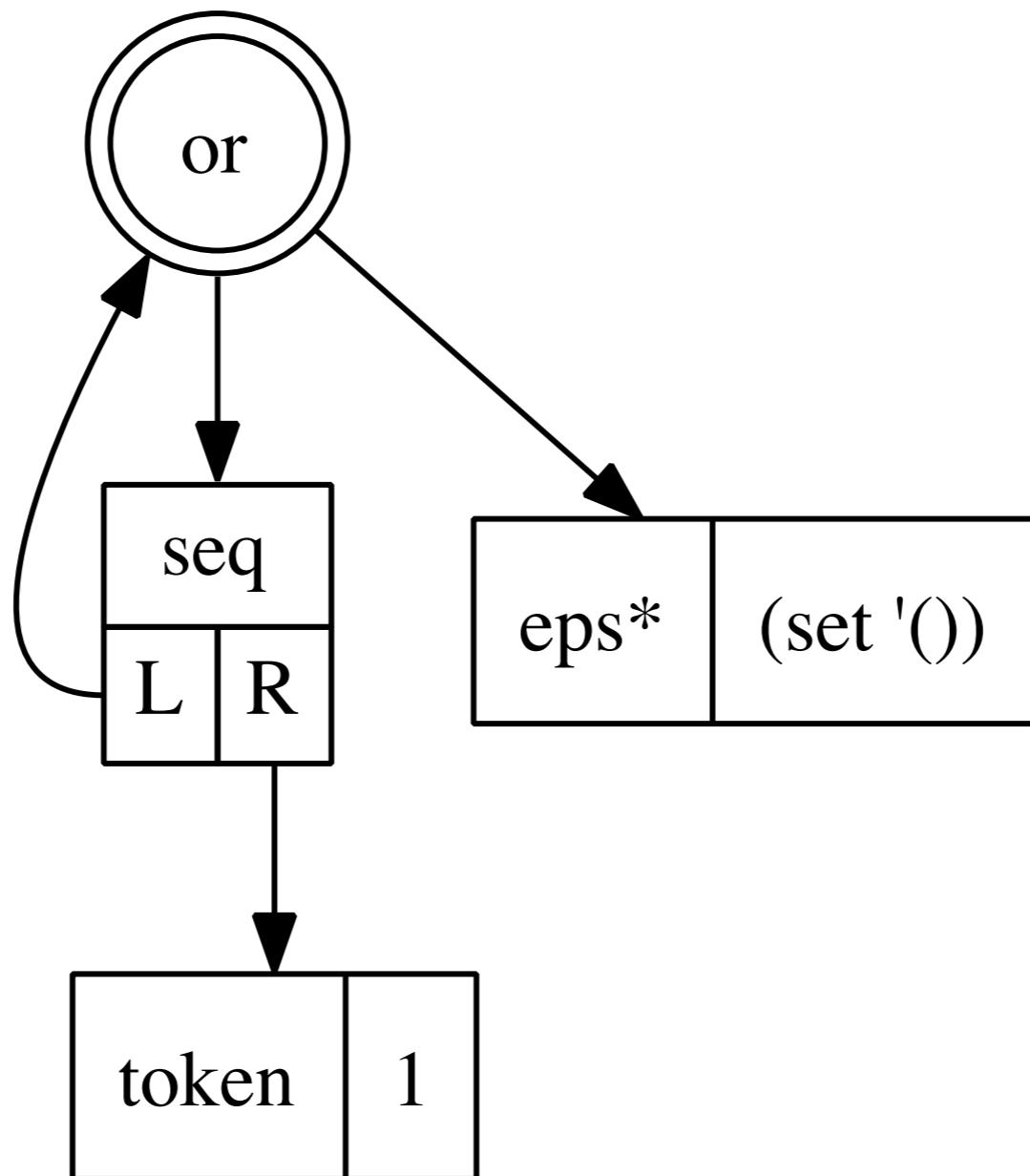
# Atrocious.

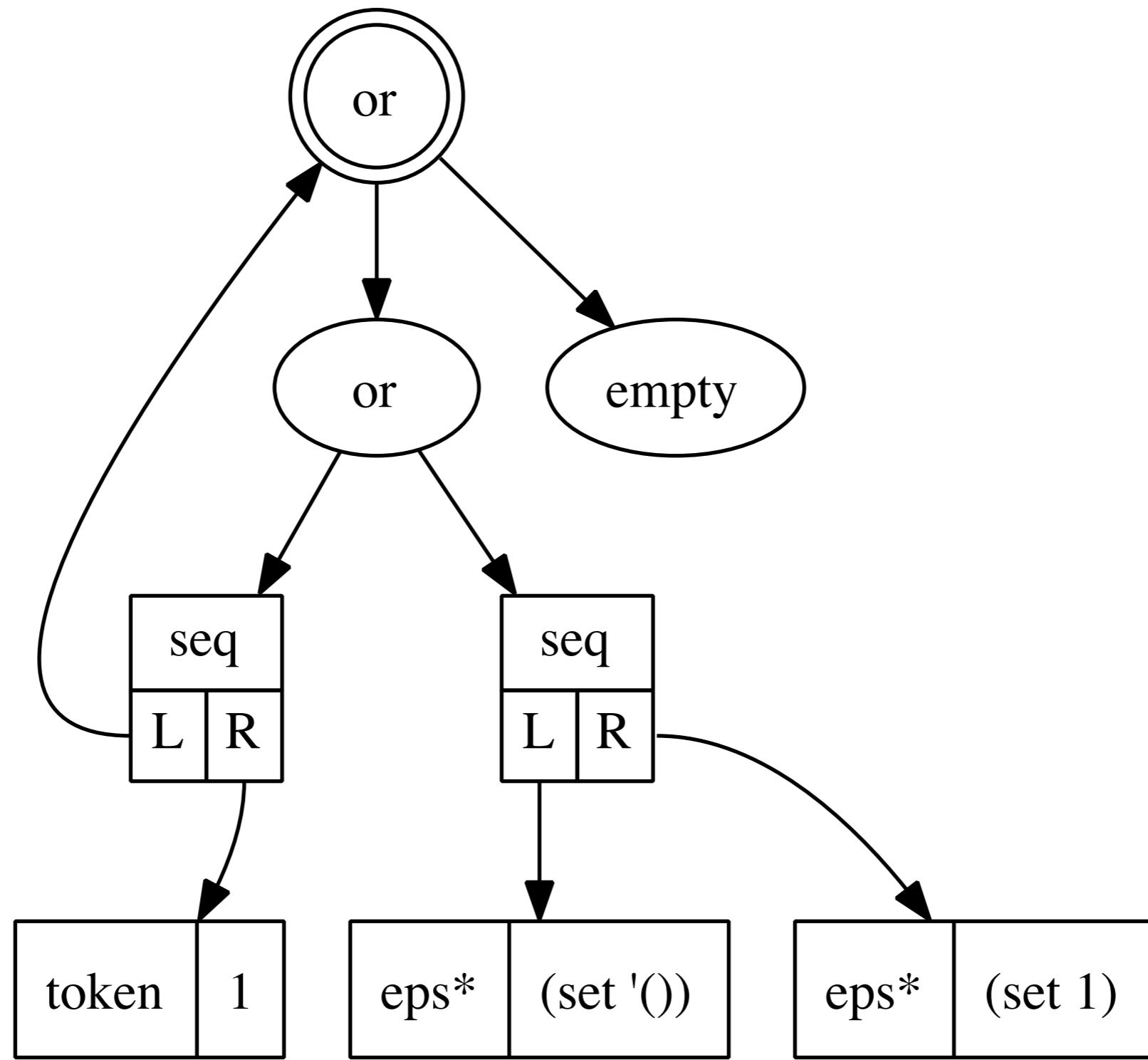


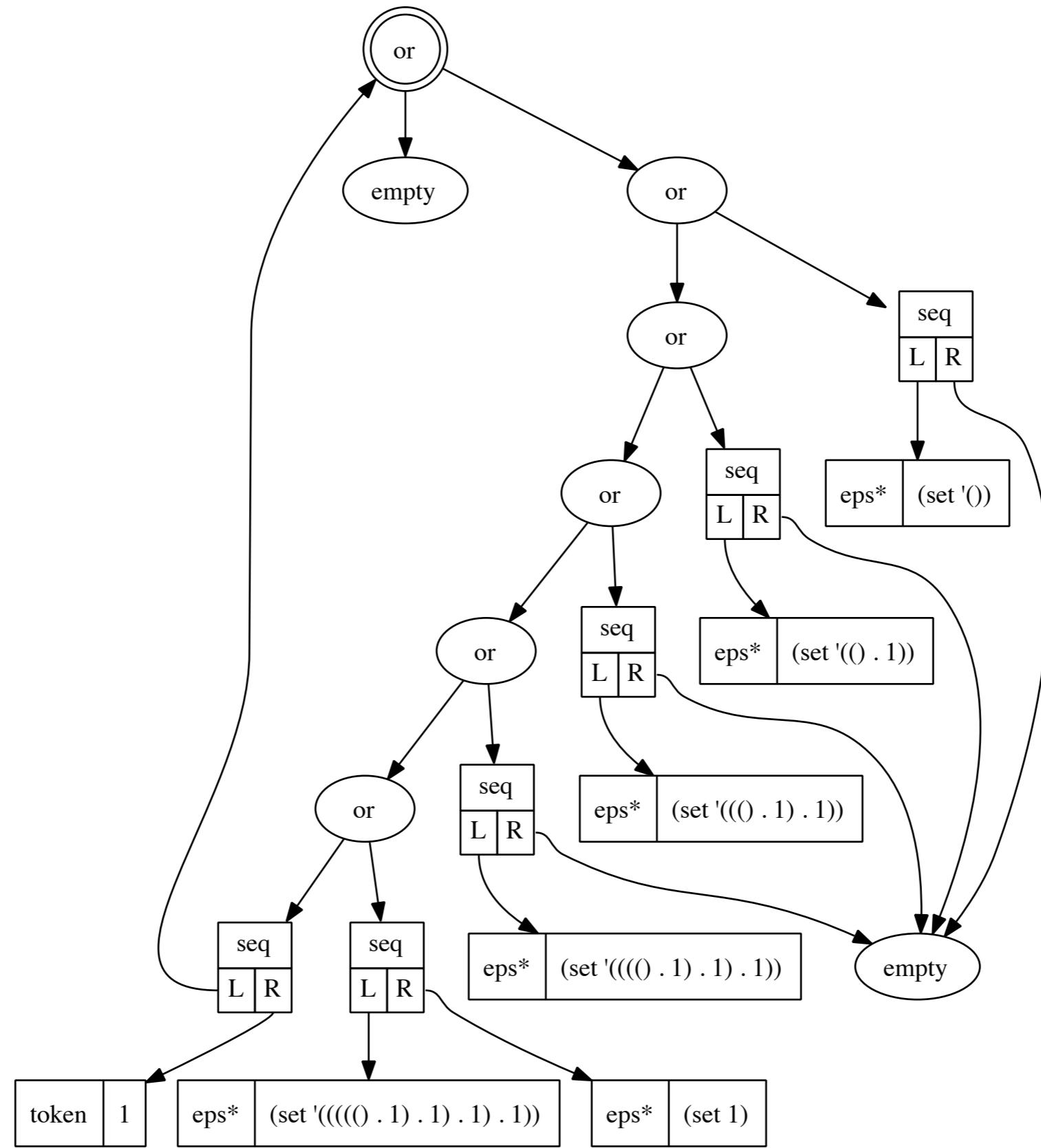
# Complexity?

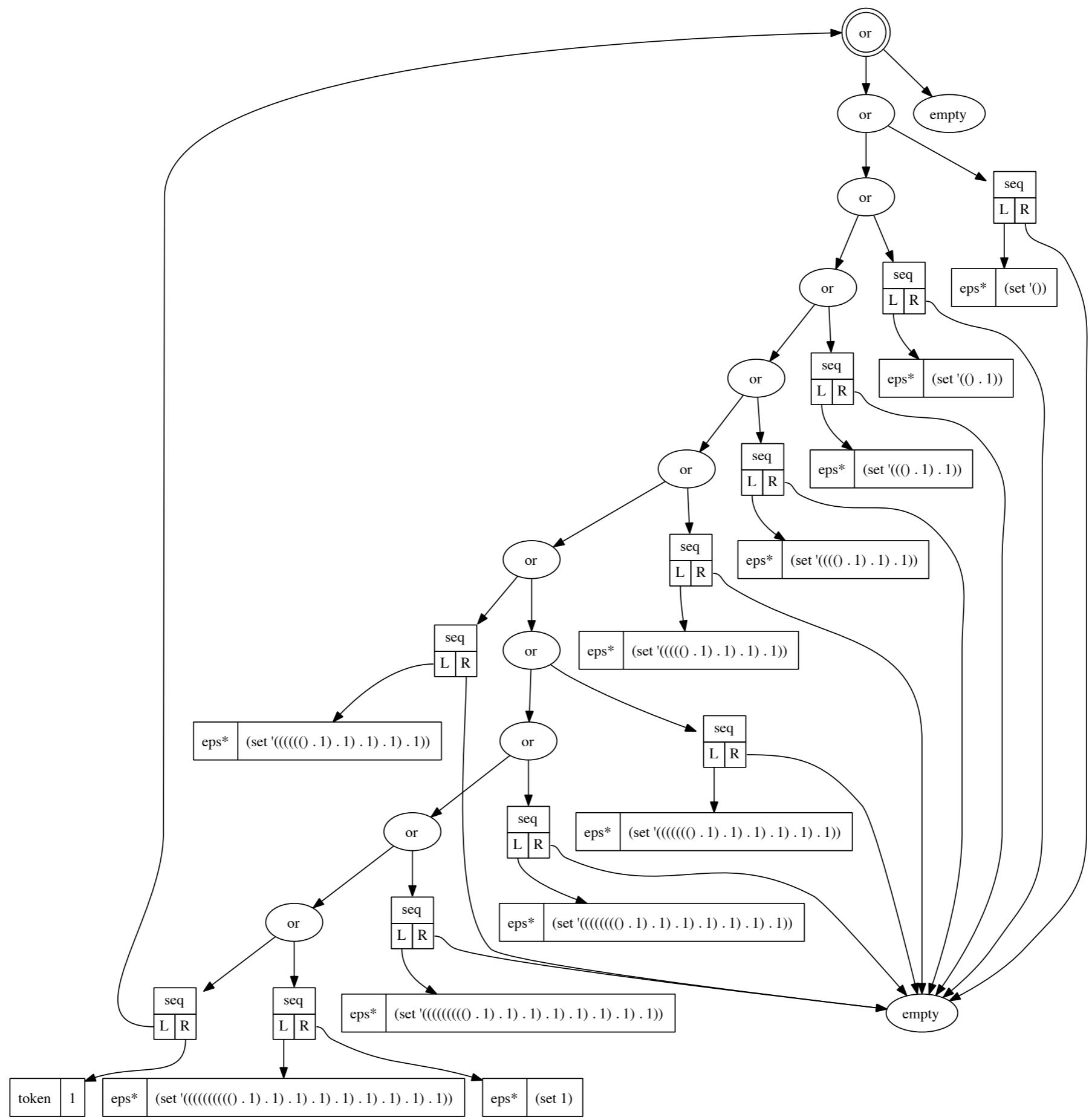
$$O(n2^n |G|)$$

```
(define OneList (lang (or (seq OneList '1)
                           (eps* (set '()))))))
```









# Needs compaction.

$$p \cdot \emptyset = \emptyset$$

$$\emptyset \circ p = p \circ \emptyset \Rightarrow \emptyset$$

$$\emptyset \cup p = p \cup \emptyset \Rightarrow p$$

$$(\epsilon \downarrow \{t_1\}) \circ p \Rightarrow p \rightarrow \lambda t_2.(t_1, t_2)$$

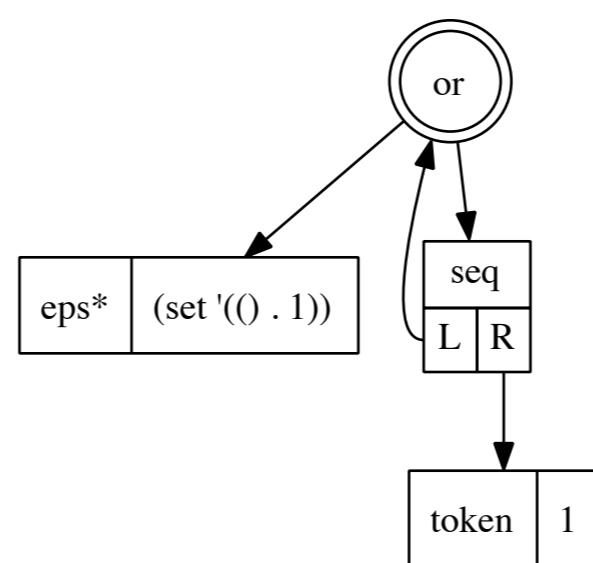
$$p \circ (\epsilon \downarrow \{t_2\}) \Rightarrow p \rightarrow \lambda t_1.(t_1, t_2)$$

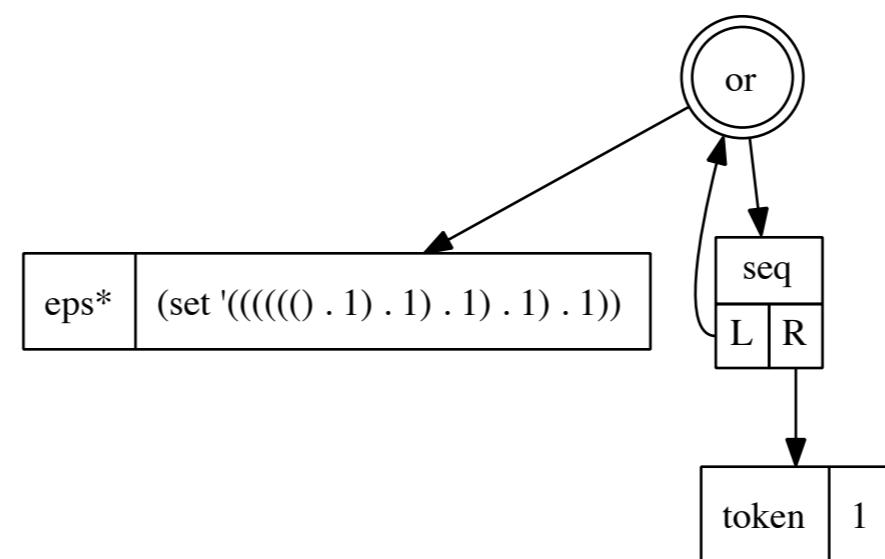
$$(\epsilon \downarrow \{t_1, \dots, t_n\}) \rightarrow f \Rightarrow \epsilon \downarrow \{f(t_1), \dots, f(t_n)\}$$

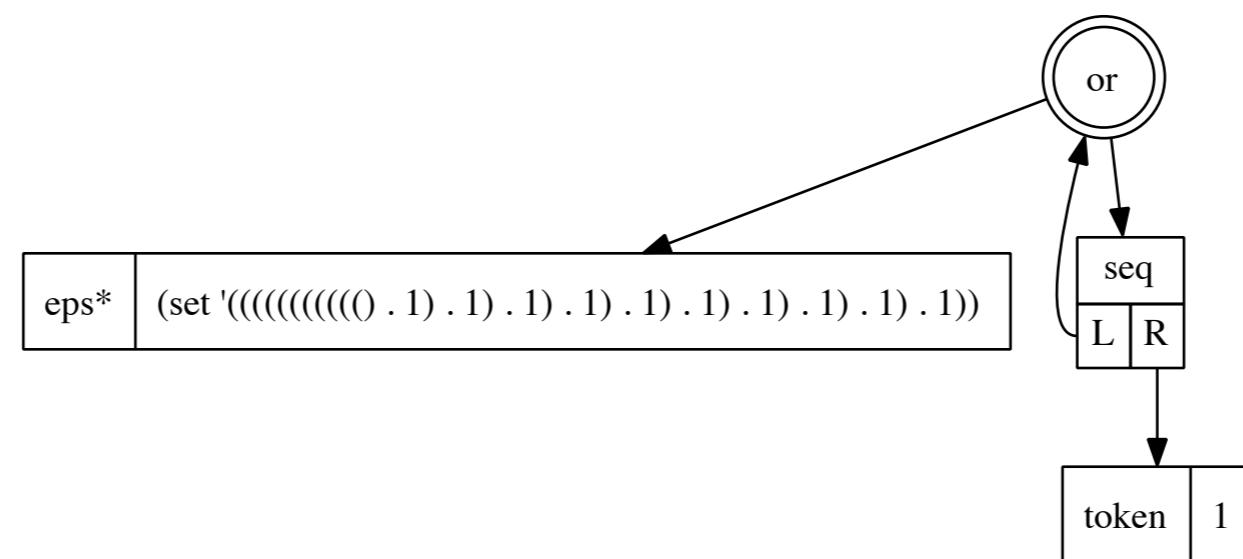
$$((\epsilon \downarrow \{t_1\}) \circ p) \rightarrow f \Rightarrow p \rightarrow \lambda t_2.(t_1, t_2)$$

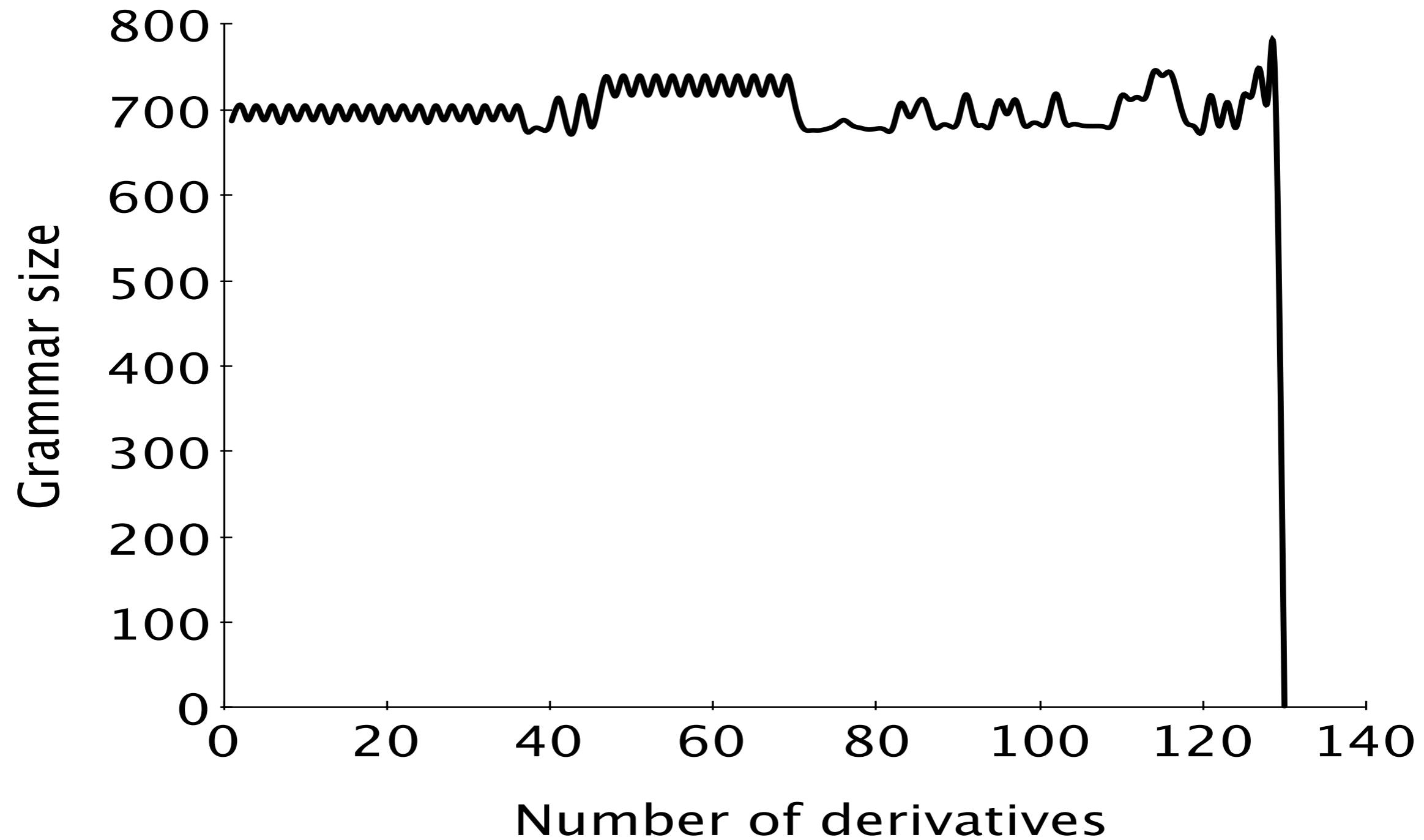
$$(p \rightarrow f) \rightarrow g \Rightarrow p \rightarrow (g \circ f)$$

$$\emptyset^{\star} \Rightarrow \epsilon \downarrow \{\langle \rangle\}.$$









# Practice

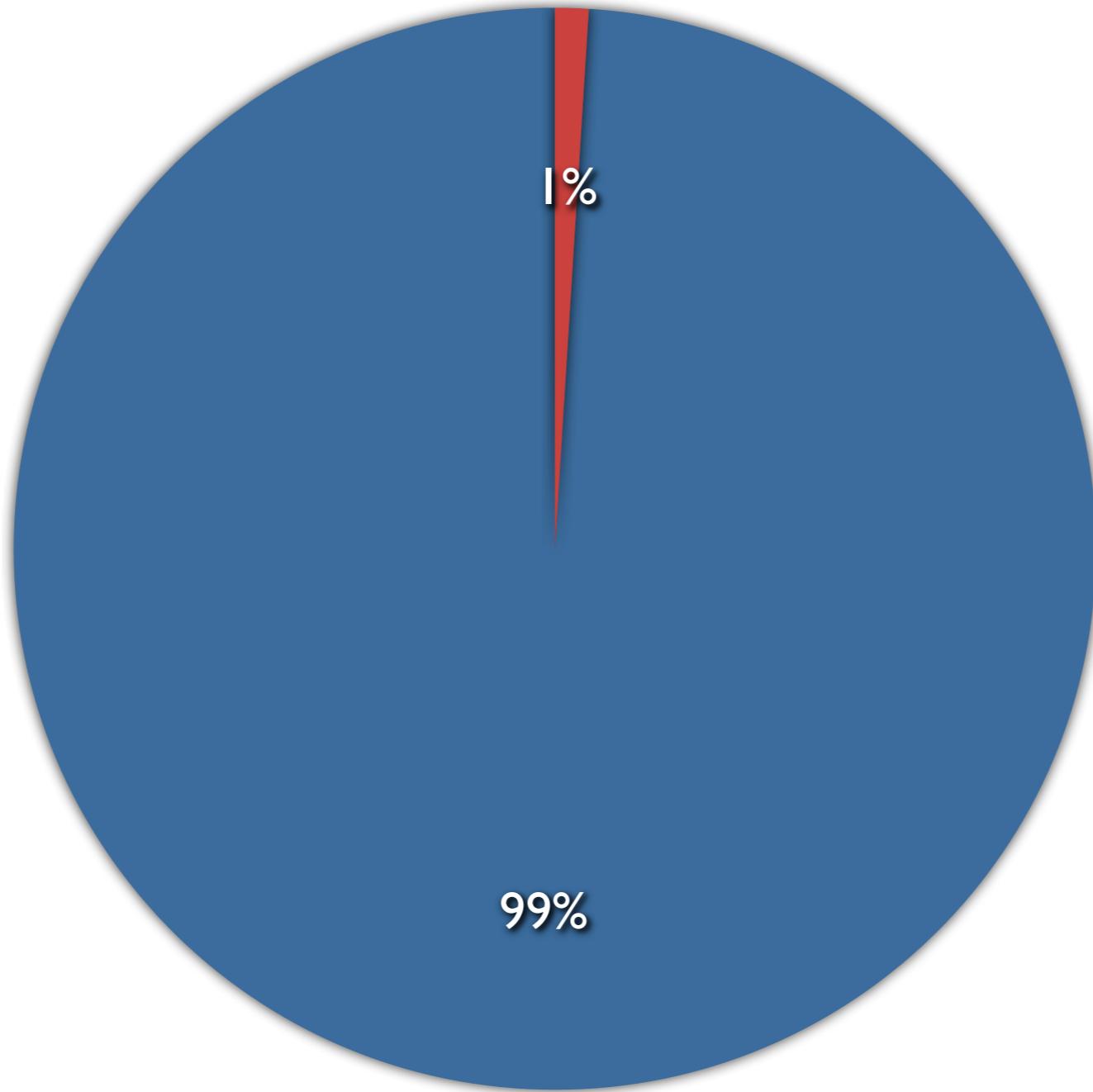
$\approx O(nG)$

# Performance

Good enough.

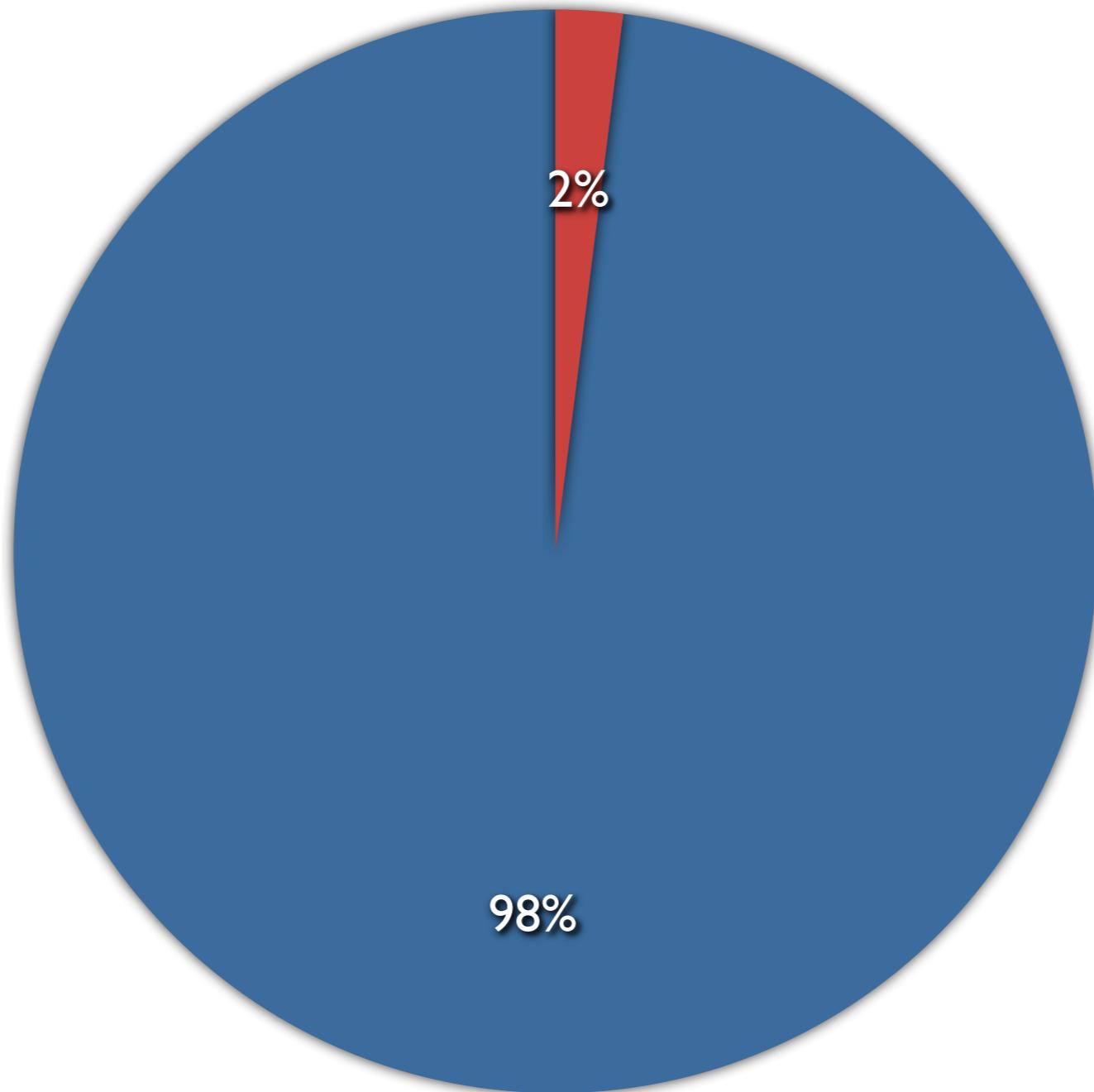
● Parsing

● Analysis



● Parsing

● Analysis



# Average complexity?

$O(n|G|)$ ?

# What we need

- Allow non-blocking I/O
- Act as library within the language
- Handle all CFGs: parse forest
- Draw on regular expressions
- Cacheable partial parses
- Be easy to implement

# What we need



Allow non-blocking I/O

- Act as library within the language
- Handle all CFGs: parse forest
- Draw on regular expressions
- Cacheable partial parses
- Be easy to implement

# What we need



Allow non-blocking I/O

Act as library within the language

- Handle all CFGs: parse forest
- Draw on regular expressions
- Cacheable partial parses
- Be easy to implement

# What we need



Allow non-blocking I/O



Act as library within the language



Handle all CFGs: parse forest

- Draw on regular expressions
- Cacheable partial parses
- Be easy to implement

# What we need

- Allow non-blocking I/O
- Act as library within the language
- Handle all CFGs: parse forest
- Draw on regular expressions
  - Cacheable partial parses
  - Be easy to implement

# What we need

- Allow non-blocking I/O
- Act as library within the language
- Handle all CFGs: parse forest
- Draw on regular expressions
- Cacheable partial parses
  - Be easy to implement

# What we need

- Allow non-blocking I/O
- Act as library within the language
- Handle all CFGs: parse forest
- Draw on regular expressions
- Cacheable partial parses
- Be easy to implement

# Would be nice

- Dynamic reconfigurability
- Beyond context-free languages
- Built into the language
- Support compositionality

# Would be nice



Dynamic reconfigurability

- Beyond context-free languages
- Built into the language
- Support compositionality

# Would be nice



Dynamic reconfigurability

Beyond context-free languages

- Built into the language
- Support compositionality

# Would be nice

- Dynamic reconfigurability
- Beyond context-free languages
- Built into the language
- Support compositionality

# Would be nice

-  Dynamic reconfigurability
-  Beyond context-free languages
-  Built into the language
-  Support compositionality

# Questions