# On the Complexity and Performance of Parsing with Derivatives

Michael D. Adams,* Celeste Hollenbeck, Matthew Might

*    Seeking a tenure-track position in the upcoming hiring cycle

Parsing with Derivatives
can be fast
… and it's not too hard

ANTLR

Parsec

SLR

LR(k)

Yacc

LL(k)

Happy

CYK

Dypgen

# Isn't parsing a solved problem?

Packrat

Early

GLR

LALR(k)

Bison

PEG

Elkhound

Compilers

Interpreters

Domain-specific langauges

Syntax highlighting

Natural languages

Command lines

Configuration files

Serialization

Query languages

Network protocols

# *buf++

Apache
(HTTP - RFC 2616)

lighttpd
(HTTP - RFC 2616)

Courier
(IMAP - RFC 3501)

Freeenode IRCD
(IRC - RFC 2812)

```
$ ls /etc

/etc/fstab
/etc/crontab
/etc/hosts
/etc/passwd
/etc/protocols
/etc/services
/etc/shadow
/etc/sudoers
...
```

strtok

# Yacc    `*buf++`    strtok

Create .y file

Define %union lval

Define tokens/types

Define yacc rules

Compile yacc grammar     `*buf++`       `strtok`

Create .l lex file

#include "y.tab.h"

Define states

Define lex rules

Deal with yywrap()

Call yyparse()

# Parsing with Derivatives

(Might et al. 2011)

$$D_c(L) = \{\ w \mid cw \in L\ \}$$

$$D_o(D_o(D_f(\ L\ )))= \{\qquad \varepsilon \qquad\qquad \}$$

$$D_c(\boldsymbol{\varepsilon}) = \varnothing \qquad D_c(c) = \boldsymbol{\varepsilon}$$

$$D_c(\varnothing) = \varnothing \qquad D_c(c') = \varnothing \text{ if } c \neq c'$$

$$D_c(l_1 \cup l_2) = D_c(l_1) \cup D_c(l_2)$$

$$D_c(l_1 >> l_2) = (D_c(l_1) >> l_2) \cup D_c(l_2) \text{ if } \boldsymbol{\varepsilon} \in l_1$$

$$D_c(l_1 >> l_2) = D_c(l_1) >> l_2 \text{ if } \boldsymbol{\varepsilon} \notin l_1$$

```
(define/memoize (derive c l)
  #:order ([l #:eq] [c #:equal])
  (match l
    [(empty) l]
    [(eps) (empty)]
    [(token pred class)
     (if (pred c) (eps) (empty))]
    [(alt l1 l2) (alt (derive c l1)
                      (derive c l2))
    [(seq (and (nullablep) ...
     (alt (derive c l2)
     ...
    [(seq l1 l2) (seq (derive c l1) l2)]))
```

$$D_c(e) = \emptyset \qquad D_c(c) = \epsilon$$

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

$$D_c(l_1 \cup l_2) = D_c(l_1) \cup D_c(l_2)$$

$$D(l_1 \gg l_2) = (D(l_1) \gg l_2) \cup D(l_2) \text{ if } \epsilon \in l_1$$

$$D(l_1 \gg l_2) = D(l_1) \gg l_2 \text{ if } \epsilon \notin l_1$$

# Performance

~~$O(2^{2n})$~~

$O(n^3)$

~~Parsing 31 lines of Python:~~
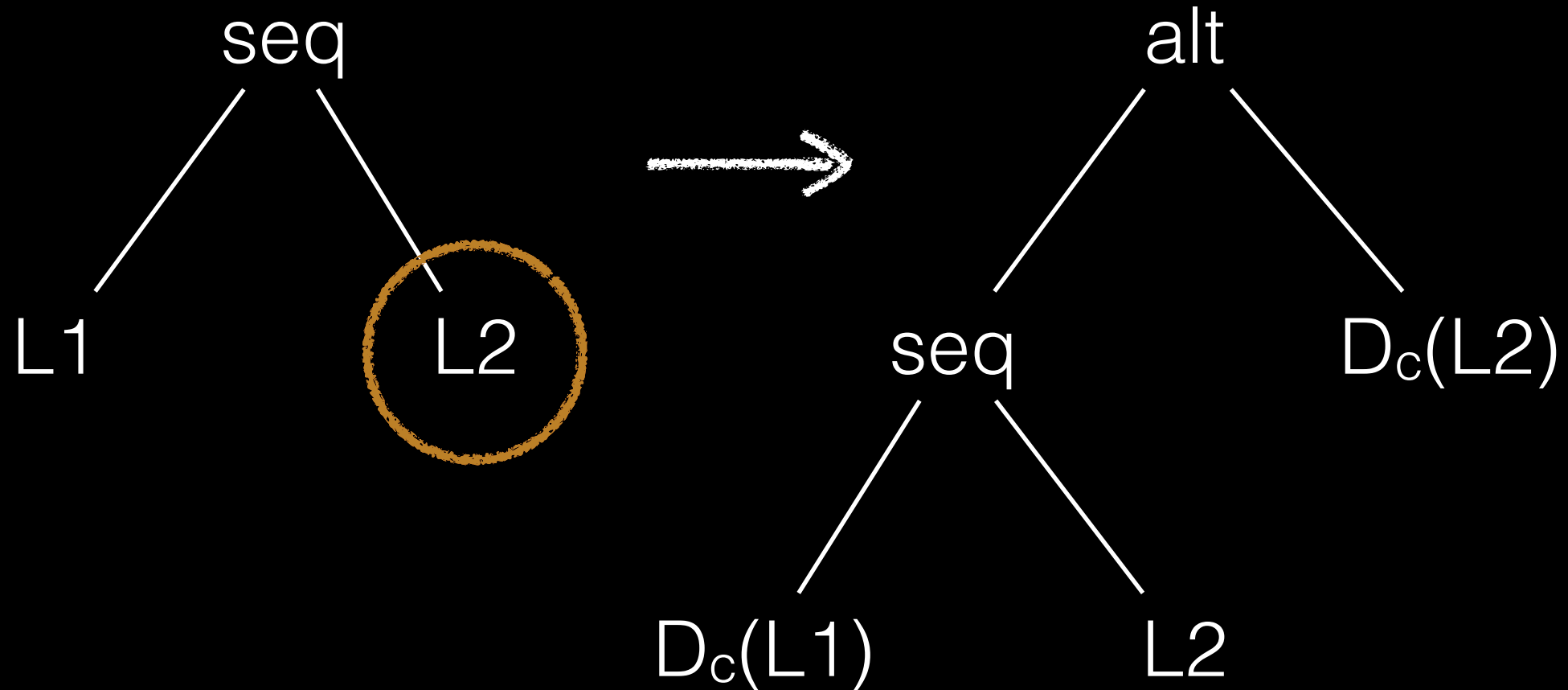~~24,000x slower than Bison~~
~~3 minutes~~
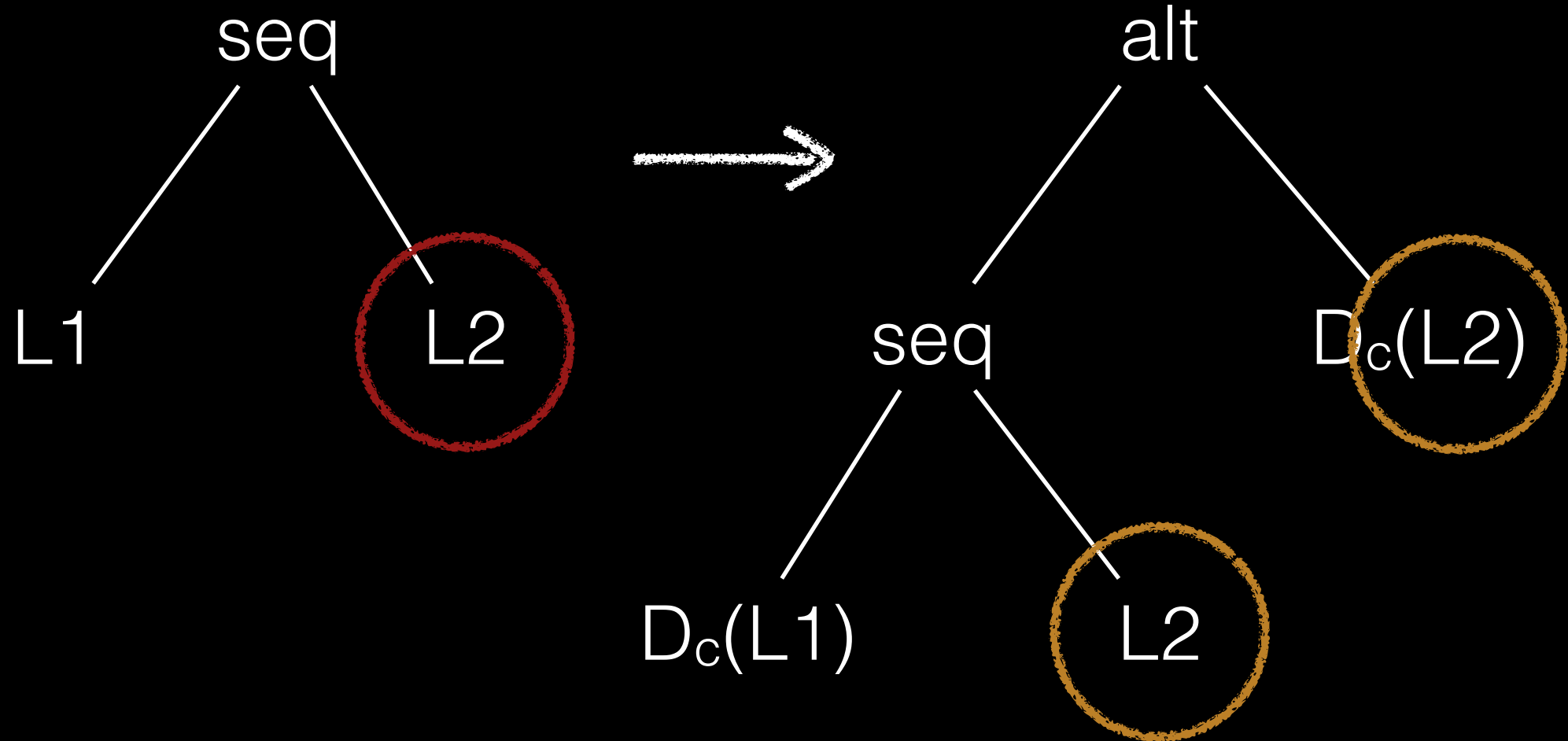25x slower than Bison
2 seconds

# Complexity Analysis

```
(define/memoize (D c L)
  #:order ([L #:eq] [c #:equal])
  (match L
    [(empty) (empty)]
    [(eps* T) (empty)]
    [(char a) (if (equal? a c)
                  (eps* (set c))
                  (empty))]
    [(red L f) (red (D c L) f)]
    [(alt L1 L2) (alt (D c L1) (D c L2))]
    [(seq L1 L2)
     (if (nullable? L1)
         (alt (seq (D c L1) L2) (D c L2))
         (seq (D c L1) L2))]))
```

```
[(seq L1 L2)
 (if (nullable? L1)
  (alt (seq (D c L1) L2) (D c L2))
```

seq → alt

L1    L2    seq         $D_c$(L2)

         $D_c$(L1)    L2

```
[(seq L1 L2)
 (if (nullable? L1)
  (alt (seq (D c L1) L2) (D c L2))
```

seq
├── L1
└── L2

→

alt
├── seq
│   ├── D_c(L1)
│   └── L2
└── D_c(L2)

# Time

$\leq$ k · # of allocated nodes

$\leq$ k · # of node names

$\leq$ O(n$^3$)

$$L \qquad\qquad L$$

$$D_f(L) \qquad\qquad Lf$$

$$D_o(D_f(L)) \qquad\qquad Lfo$$

$$D_o(D_o(D_f(L))) \qquad Lfoo$$

$$\# \text{ of names} \leq |G| \cdot n^2$$

\# of initial grammar nodes

\# of possible substrings of input

```
[(seq L1 L2)
 (if (nullable? L1)
   (alt (seq (D c L1) L2) (D c L2))
```



One problem…

# of names ≤

$La_1 \cdots a_k$

$La_1 \cdots a_i \bullet a_{i-1} \cdots a_k$

$$O(n^3)$$

$|G| \cdot n^2 +$

$|G| \cdot n^2 \cdot n$

# of initial grammar nodes

# of possible substrings of input

# of possible positions for bullet
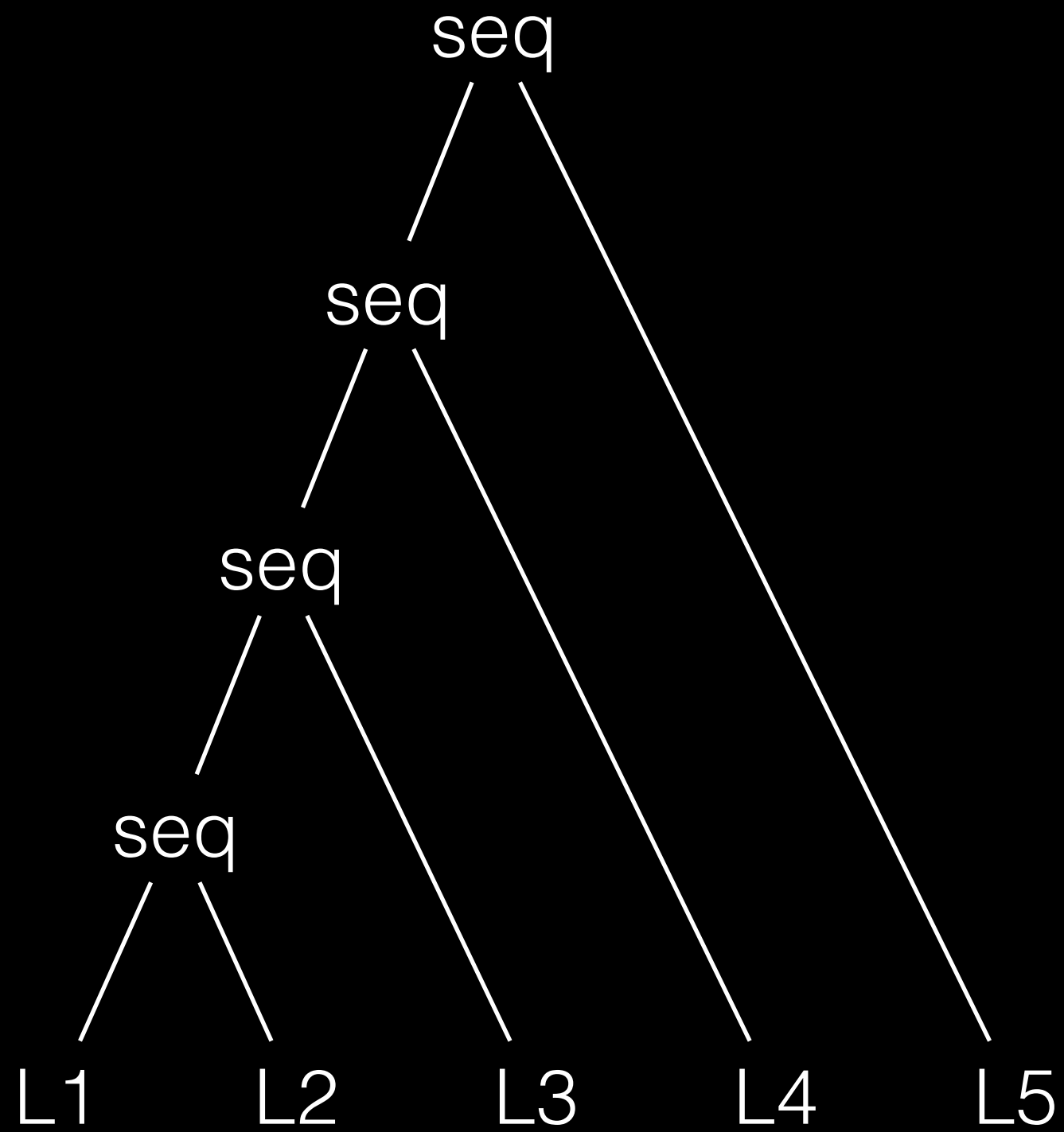
# Performance

# Optimizations

Compaction:
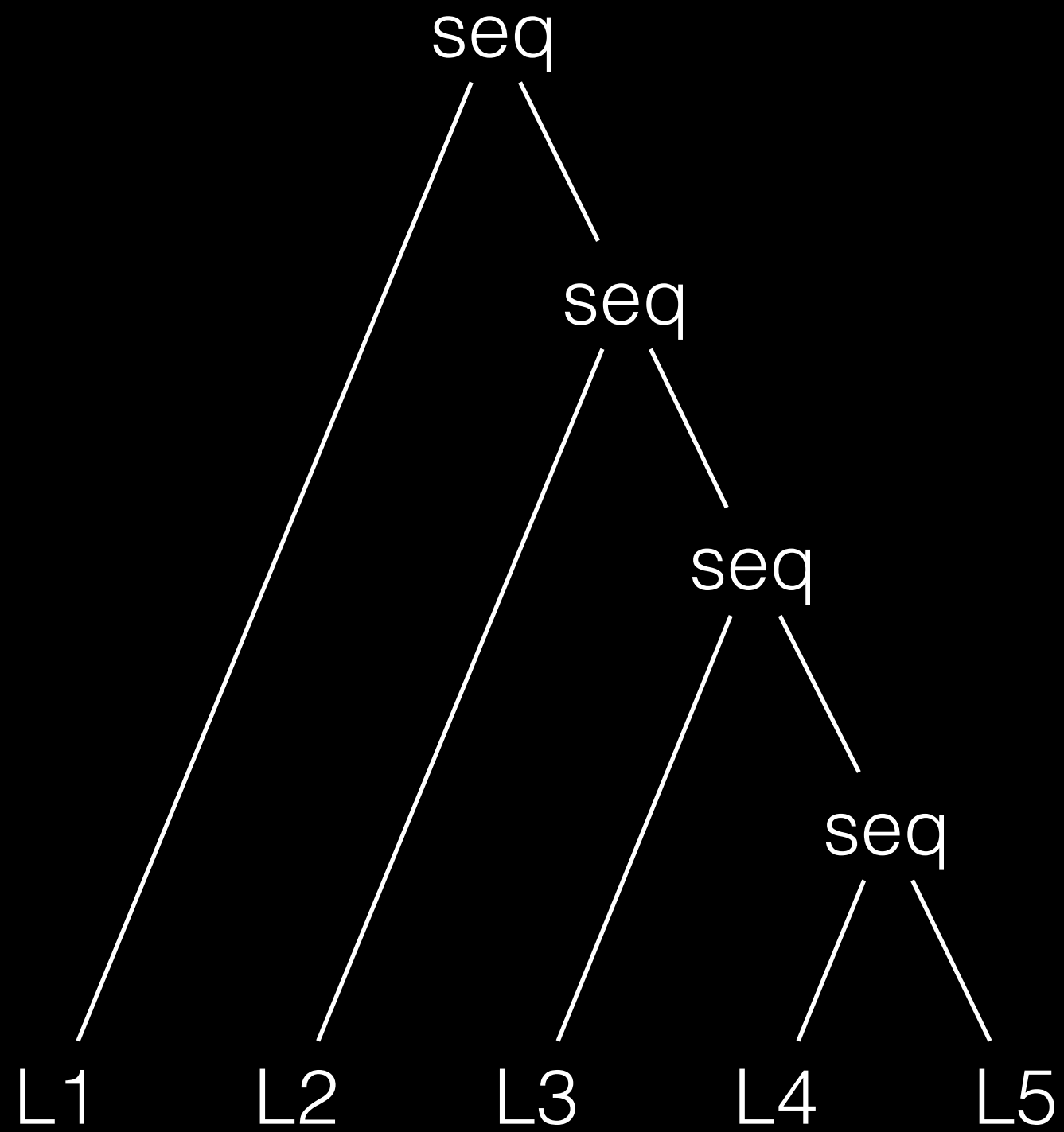 - Pre-processing
 - Canonicalization
 - Smart Constructors

Nullability:
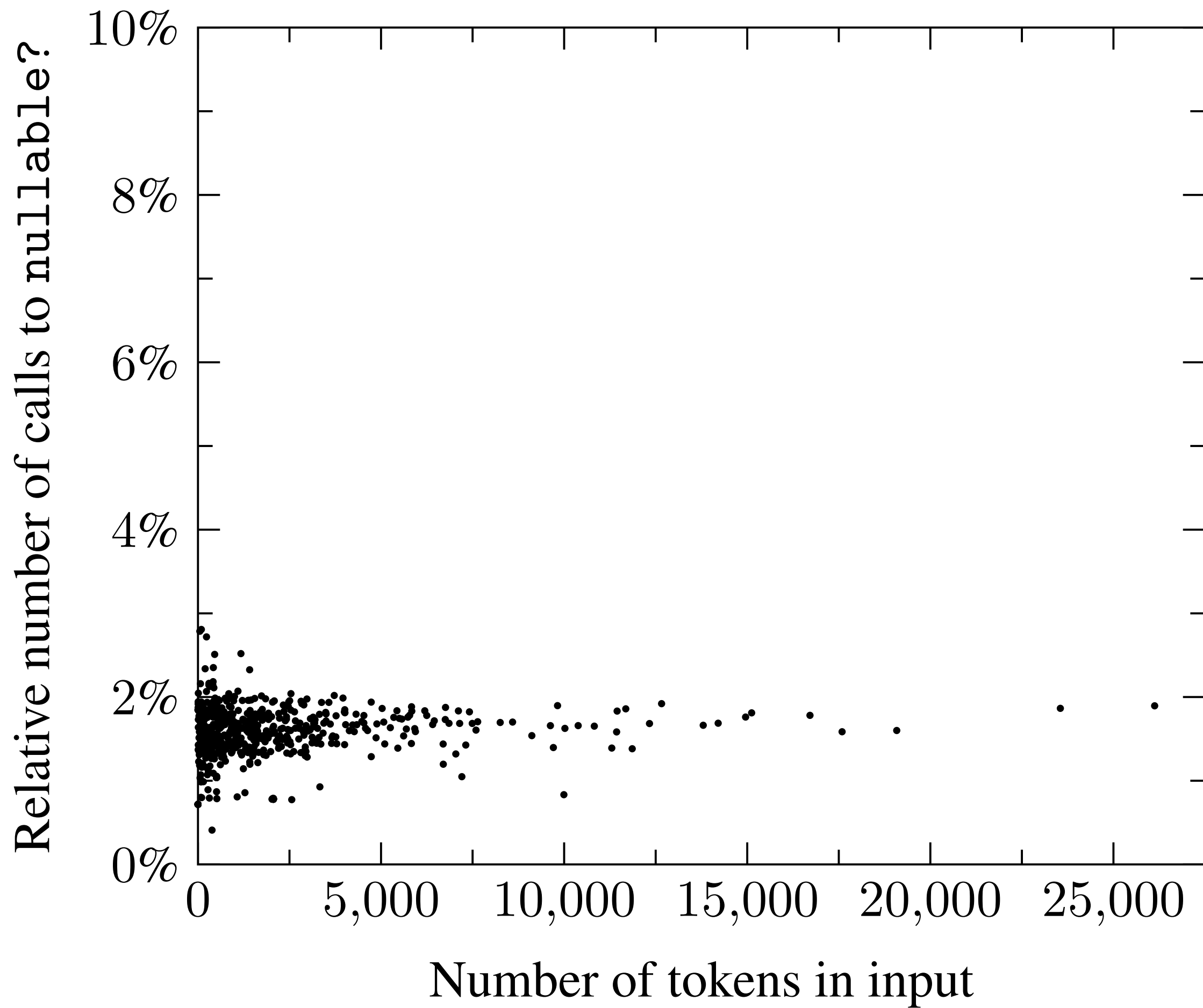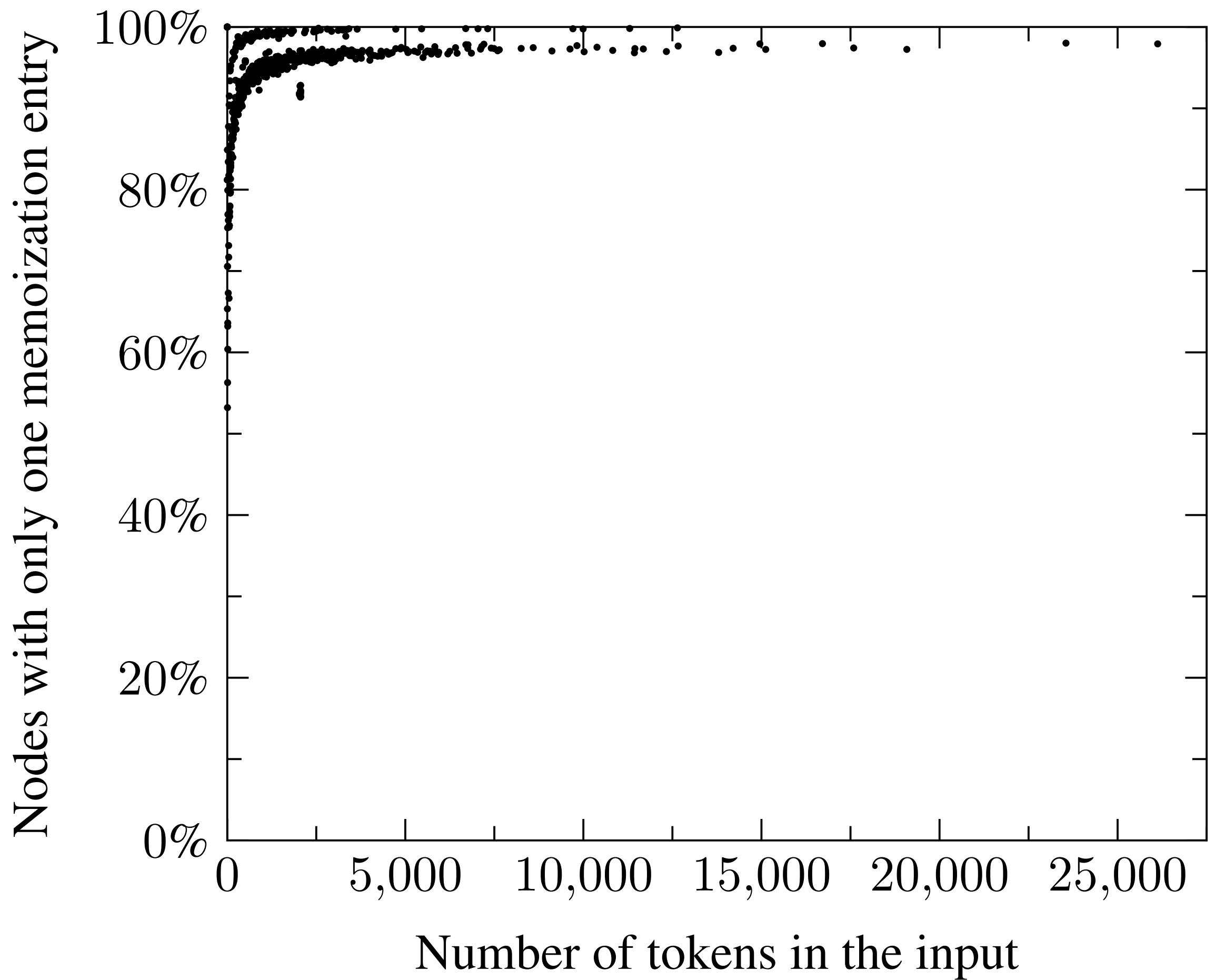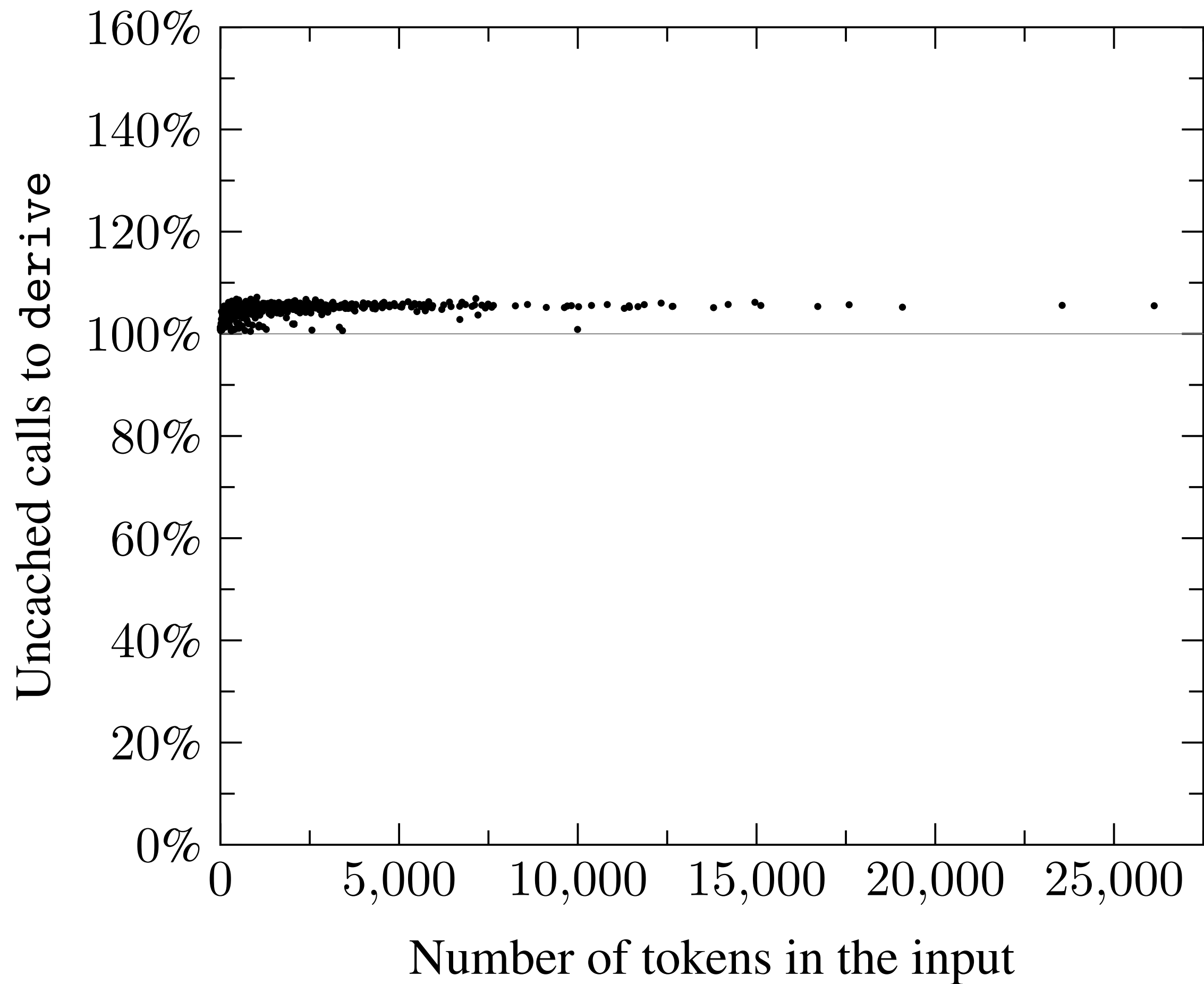 - Demand Based
 - Field Instead of Hashtable

Memoization:
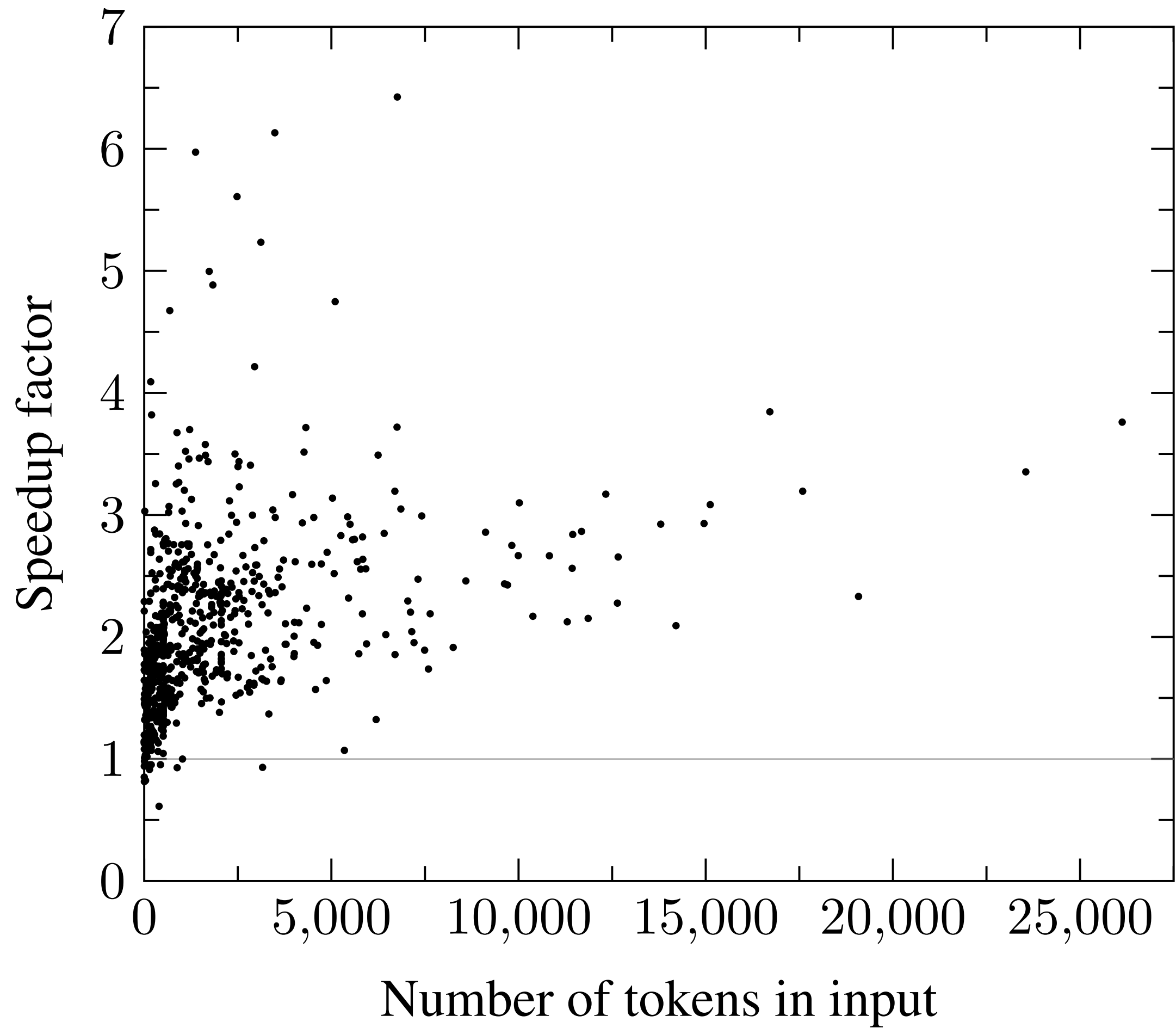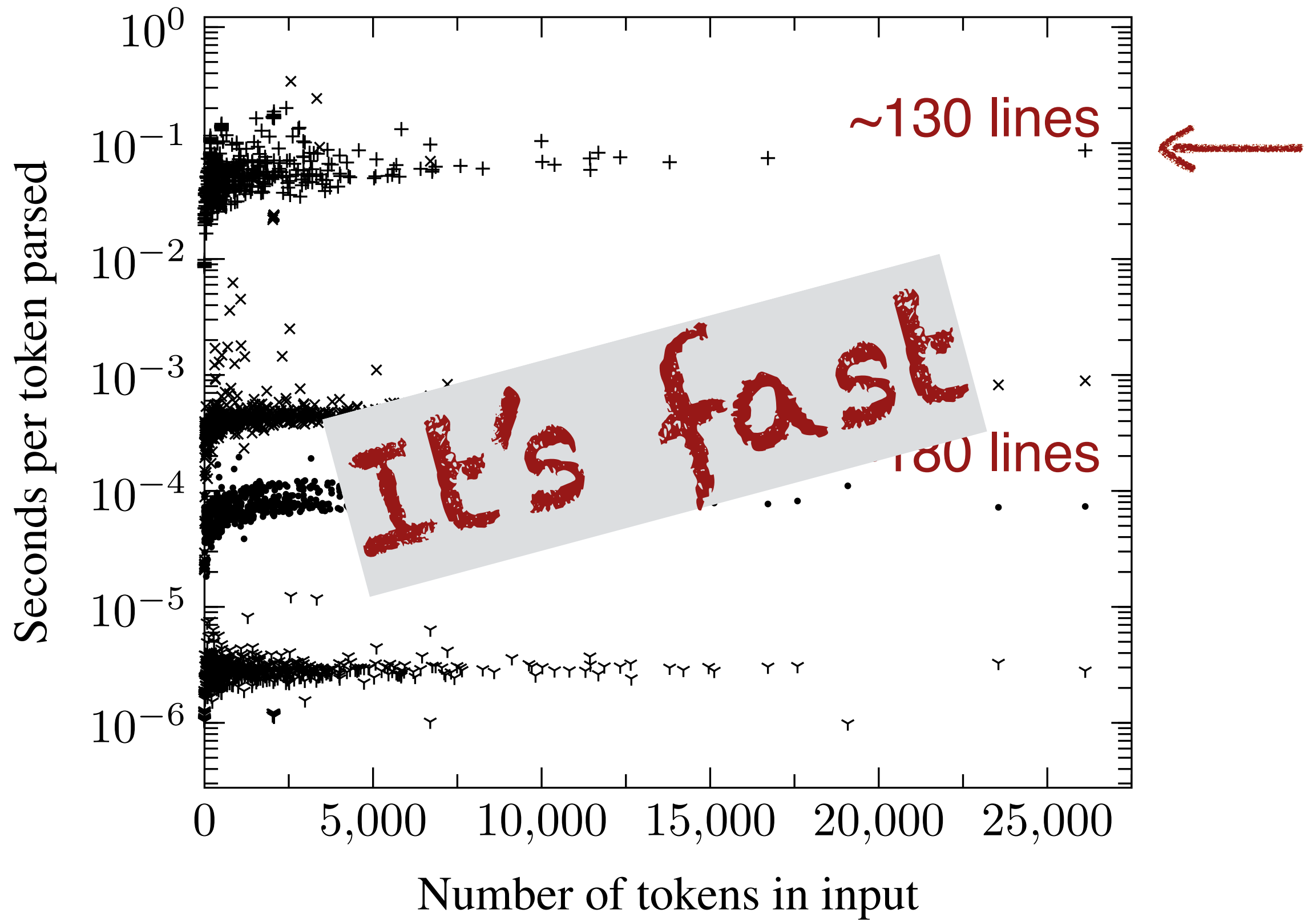 - One-entry/Forgetful Hashtables
 - Field Instead of Hashtable

# On the Complexity and Performance of Parsing with Derivatives

$O(2^{2^n})$     $O(n^3)$

Michael D. Adams * Celeste Hollenbeck, Matthew Might

24,000x     25x

Implementation available at
http://michaeldadams.org/papers/derivatives2/

* Seeking a tenure-track position in the upcoming hiring cycle