

PARSING WITH REGULAR EXPRESSIONS AND EXTENSIONS TO KLEENE ALGEBRA

Niels Bjørn Bugge Grathwohl

DIKU, November 4th 2015

PhD Thesis defense



STRING REWRITING

1,John,john@gmail.com,male,123456,DK
2,Benny,benny@hotmail.com,male,98234,UK

→

John 123456
Benny 98234

Want:

- *Streaming* – i.e., output while reading input.
- *Fast* – several Gbps throughput per core.
- Linear running time in the size of the input.

```
main := (row /\n/)*  
col  := /[^\n]*/  
row  := ~(col /,/ ) col "\t" ~/,/ ~(col /,/ )  
      ~(col /,/ ) col ~/,/      ~col
```

Program is essentially a *regular expression* with outputs.

Regular expression syntax

$$E ::= 0 \mid 1 \mid a \mid E_1 + E_2 \mid E_1E_2 \mid E_1^*$$

$(a \in \Sigma)$

Examples

$(\Sigma = \{a, b\})$

a

$(ab)^* + (a + b)^*$

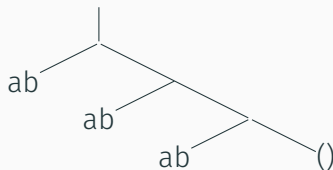
$(a + b)^*$

WHAT IS REGULAR EXPRESSION “MATCHING”?

Expression $(ab)^* + (a + b)^*$

Input $s = ababab$

- **acceptance testing**—is input string member of language?
Answer: “Yes!”
- **subgroup matching**—substrings in input for subterms in expression.
Answer: $[0, 6], [4, 2]$
- **parsing**—what is the parse tree of the input?



Input s matches E iff $s \in \mathcal{L}[E]$.

Language interpretation

$$\mathcal{L}[0] = \emptyset$$

$$\mathcal{L}[1] = \{\epsilon\}$$

$$\mathcal{L}[a] = \{a\}$$

$$\begin{aligned} \mathcal{L}[E + F] &= \{s \mid s \in \mathcal{L}[E]\} \\ &\cup \{t \mid t \in \mathcal{L}[F]\} \end{aligned}$$

$$\mathcal{L}[EF] = \{s \cdot t \mid s \in \mathcal{L}[E], t \in \mathcal{L}[F]\}$$

$$\mathcal{L}[E^*] = \mathcal{L}[E]^*$$

Example

$$\begin{aligned} & \mathcal{L}[(ab)^* + (a + b)^*] \\ = & \mathcal{L}[(ab)^*] \cup \mathcal{L}[(a + b)^*] \\ = & \mathcal{L}[ab]^* \cup \mathcal{L}[a + b]^* \\ = & \{ab\}^* \cup \{a, b\}^* \\ = & \{\epsilon, ab, abab, \dots\} \cup \{\epsilon, a, b, ab, ba, aba, \dots\} \\ = & \{\epsilon, a, b, aa, ab, aaa, aab, \dots\} \end{aligned}$$

Construct parse tree from input s such that *flattening* of parse tree is s .

Type interpretation [FC'04;HN'11]

$$\mathcal{T}[0] = \emptyset$$

$$\mathcal{T}[1] = \{()\}$$

$$\mathcal{T}[a] = \{a\}$$

$$\begin{aligned} \mathcal{T}[E + F] &= \{\text{inl } v \mid v \in \mathcal{T}[E]\} \\ &\cup \{\text{inr } w \mid w \in \mathcal{T}[F]\} \end{aligned}$$

$$\mathcal{T}[EF] = \mathcal{T}[E] \times \mathcal{T}[F]$$

$$\mathcal{T}[E^*] = \{[v_1, \dots, v_n] \mid n \geq 0, v_i \in \mathcal{T}[E]\}$$

Values in $\mathcal{T}[E]$ are *parse trees*.

Example

$\mathcal{T}[(ab)^* + (a + b)^*]$ contains the parse trees:

- $\text{inl} [(a, b), (a, b), (a, b)]$
- $\text{inr} [\text{inl } a, \text{inr } b, \text{inl } a, \text{inr } b, \text{inl } a, \text{inr } b]$

which are *not* in $\mathcal{T}[(a + b)^*]$!

So

$$\mathcal{T}[(ab)^* + (a + b)^*] \neq \mathcal{T}[(a + b)^*],$$

whereas

$$\mathcal{L}[(ab)^* + (a + b)^*] = \mathcal{L}[(a + b)^*]$$

One input string can be parsed in multiple ways: **ababab**
under $E = (ab)^* + (a + b)^*$ can be parsed *both* as

inl $[(a, b), (a, b), (a, b)]$

and

inr [inl a, inr b, inl a, inr b, inl a, inr b]

Disambiguation policy: the **left-most** option is always prioritized. “Greedy parsing.”

AMBIGUITY

One input string can be parsed in multiple ways: **ababab**
under $E = (ab)^* + (a + b)^*$ can be parsed *both* as

inl [(a, b), (a, b), (a, b)]

and

inr [inl a, inr b, inl a, inr b, inl a, inr b]

Disambiguation policy: the **left-most** option is always prioritized. “Greedy parsing.”

Bit-coded parse trees: only store *choices*.

Parse tree as stream of bits; **meaningless** without expression!

Example

$E = (ab)^* + (a + b)^*$, **ababab**:

inl [(a, b), (a, b), (a, b)]

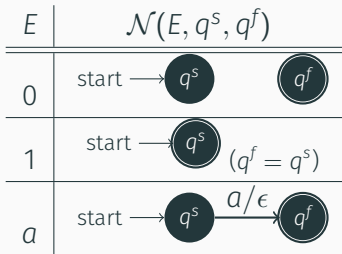
00001

inr [inl a, inr b, inl a, inr b, inl a, inr b]

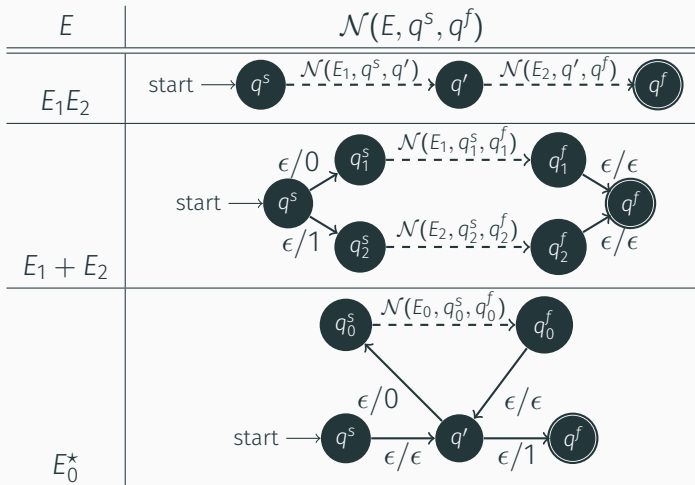
10001000100011

FINITE STATE TRANSUCERS

- Thompsons FSTs with input alphabet Σ , output alphabet $\{0, 1\}$.
- Construction:



FINITE STATE TRANSDUCTIONS



Theorem (Brüggemann-Klein 1993, GHNR 2013)

1-to-1 correspondence between

- parse trees for E ,
- paths in Thompson FST for E ,
- bit-coded parse trees.

Constructing the parse tree corresponds to finding a path through the FST.

Optimally streaming parsing

Output the longest common prefix of possible parse trees after reading each input symbol.

Example

$$E = (aaa + aa)^*$$

Possible parse tree prefixes after **aaaa**:

$$\{01011, 000\dots\}$$

Possible parse tree prefixes after **aaaaa**:

$$\{00011, 0000\dots\}$$

GREEDY PARSING

	Time	Space	Aux	Answer
Parse (3-p) ¹	$O(mn)$	$O(m)$	$O(n)$	greedy parse
Parse (2-p) ²	$O(mn)$	$O(m)$	$O(n)$	greedy parse
Parse (str.) ³	$O(mn + 2^{m \log m})$	$O(m)$	$O(n)$	greedy parse

(n size of input, m size of expression)

¹Frisch, Cardelli (2004)

²Grathwohl, Henglein, Nielsen, Rasmussen (2013)

³Grathwohl, Henglein, Rasmussen (2014)

GREEDY PARSING

	Time	Space	Aux	Answer
Parse (3-p) ¹	$O(mn)$	$O(m)$	$O(n)$	greedy parse
Parse (2-p) ²	$O(mn)$	$O(m)$	$O(n)$	greedy parse
Parse (str.) ³	$O(mn + 2^{m \log m})$	$O(m)$	$O(n)$	greedy parse

(n size of input, m size of expression)

¹Frisch, Cardelli (2004)

²Grathwohl, Henglein, Nielsen, Rasmussen (2013)

³Grathwohl, Henglein, Rasmussen (2014)

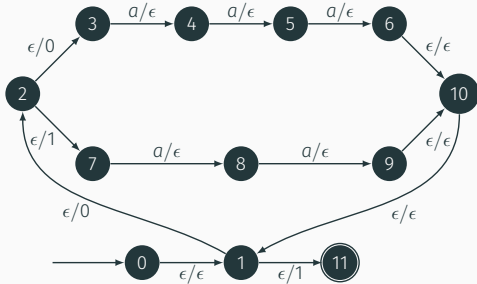
Optimally streaming algorithm

- Preprocessing step of FST: compute *coverage* of state sets.
- Maintain a *path tree* during FST simulation, recording the path taken to each state in the FST.
 - Prune states that are covered by higher-prioritized states.
- Output on the stem of the path tree is longest common prefix of any succeeding parse.

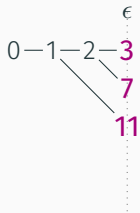
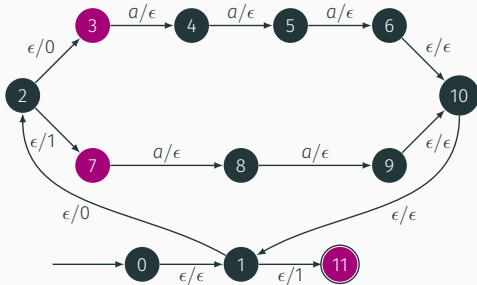
Theorem (GHR'14)

Optimally streaming algorithm computes the optimally streaming parsing function in time $O(mn + 2^{m \log m})$.

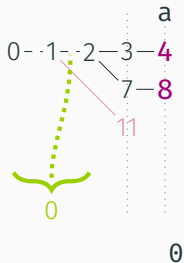
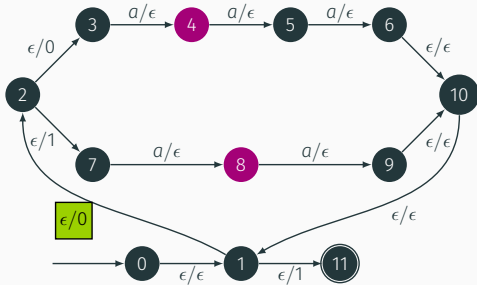
PATH TREE EXAMPLE: $(aaa + aa)^*$



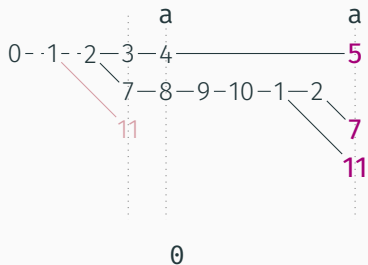
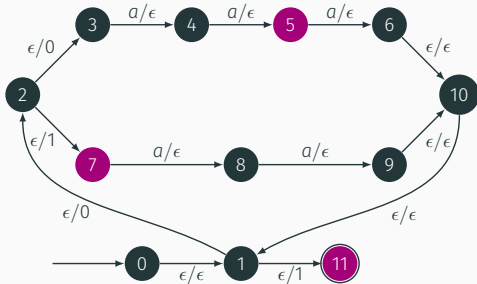
PATH TREE EXAMPLE: $(aaa + aa)^*$



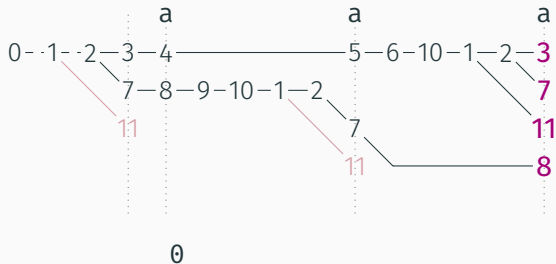
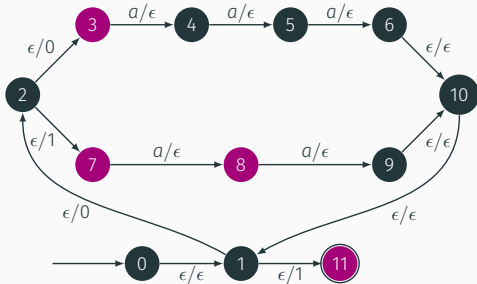
PATH TREE EXAMPLE: $(aaa + aa)^*$



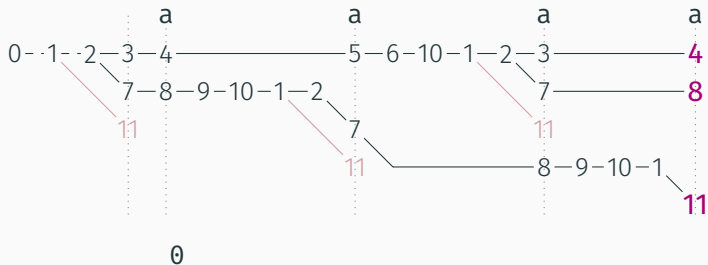
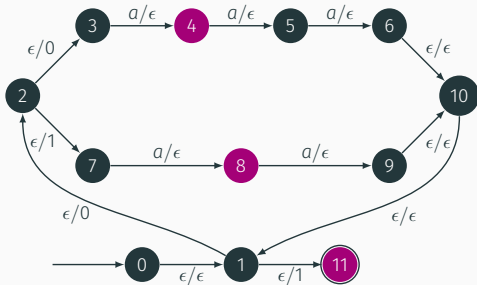
PATH TREE EXAMPLE: $(aaa + aa)^*$



PATH TREE EXAMPLE: $(aaa + aa)^*$



PATH TREE EXAMPLE: $(aaa + aa)^*$



Observation

Approach is not limited to Thompson FSTs outputting bit-coded parse trees.

Kleenex is a surface language for specifying FSTs and their output:

- *grammar* with greedy disambiguation;
- embedded *output actions*.

- Essentially optimally streaming behaviour.
- Linear running time in size of input string.
- *Fast*. >1 Gbps common.

```
main := (num /\n/)*  
num  := digit{1,3} ("," digit{3})*  
digit := /[0-9]/
```

"1000000000000" → "100,000,000,000"

- Problem: need to read entire number; no bounded lookahead!
- But: each newline ends a number, so output.
- Optimal streaming gives this for free!

Path tree algorithm is “NFA simulation with path trees as state sets.”

Compilation of FSTs? Analogy to NFA-DFA determinization with subset construction?

DETERMINIZATION

Path tree algorithm is “NFA simulation with path trees as state sets.”

Compilation of FSTs? Analogy to NFA-DFA determinization with subset construction?

Problem: **Inifinite number of path trees!**

DETERMINIZATION

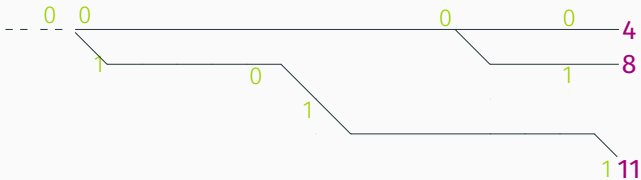
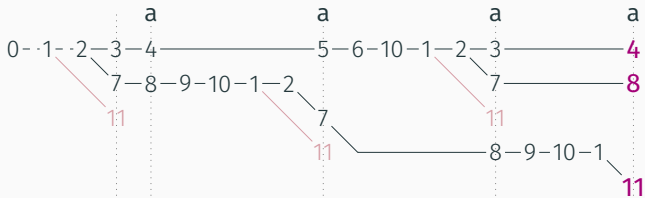
Path tree algorithm is “NFA simulation with path trees as state sets.”

Compilation of FSTs? Analogy to NFA-DFA determinization with subset construction?

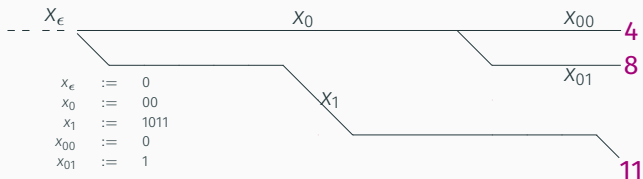
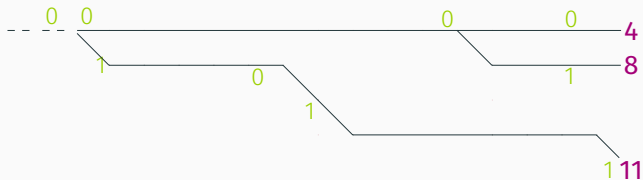
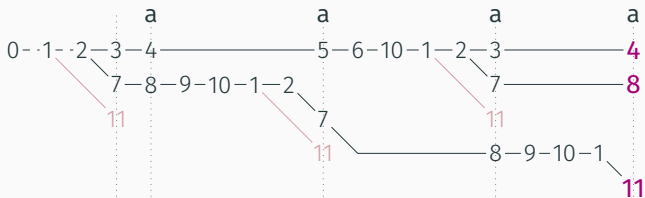
Problem: **Infinite number of path trees!**

Solution: *contract* unary paths in path trees and store output in registers.

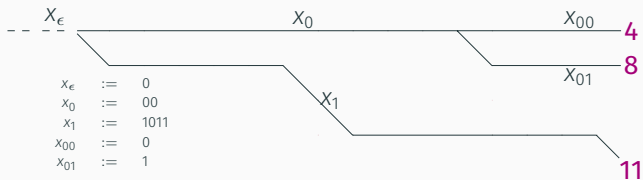
DETERMINIZATION



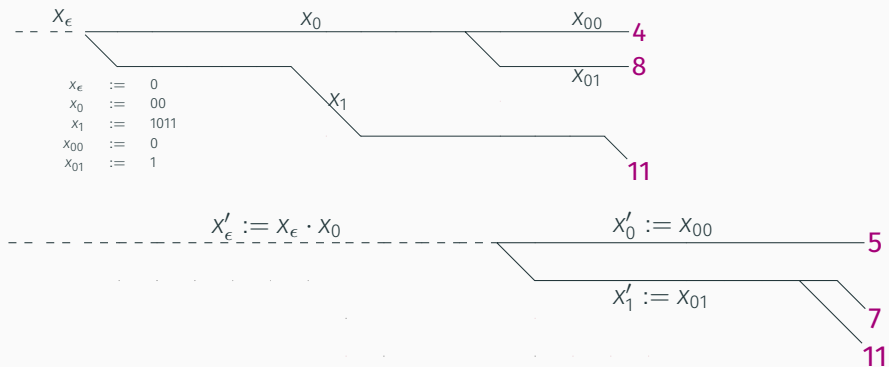
DETERMINIZATION



DETERMINIZATION



DETERMINIZATION



DETERMINIZATION

- Streaming string transducer:
 - deterministic finite automata,
 - each state equipped with fixed number of registers containing strings
 - registers updated on transition by affine function;
 - Alur, D'Antoni, Raghothaman (2015).

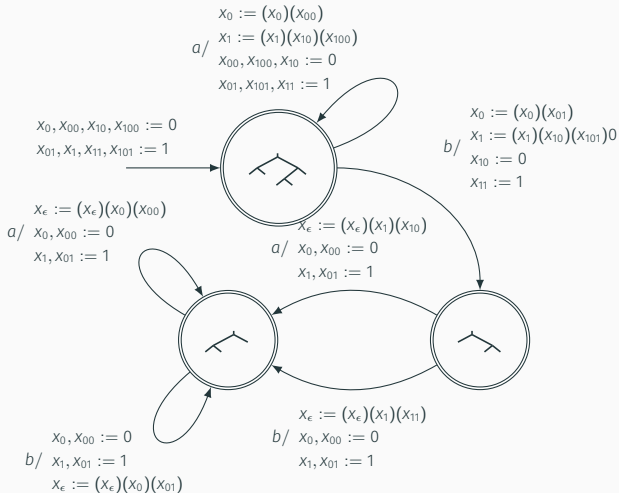
- Streaming string transducer:
 - deterministic finite automata,
 - each state equipped with fixed number of registers containing strings
 - registers updated on transition by affine function;
 - Alur, D'Antoni, Raghothaman (2015).

Theorem

FSTs with greedy order semantics correspond to SSTs.

- States are contracted path trees.
- Edges in contracted path trees \cong registers in SST.

DETERMINIZATION



Haskell implementation

Kleenex source \rightarrow FST \rightarrow SST \rightarrow C

C code compiled with GCC/clang

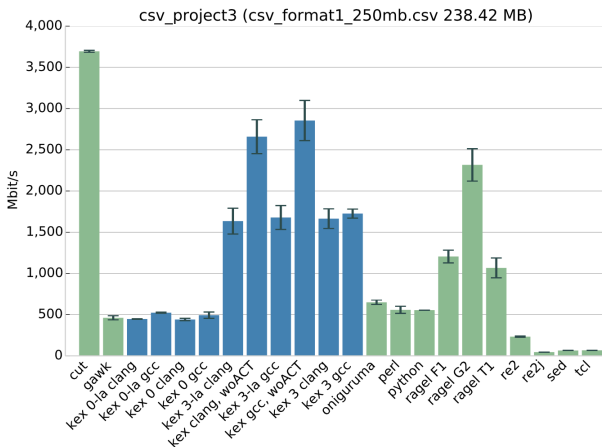
Performance comparison with regular expression libraries:

- AWK, Perl, Python, Sed, Tcl
- RE2/RE2j
- Ragel state machine compiler

<https://github.com/diku-kmc/kleenexlang>

PERFORMANCE

```
main := (row /\n/)*  
col  := /[^\n]*/  
row  := ~(col /,/ ) col "\t" ~/,/ ~(col /,/ )  
      ~(col /,/ ) col ~/,/      ~col
```



FUTURE WORK

- Program fragments as output actions
- Memoization techniques à la NFA/DFA memoization in RE2.
- Applications – bioinformatics, finance, log digging,;
- Parallel processing: read >8 bits in parallel;
- Approximate matching – necessary in biological applications;
- Expressiveness, visibly pushdown automata;
- Automatically generate interfaces for various programming languages.

KLEENE ALGEBRA

Kleene algebra

A structure $(K, +, \cdot, *, 0, 1)$:

- A set of elements K ,
- binary operators $+$ and \cdot ,
- unary operator $*$,
- special elements 0 and 1 ,

that satisfies the *Kleene algebra axioms*.

Semiring

$$x \cdot (y \cdot z) = (x \cdot y) \cdot z \qquad x + (y + z) = (x + y) + z$$

$$1 \cdot x = x = x \cdot 1$$

$$0 + x = x = x + 0$$

$$x + y = y + x$$

$$x \cdot (y + z) = x \cdot y + x \cdot z$$

$$(x + y) \cdot z = x \cdot z + y \cdot z$$

$$0 \cdot x = 0 = x \cdot 0$$

- idempotence: $x + x = x$
- partial order: $x \leq y \iff x + y = y$

Kleene algebra

Idempotent semiring with $*$ operator:

$$1 + x \cdot x^* \leq x^*$$

$$1 + x^* \cdot x \leq x^*$$

$$b + a \cdot x \leq x \implies a^* \cdot b \leq x$$

$$b + x \cdot a \leq x \implies b \cdot a^* \leq x$$

Any structure with these operators that satisfies the axioms is a Kleene algebra.

- Languages: $(L, \cup, \cdot, *, \emptyset, \{\epsilon\})$.
 - L is set of strings over an alphabet,
 - \cup is set union,
 - \cdot is string concatenation,
 - $*$ is repetition of strings,
 - partial order \leq is subset inclusion \subseteq .

Language interpretation of regular expressions from before.

- Relation model, tropical semiring, ...

“Regular expressions:” syntax to describe elements in a Kleene algebra.

Canonical interpretation

The *canonical interpretation* of a term E is the regular language interpretation:

$$L_{\Sigma}(x) = \{x\} \quad L_{\Sigma}(e_0 + e_1) = L_{\Sigma}(e_0) \cup L_{\Sigma}(e_1)$$

$$L_{\Sigma}(0) = \emptyset \quad L_{\Sigma}(e_0 e_1) = \{vw \mid v \in L_{\Sigma}(e_0), w \in L_{\Sigma}(e_1)\}$$

$$L_{\Sigma}(1) = \{\epsilon\} \quad L_{\Sigma}(e^*) = \bigcup_{n \geq 0} L_{\Sigma}(e^n).$$

Polynomials

Given idempotent semiring C and a set of variables X , form *polynomials* over C and X :

$$0 \quad a \quad ax^2 + bxy^3 + 1 \quad 1 + a + ax + by$$

System of polynomial inequalities

$$\begin{aligned} 1 + aB + bA &\leq S \\ A + aS + bAA &\leq A \\ bS + aBB &\leq B \end{aligned}$$

Solution: valuation of variables in X such that the inequalities are satisfied.

A semiring C is *algebraically closed* if all finite systems of polynomials have *least* solutions.

Definition

A *Chomsky algebra* is a an algebraically closed idempotent semiring.

Context-free languages over symbols from X are not Kleene algebras, but they are Chomsky algebras.

Context-free grammar corresponds to system of polynomial inequalities:

$$S \rightarrow \epsilon \mid aB \mid bA$$

$$A \rightarrow aS \mid bAA$$

$$B \rightarrow bS \mid aBB$$

$$1 + aB + bA \leq S$$

$$aS + bAA \leq A$$

$$bS + aBB \leq B$$

- Regular expressions: denote elements in Kleene algebra.
- μ -terms: denote elements in Chomsky algebra.

μ -terms

T_X are μ -terms over an alphabet X :

$$t ::= 0 \mid 1 \mid x \mid t + t \mid t \cdot t \mid \mu x.t \quad x \in X$$

n -fold composition

$$0x.t \equiv 0 \qquad (n + 1)x.t \equiv t[x/nx.t]$$

Examples

$$0x.axb + 1 = 0$$

$$1x.axb + 1 = a(0x.axb + 1)b + 1 = a0b + 1 = 1$$

$$2x.axb + 1 = a(1x.axb + 1)b + 1 = ab + 1$$

Given interpretation of literals: $\sigma : X \rightarrow C$, *interpretation of μ -terms over Chomsky algebra C .*

Function $\sigma : TX \rightarrow C$ where:

$$\sigma(0) = 0 \quad \sigma(a + b) = \sigma(a) + \sigma(b)$$

$$\sigma(1) = 1 \quad \sigma(a \cdot b) = \sigma(a) \cdot \sigma(b)$$

$$\sigma(\mu x.t) = \text{least } a \in C \text{ such that } \sigma[x/a](t) \leq a$$

Canonical interpretation

Canonical interpretation as *context-free languages*:

$$\begin{aligned} L_X(x) &= \{x\} & L_X(t_0 + t_1) &= L_X(t_0) \cup L_X(t_1) \\ L_X(0) &= \emptyset & L_X(t_0 \cdot t_1) &= \{vw \mid v \in L_X(t_0), w \in L_X(t_1)\} \\ L_X(1) &= \{\epsilon\} & L_X(\mu x.t) &= \bigcup_{n \geq 0} L_X(nx.t). \end{aligned}$$

$\sum_{n \geq 0} t_n$ denotes *supremum* with respect to partial order \leq

μ -continuity

A Chomsky algebra C is μ -continuous if

$$\sigma(a(\mu x.t)b) = \sum_{n \geq 0} \sigma(a(nx.t)n)$$

for any interpretation σ over C .

Canonical interpretation as context-free language is μ -continuous:

$$L_X(\mu x.t) = \bigcup_{n \geq 0} L_X(nx.t).$$

Theorem

The following are equivalent:

- (i) $s = t$ holds in all μ -continuous Chomsky algebras,
- (ii) $L_X(s) = L_X(t)$ holds in the canonical interpretation as a context-free language over variables X .

AXIOMATIZATION

Two context free languages $L_X(s)$ and $L_X(t)$ are equivalent if and only if $s = t$ is provable from the axioms of μ -continuous Chomsky algebra:

Axioms

$$x \cdot (y \cdot z) = (x \cdot y) \cdot z$$

$$1 \cdot x = x = x \cdot 1$$

$$x \cdot (y + z) = x \cdot y + x \cdot z$$

$$(x + y) \cdot z = x \cdot z + y \cdot z$$

$$0 \cdot x = 0 = x \cdot 0$$

$$x + (y + z) = (x + y) + z$$

$$0 + x = x = x + 0$$

$$x + y = y + x$$

$$x + x = x$$

$$a(\mu x.t)b = \sum_{n \geq 0} a(nx.t)b$$

μ -continuity axiom is **infinitary**:

$$a(nx.t)b \leq a(\mu x.t)b, \quad n \geq 0$$
$$\left(\bigwedge_{n \geq 0} (a(nx.t)b \leq w) \right) \implies a(\mu x.t)b \leq w$$

- Equivalence of context-free languages is **undecidable**.
- To use inference, one must establish infinitely many premises.

SUMMARY, FURTHER DIRECTIONS

- Extend Chomsky algebra with test symbols, analogously to Kleene algebra with tests.
- Coalgebraic treatment of Chomsky algebra?
- Applications to program verification, like Kleene algebra?
- “Visibly pushdown” Chomsky algebra?
- KAT+B! is an extension to Kleene algebra with tests adding mutable state:
 - elements correspond to square matrices with regular language entries.
 - extend Kleenex with mutable state?

THANK YOU