

# Underlying Principles and Recurring Ideas of Formal Grammars

after a paper by Alexander Okhotin\*

talk by Tobias Denkinger

Freitagsseminar 2018-04-20

\* [Okhotin 18]

## Overview

1. Grammar types
2. Parse trees
3. Language equations
4. Parsing
5. Characterisation results

# Grammar types I

- ordinary grammar (aka context-free grammar) [Chomsky 56]

... is a tuple  $(N, \Sigma, R, S)$  of finite sets

with rules (in  $R$ ) of the form

$$A \rightarrow u_0 B_1 u_1 \cdots B_k u_k \quad u_i \in \Sigma^*, B_i \in N$$

# Grammar types I

- ordinary grammar (aka context-free grammar) [Chomsky 56]

... is a tuple  $(N, \Sigma, R, S)$  of finite sets

with rules ( $\in R$ ) of the form

$$A \rightarrow u_0 B_1 u_1 \cdots B_k u_k \quad u_i \in \Sigma^*, B_i \in N$$

characterised by 3 elements:

- constituents (CON)                          here: strings
- operations on constituents (C-OP)        here: concatenation •
- logical operations (L-OP)                    here: disjunction +

## Grammar types II

- linear grammars:  $A \rightarrow uBv \quad u, v \in \Sigma^*, A, B \in N$

C-OP: only concatenation of fixed strings to constituents

## Grammar types II

- linear grammars:  $A \rightarrow uBv \quad u, v \in \Sigma^*, A, B \in N$   
C-OP: only concatenation of fixed strings to constituents
- unambiguous grammars:
  - C-OP:  $K \cdot L$  only if  $\nexists w: \exists! (u, v): (u \cdot v = w) \wedge (u \in K) \wedge (v \in L)$
  - L-OP:  $K + L$  only if  $K \cap L = \emptyset$

## Grammar types II

- linear grammars:  $A \rightarrow uBv \quad u, v \in \Sigma^*, A, B \in N$   
C-OP: only concatenation of fixed strings to constituents
- unambiguous grammars:
  - C-OP:  $K \cdot L$  only if  $\forall w: \exists! (u, v): (u \cdot v = w) \wedge (u \in K) \wedge (v \in L)$
  - L-OP:  $K + L$  only if  $K \cap L = \emptyset$
- conjunctive grammar [Okhotin 01, 13]:  
L-OP: + and conjunction &

## Grammar types II

- linear grammars:  $A \rightarrow uBv \quad u, v \in \Sigma^*, A, B \in N$   
C-OP: only concatenation of fixed strings to constituents
- unambiguous grammars:
  - C-OP:  $K \cdot L$  only if  $\nexists w: \exists! (u, v): (u \cdot v = w) \wedge (u \in K) \wedge (v \in L)$
  - L-OP:  $K + L$  only if  $K \cap L = \emptyset$
- conjunctive grammar [Okhotin 01, 13]:  
L-OP: + and conjunction &
- Boolean grammars [Okhotin 04, 13]: L-OP: +, &, negation

## Grammar types III

- M-component grammars [Vijay-Shanker+87, Seki+91]

$$B_0(\alpha_1, \dots, \alpha_{m_0}) \leftarrow B_1(x_{1,1}, \dots, x_{1,m_1}) \cdots B_k(x_{k,1}, \dots, x_{k,m_k})$$

CON: ( $\leq M$ )-tuples of strings  $[m_i \leq M]$

C-OP: "combination" of string tuples

## Grammar types III

- $M$ -component grammars [Vijay-Shanker + 87, Seki + 91]

$$B_0(\alpha_1, \dots, \alpha_{m_0}) \leftarrow B_1(x_{1,1}, \dots, x_{1,m_1}) \cdots B_k(x_{k,1}, \dots, x_{k,m_k})$$

CON:  $(\leq M)$ -tuples of strings  $[m_i \leq M]$

C-OP: "combination" of string tuples

- well-nested  $M$ -component grammars:

CON:  $(\leq M)$ -tuples of strings

C-OP: "combination" of string tuples, but unrelated variables  
may not cross each other

## Grammar types III

- M-component grammars [Vijay-Shanker+87, Seki+91]

$$B_0(x_1, \dots, x_{m_0}) \leftarrow B_1(x_{1,1}, \dots, x_{1,m_1}) \cdots B_k(x_{k,1}, \dots, x_{k,m_k})$$

CON: ( $\leq M$ )-tuples of strings  $[m_i \leq M]$

C-OP: "combination" of string tuples

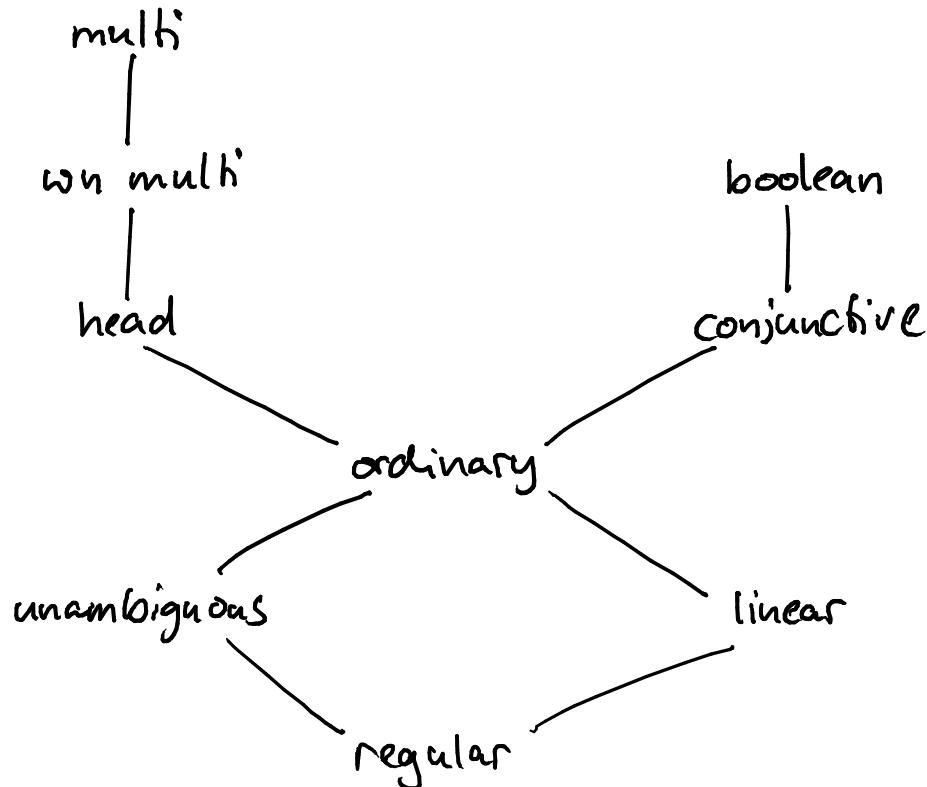
- well-nested M-component grammars:

CON: ( $\leq M$ )-tuples of strings

C-OP: "combination" of string tuples, but unrelated variables  
may not cross each other

- head grammars = well-nested 2-component grammars

# Hierarchy of grammar formalisms (by inclusion)



# Parse trees I

- ordinary grammar

$$\frac{B(u) \quad C(v)}{A(uv)} (A \rightarrow BC)$$

# Parse trees I

- ordinary grammar

$$\frac{B(u) \quad C(v)}{A(uv)} (A \rightarrow BC)$$

example:

$$\frac{\frac{\frac{S(\epsilon)}{(S \rightarrow c)}}{(S \rightarrow ab)} (S \rightarrow SS)}{S(cab)}$$

$(S \rightarrow \epsilon)$   
 $S(\epsilon)$   
 $(S \rightarrow aSb)$   
 $(S \rightarrow SS)$

# Parse trees I

- ordinary grammar

$$\frac{B(u) \quad C(v)}{A(uv)} (A \rightarrow BC)$$

example:

$$\frac{\frac{\frac{S(\epsilon)}{(S \rightarrow c)}}{S(c)} \frac{\frac{S(ab)}{(S \rightarrow aSb)}}{S(ab)}}{S(cab)} (S \rightarrow SS)$$

- conjunctive grammar

$$\frac{B(u) \quad C(v) \quad D(x) \quad E(y)}{A(uv)} (A \rightarrow BC \& DE)$$

# Parse trees I

- ordinary grammar

$$\frac{B(u) \quad C(v)}{A(uv)} (A \rightarrow BC)$$

example:

$$\frac{\overbrace{S(\epsilon)}^{(S \rightarrow \epsilon)} \quad \overbrace{S(c)}^{(S \rightarrow c)}}{\overbrace{S(ab)}^{(S \rightarrow ab)}} \frac{(S \rightarrow aSb)}{(S \rightarrow SS)} S(cab)$$

- conjunctive grammar

$$\frac{B(u) \quad C(v) \quad D(x) \quad E(y)}{A(uv)} (A \rightarrow BC \& DE)$$

example:

$$\frac{\overbrace{A(\epsilon)}^{\text{---}} \quad \overbrace{B(\epsilon)}^{\text{---}} \quad \overbrace{D(\epsilon)}^{\text{---}} \quad \overbrace{C(\epsilon)}^{\text{---}}}{\overbrace{A(a)}^{\text{---}} \quad \overbrace{B(bc)}^{\text{---}} \quad \overbrace{D(ab)}^{\text{---}} \quad \overbrace{C(c)}^{\text{---}}} \frac{}{(S \rightarrow AB \& DC)} S(abc)$$

## Parse trees II

- multi-component

$$\frac{A(t, u) \quad B(v, w)}{S(truw)} (S(x_1 y_1 x_2 y_2) \rightarrow A(x_1, x_2) \ B(y_1, y_2))$$

# Parse trees II

- multi-component

$$\frac{A(t, u) \quad B(v, w)}{S(truw)} (S(x_1 y_1 x_2 y_2) \rightarrow A(x_1, x_2) B(y_1, y_2))$$

example:

$$\begin{array}{c} \overline{A(\varepsilon, \varepsilon)} \\ \overline{A(a, c)} \\ \overline{A(aa, cc)} \\ \hline S(aabcccd) \end{array} \qquad \begin{array}{c} \overline{B(\varepsilon, \varepsilon)} \\ \overline{B(b, d)} \\ \hline (S(x_1 y_1 x_2 y_2) \rightarrow \dots) \end{array}$$

## Language Equations

- ordinary grammars

$$G: S \rightarrow \epsilon$$

$$S \rightarrow aSb$$

$$S \rightarrow SS$$

## Language Equations

- ordinary grammars

$$G: S \rightarrow \epsilon \quad S \rightarrow aSb \quad S \rightarrow SS$$

$$S = \{\epsilon\} \cup (\{a\} \cdot S \cdot \{b\}) \cup (S \cdot S)$$

$L(G)$  is the least fixed point of the equation (system)

## Language Equations

- ordinary grammars

$$G: S \rightarrow \epsilon \quad S \rightarrow aSb \quad S \rightarrow SS$$

$$S = \{\epsilon\} \cup (\{a\} \cdot S \cdot \{b\}) \cup (S \cdot S)$$

$L(G)$  is the least fixed point of the equation (system)

- head/multi-component grammars

→ use sets of tuples of strings as domain

## Language Equations

- ordinary grammars

$$G: S \rightarrow \epsilon \quad S \rightarrow aSb \quad S \rightarrow SS$$

$$S = \{\epsilon\} \cup (\{a\} \cdot S \cdot \{b\}) \cup (S \cdot S)$$

$L(G)$  is the least fixed point of the equation (system)

- head/multi-component grammars

→ use sets of tuples of strings as domain

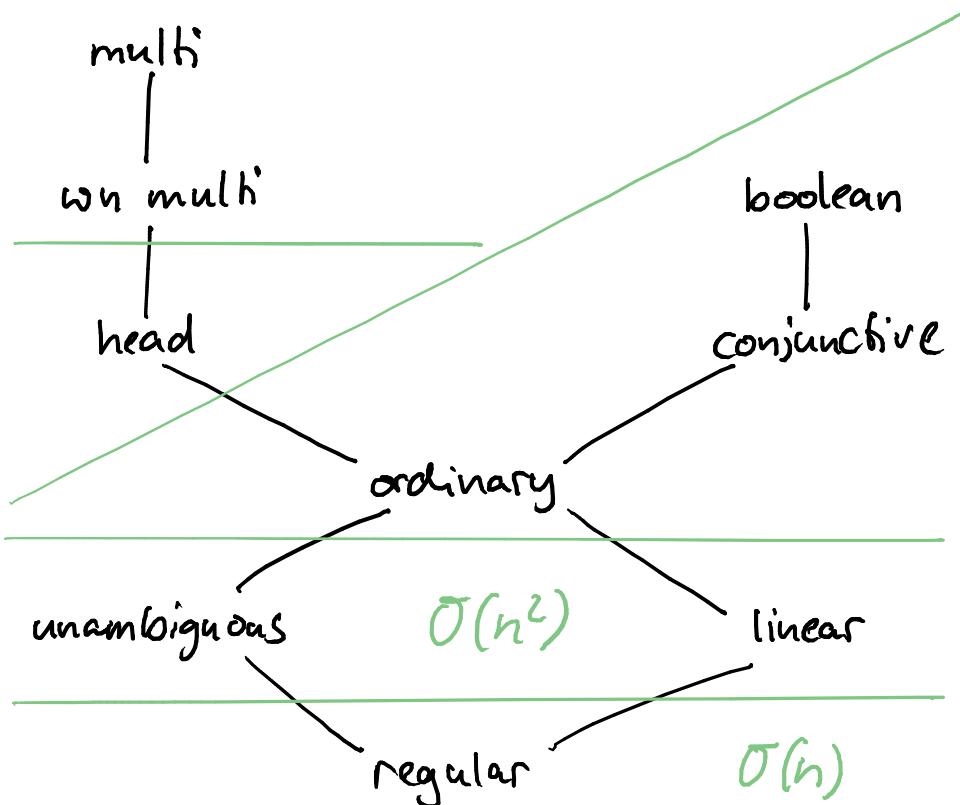
- conjunctive/Boolean grammars

→  $\&$  becomes  $\cap$ ,  $\neg S$  becomes  $\Sigma^* \setminus S$

# Parsing

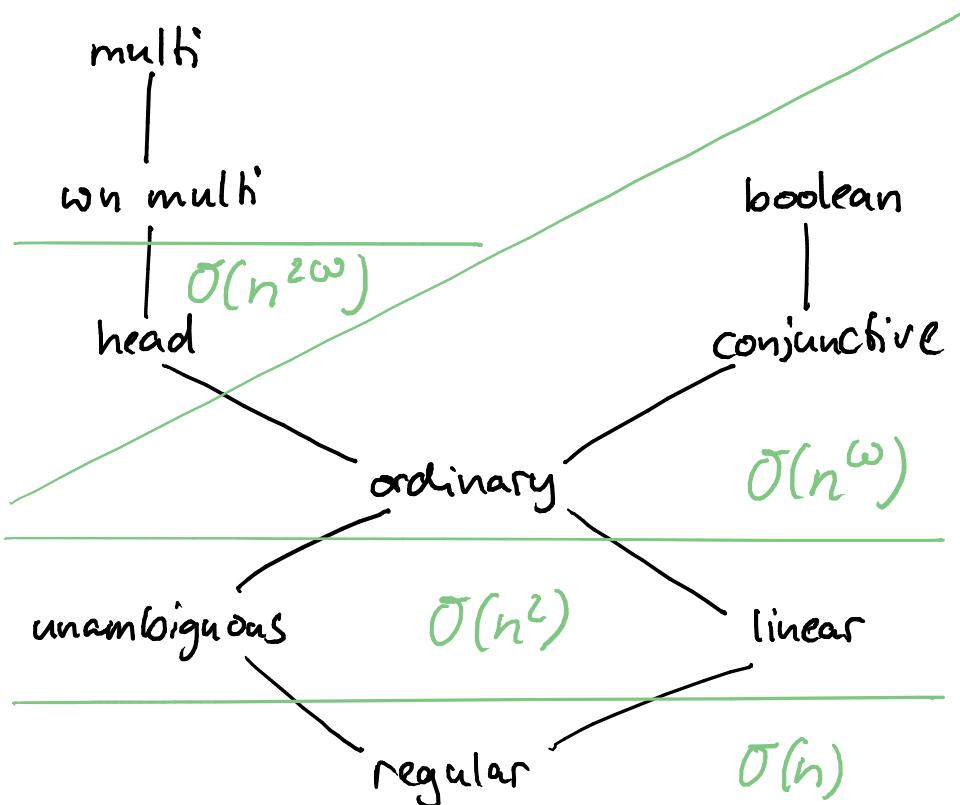
- use dynamic programming algorithm (CKY-like)

# Parsing



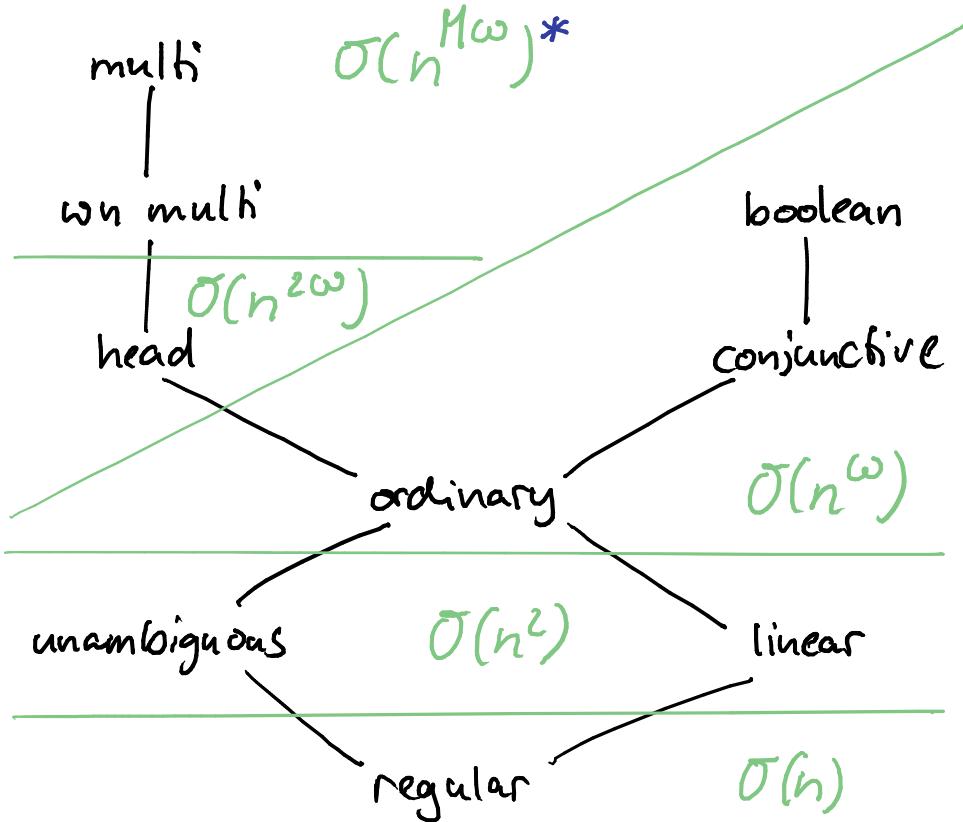
- use dynamic programming algorithm (CKY-like)

# Parsing



- use dynamic programming algorithm (CKY-like)
- two  $n \times n$  matrices can be multiplied in time  $O(n^\omega)$   
[Valiant 75]  
 $\omega \leq 2.38$

# Parsing



- use dynamic programming algorithm (CKY-like)

- two  $n \times n$  matrices can be multiplied in time  $O(n^\omega)$

[Valiant 75]

$\omega \leq 2.38$

\* [Cohen+16]

## Characterisation results (a selection)

- by logic (Büchi results): find  $X$  s.t.  $\nexists G \exists \varphi$  with  
 $L(G) = L(\varphi)$        $\varphi$  is an  $X$ -formula

## Characterisation results (a selection)

- by logic (Büchi results): find  $X$  s.t.  $\forall h \exists \varphi$  with  
 $L(G) = L(\varphi)$        $\varphi$  is an  $X$ -formula
- by automata: find  $Y$  s.t.  $\forall h \exists M$  with  
 $L(G) = L(M)$        $M$  is an  $Y$ -automaton

## Characterisation results (a selection)

- by logic (Büchi result): find  $X$  s.t.  $\forall h \exists \varphi$  with  
 $L(G) = L(\varphi)$        $\varphi$  is an  $X$ -formula
- by automata: find  $Y$  s.t.  $\forall h \exists M$  with  
 $L(G) = L(M)$        $M$  is an  $Y$ -automaton
- homomorphic characterisation (Chomsky-Schützenberger result)  
find  $Z$  s.t.  $\forall h \exists h, R, D$  with  
 $L(G) = h(R \cap D)$        $R \in \text{REG}$ ,  $h \in \text{HOM}$ ,  
     $D$  is a  $Z$ -Dyck language

# References I (in order of appearance)

- [Okhotin 18] A. Okhotin: Underlying Principles and Recurring Ideas of Formal Grammars, 2018. DOI.
- [Chomsky 56] N. Chomsky: Three models for the description of language, 1956. DOI.
- [Okhotin 01] A. Okhotin: Conjunctive Grammars, 2001.
- [Okhotin 04] A. Okhotin: Boolean Grammars, 2004. DOI.
- [Okhotin 13] A. Okhotin: Conjunctive and Boolean grammars, the true general case of the context-free grammars, 2013. DOI.
- [Vijay-Shankar 87] K. Vijay-Shankar, D.J. Weir: Characterizing structural descriptions produced by various grammar formalisms, 1987. DOI.

## References II (in order of appearance)

- [Seki + 91] H. Seki, T. Matsumura, M. Fujii, T. Kasami: On multiple context-free grammars, 1991. DOI.
- [Valiant 75] L.G. Valiant: General context-free recognition in less than cubic time, 1975. DOI.
- [Cohen + 16] S.B. Cohen, D. Gildea: Parsing Linear Context-Free Rewriting Systems with Fast Matrix Multiplication, 2016. DOI.