

Multiple Context-Free Grammars: Basic Properties and Complexity

Hiroyuki Seki

NAIST

MCFG+2

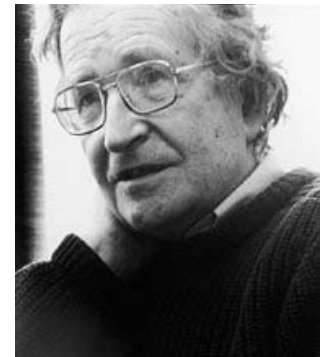
Nara, Sept 9, 2011

Chomsky hierarchy

regular →

context-free(CF) → context-sensitive(CS)

→ phrase structure(PS)



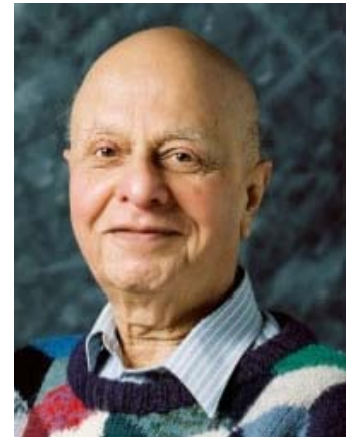
Noam Chomsky (1928-)

Mild context-sensitivity

regular \rightarrow

context-free(CF) \rightarrow context-sensitive(CS)

\rightarrow phrase structure(PS)



Arvind Joshi (1929-)

Macro grammar (Fischer 1968)
Indexed grammar (Aho 1968)
CF tree grammar (Rounds 1970)

Crossed interaction grammar
(Rivas & Eddy 2000)

Tree transducer (Engelfriet+)

**Syntax-Directed
Translation**

**Biological Sequence
Analysis**

1970

1980

1990

2000

...

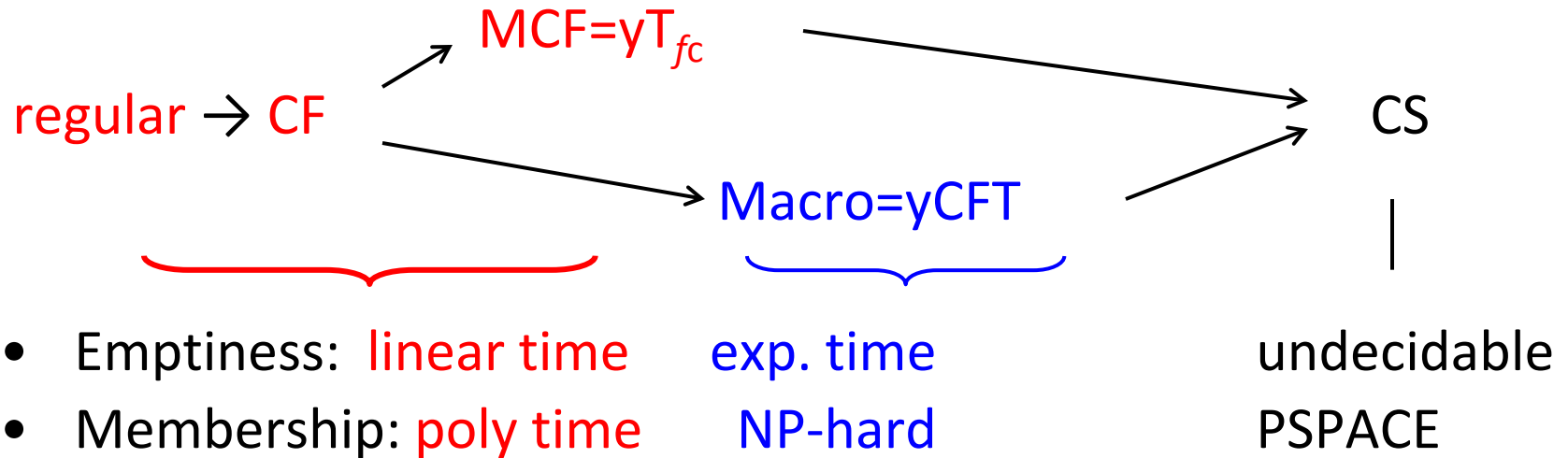
...

Tree Adjoining Grammar
(Joshi, Levy & Takahashi 1975)

**Description of Natural
Language Syntax**

Linear CF rewriting system (Vijay-Shanker, Weir & Joshi 1987)
= Multiple CF grammar (Kasami, Seki & Fujii 1987),

Why MCFG ?



MCF: multiple CF languages

yT_{fc}: finite-copying tree to string transducers

Macro: macro languages

yCFT: yield of context-free tree languages

MCFG

- Kasami, Seki & Fujii, Tech. Rep., Osaka U. 1987, also in IEICE Trans., J71-D-I(5,6), 1988.
- Seki, Matsumura, Fujii & Kasami, TCS 88(2), 1991.
- Seki, Nakanishi, Kaji, Ando & Kasami, 31st ACL, 130-139, 1993.
- Kaji, Nakanishi, Seki & Kasami, Computational Intelligence, 10(4), 440-452, 1994., etc.



Tadao Kasami (1930-2007)

Contents

- Definitions
- Basic properties
- Recognition complexity
- Appendix

CFG as a set of clauses

CFG

DHC (Definite Horn Clause)

Nonterminal
symbol A



Unary predicate A
 $A(x)$ "A can derive x."

rules



clauses

$A \rightarrow B C$

$A(xy) :- B(x), C(y).$

$A \rightarrow a$

$A(a).$

$(B \Rightarrow^* u \text{ and } C \Rightarrow^* v)$

$(B(u) \text{ and } C(v))$

implies $A \Rightarrow^* uv .$

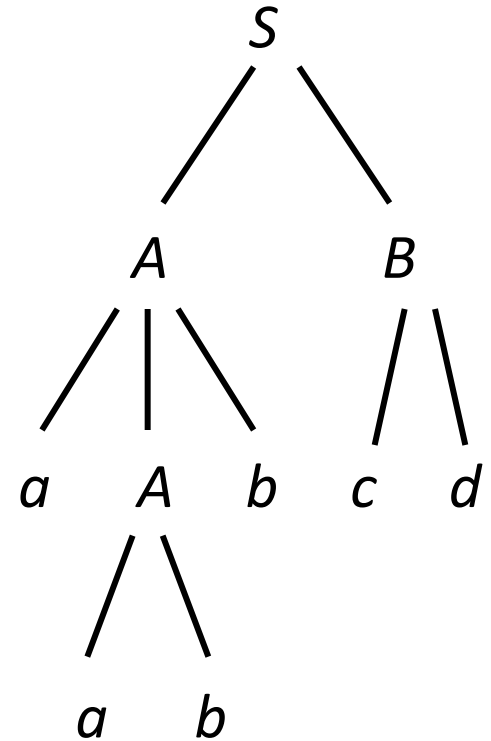
implies $A(uv) .$

Example

$S \rightarrow AB$

$A \rightarrow aAb \mid ab$

$B \rightarrow cBd \mid cd$



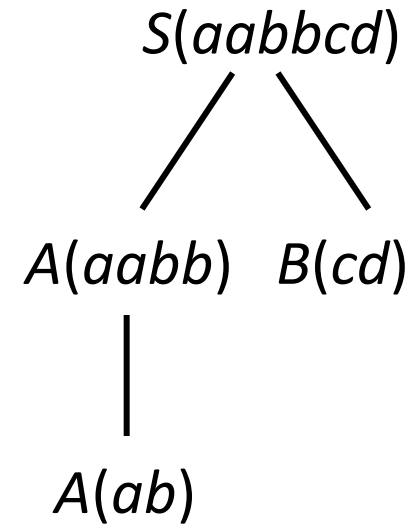
Example

$S(xy) :- A(x), B(y).$

$A(axb) :- A(x). \quad A(ab).$

$B(cxd) :- B(x). \quad B(cd).$

Derivation of $aabbcd$ from S
= Proof of $S(aabbcd).$



From CFG to MCFG

CFG rule

$A \rightarrow BC$

$A \rightarrow a$

DHC (Unary predicates only)

$A(xy) :- B(x), C(y).$

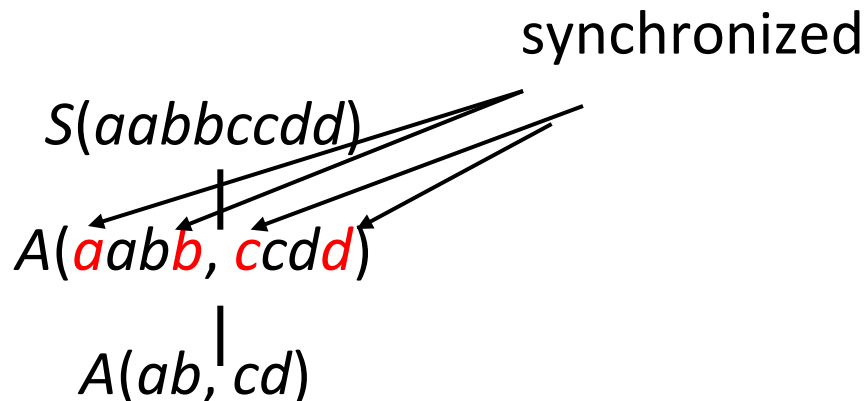
$A(a).$

↓ Extension (arity ≥ 1)

(Ex) $S(xy) :- A(x, y).$

$A(axb, cyd) :- A(x, y).$

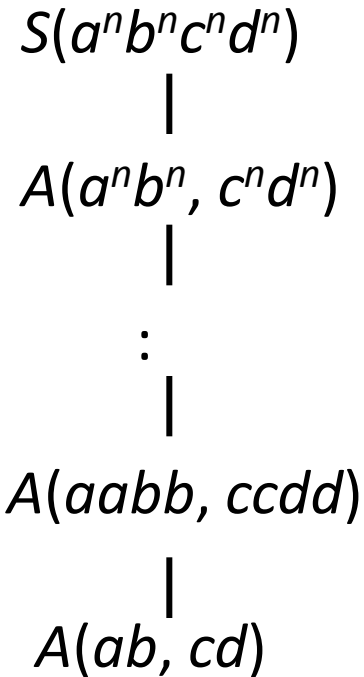
$A(ab, cd).$



Example

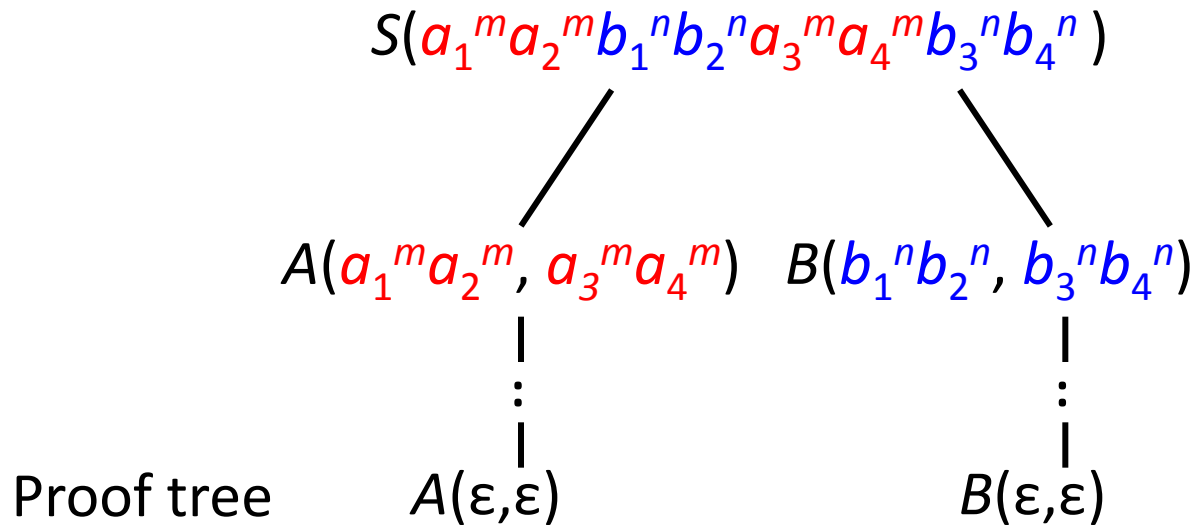
G_1 : $S(xy) :- A(x, y).$
 $A(axb, cyd) :- A(x, y).$
 $A(ab, cd).$

$$L(G_1) = \{ a^n b^n c^n d^n \mid n \geq 1 \}$$



Example

$G_2:$ $S(x_1y_1x_2y_2) :- A(x_1, x_2), B(y_1, y_2).$
 $A(a_1xa_2, a_3ya_4) :- A(x, y). A(\epsilon, \epsilon).$
 $B(b_1xb_2, b_3yb_4) :- B(x, y). B(\epsilon, \epsilon).$



$$L(G_2) = \{ a_1^m a_2^m b_1^n b_2^n a_3^m a_4^m b_3^n b_4^n \mid m, n \geq 0 \}$$

MCFG

$G=(N, T, V, P, S)$

N : predicates (nonterminals), T : terminals, V : variables,
 P : rules, $S \in N$: start predicate

- A rule $\pi \in P$ is a definite Horn clause s.t.
 - Both **head** & **body** are **linear** (w. r. t. V),
 - **Argument** of predicate in **body** is **variable**,
 - **Argument** of predicate in **head** is
string over T and variables in its body.

(Ex) $A(axb, cyd) :- A(x, y).$

$B(x_1y_1, x_2y_2) :- C(x_1, x_2), D(y_1, y_2).$

NG $A(x, ax) :- B(x).$

NG $A(x, y) :- B(x, y), C(x).$

Derivation

Let $G=(N, T, V, P, S)$ be an MCFG.

$A(w_1, \dots, w_{\text{arity}(A)})$ is *derivable* in G

if

$A(\alpha_1, \dots, \alpha_{\text{arity}(A)})$:-

$A_1(\dots), \dots, A_j(x_{j1}, \dots, x_{j \text{arity}(A_j)}), \dots, A_n(\dots) \in P,$

$A_j(u_{j1}, \dots, u_{j \text{arity}(A_j)})$ ($j \in [1..n]$) are derivable in $G,$

$w_i = \theta(\alpha_i)$ ($i \in [1..\text{deg}(A)]$)

where $\theta(x_{jk})=u_{jk}$. ($j \in [1..n], k \in [1..\text{arity}(A_j)]$)

MCFG Language

For an **MCFG** $G=(N, T, V, P, S)$,

$$L(G) := \{ w \in T^* \mid S(w) \text{ is derivable in } G \}$$

is the *multiple context-free language (mcfl)* generated by G .

Non-deleting MCFG

MCFG $G=(N, T, V, P, S)$

Non-deleting: $\forall \pi \in P, \text{Var}_{\text{head}}(\pi)=\text{Var}_{\text{body}}(\pi).$

OK $B(x_1y, x_2) :- C(x_1, x_2), D(y).$

NG $B(y, x_2) :- C(x_1, x_2), D(y).$

Lemma: \forall MCFG $G,$

\exists non-deleting MCFG G' s.t. $L(G') = L(G).$

dim, rank, deg

MCFG $G=(N, T, V, P, S)$

- rule $\pi: A_0(\dots) :- A_1(\dots), A_2(\dots), \dots, A_n(\dots)$.

$\text{dim}(\pi) := \max \{ \text{arity}(A_i) \mid 0 \leq i \leq n \}$, $\text{rank}(\pi) := n$,

$\text{deg}(\pi) := \sum_{0 \leq i \leq n} \text{arity}(A_i)$.

($\text{deg} \leq (\text{rank} + 1) * \text{dim}$)

(Ex)

$\pi_1: A(axb, cyd) :- A(x, y)$.

$\text{dim}(\pi_1)=2$, $\text{rank}(\pi_1)=1$, $\text{deg}(\pi_1)=4$.

$\pi_2: B(x_1y_1, x_2y_2) :- C(x_1, x_2), D(y_1, y_2)$.

$\text{dim}(\pi_2)=2$, $\text{rank}(\pi_2)=2$, $\text{deg}(\pi_2)=6$.

dim, rank, deg (cnt'd)

For MCFG $G=(N, T, V, P, S)$,

$$\{\text{dim, rank, deg}\}(G) := \max_{\pi \in P} \{\text{dim, rank, deg}\}(\pi)$$

- q -MCFG(r): MCFG with $\text{dim} \leq q$ and $\text{rank} \leq r$
- q -MCFG: MCFG with $\text{dim} \leq q$
- MCFG(r): MCFG with $\text{rank} \leq r$

Notation

(Class of) Grammars	CFG	MCFG	q -MCFG(r)	...
(Class of) Languages	CFL	MCFL	q -MCFL(r)	...

Properties of MCFG

(Also see [Vijay-Shanker, Weir & Joshi 1987])

- (Generative power) $\text{CFL} \subsetneq \text{MCFL} \subsetneq \text{CSL}$.
- Every MCFL is semilinear.
- (Closure property)
 - full AFL.
 - Not closed under intersection.
- (Decidability)
 - Emptiness ($L(G)=\emptyset?$): $O(|G|)$ -time decidable.
 - Recognition ($w \in L(G)?$): poly-time decidable.
 - Inclusion ($L(G_1) \subseteq L(G_2)?$): Undecidable.

full AFL : a class of languages closed by homomorphism, inv. homomorphism, intersection with regular sets, union, Kleene closure

Proof example (exercise)

(Closure under intersection with regular language.)

CFL case:

For given CFG in Chomsky normal form $G=(N, T, P, S)$ &

FA (finite automaton) $M=(Q, T, \delta, p_0, F)$

(Q : set of states, δ : transition function, p_0 : initial state,
 F : set of final states),

construct CFG $G'=(N \times Q \times Q \cup \{S'\}, T, P', S')$ where

$\cap R$ Closure: CFG Case

- $\forall A \rightarrow BC \in P, \forall q_1, q_2, q_3 \in Q :$
 $A[q_1, q_3] \rightarrow B[q_1, q_2] C[q_2, q_3] \in P'$
- $\forall A \rightarrow a \in P, \forall q_1, q_2 \in Q$ s.t. $\delta(q_1, a) = q_2 :$
 $A[q_1, q_2] \rightarrow a \in P'$

Correctness:

$$A[q_1, q_2] \Rightarrow_{G'}^* w \quad \text{iff} \\ A \Rightarrow_G^* w \quad \text{and} \quad \delta^*(q_1, w) = q_2$$

$\cap R$ Closure: MCFG Case (1/2)

$$A(x_1y_1, y_2x_2) :- B(x_1, x_2), C(y_1, y_2) \in P$$

$$A[?, ?, ?, ?] (x_1y_1, y_2x_2) :-$$

$$B[q_1, q_2, q_3, q_4] (x_1, x_2), C[r_1, r_2, r_3, r_4] (y_1, y_2) \in P'$$

$$A[q_1, r_2, r_3, q_4] (x_1y_1, y_2x_2) :-$$

$$B[q_1, q_2, q_3, q_4] (x_1, x_2), C[q_2, r_2, r_3, q_3] (y_1, y_2) \in P'$$



$\cap R$ Closure: MCFG Case (2/2)

Correctness:

$A[q_{1s}, q_{1e}, \dots, q_{ns}, q_{ne}](w_1, \dots, w_n)$ provable in G' iff

$A(w_1, \dots, w_n)$ derivable in G

and

$\delta^*(q_{js}, w) = q_{je} \ (j \in [1..n])$

Properties of q -MCFG(r)

- q -MCFL(r), MCFL(r), q -MCFL ($r \geq 2, q \geq 1$):
substitution closed full AFL.
- q -MCFL = $\text{yT}_{\text{FIN}(q)}$ [Weir92]
(yield of tree transducers with copy bound q)
- q -MCFL(1) = $\text{ETOL}_{\text{FIN}(q)}$
(DTOL with copy bound q)
- $\text{TAL} \not\subseteq 2\text{-MCFL}(2)$

Pumping Lemma for CFL

$\forall L \in \text{CFL} \exists n \geq 1 \forall z \in L (|z| \geq n)$

\exists partition $z = uvwxy$ such that $|vx| \geq 1$

$\forall i \geq 0 : z_i = uv^iwx^iy \in L.$

Pumping Lemma for MCFL

For **general** MCFL :

$$\forall L \in q\text{-MCFL} \exists n \geq 1 \exists z \in L (|z| \geq n)$$

$$\exists \text{ partition } z = u_1 v_1 w_1 s_1 u_2 \dots u_q v_q w_q s_q u_{q+1}$$

$$\sum |v_j s_j| \geq 1$$

$$\forall i \geq 0 : z_i = u_1 v_1^i w_1 s_1^i u_2 \dots u_q v_q^i w_q s_q^i u_{q+1} \in L.$$

For **well-nested** MCFL [Ka09] :

$$\forall L \in q\text{-wnMCFL} \exists n \geq 1 \forall z \in L (|z| \geq n)$$

$$\exists \text{ partition } z = \dots \text{ (same as above) } \dots$$

Recently, it was proved that the strong version does not hold for general MCFL.

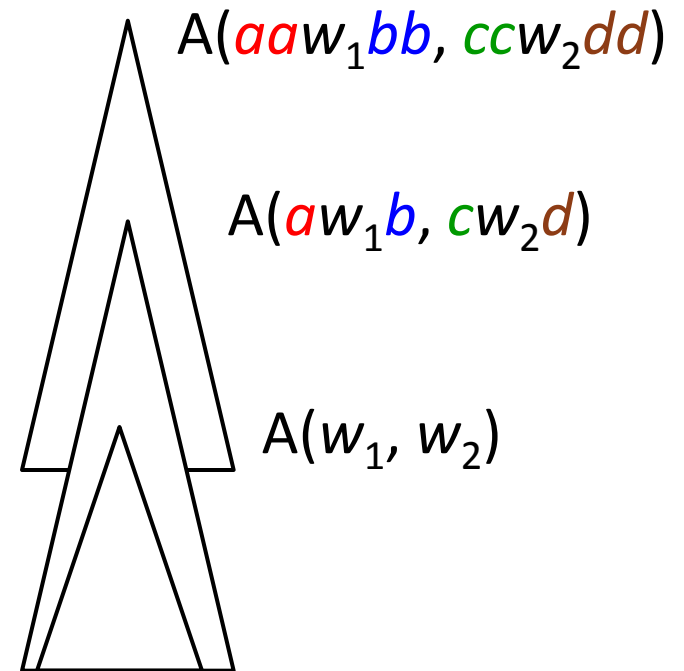
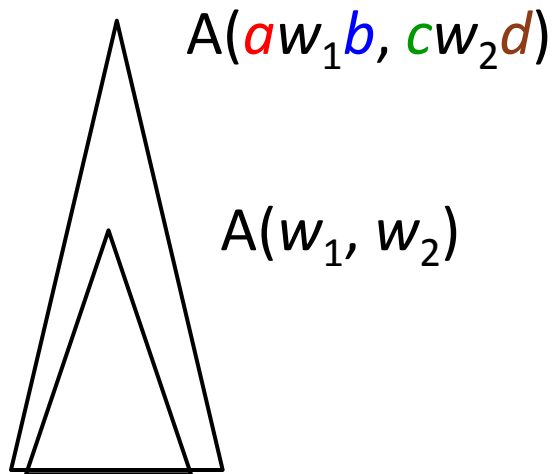
[Ka09]

context-free languages, DLT09, LNCS 5583.

Pumping lemma

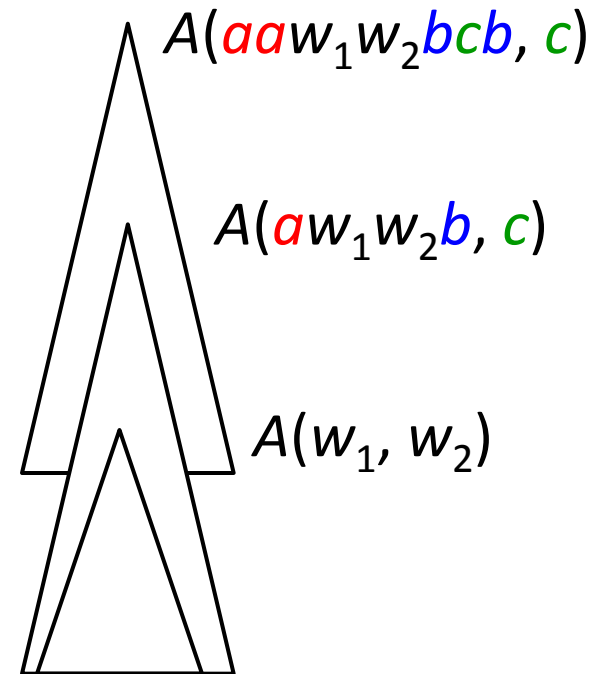
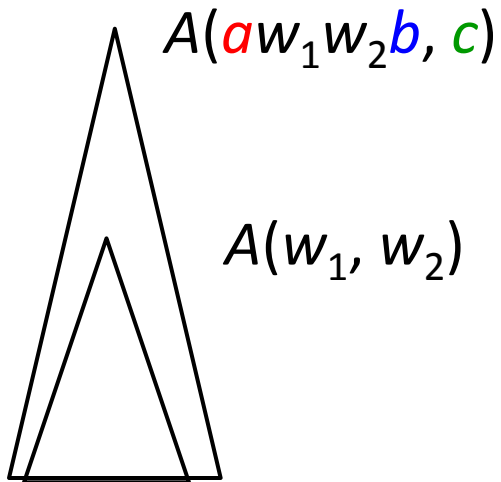
An easy case.

$$u_1 \underline{v_1^i} w_1 \underline{s_1^i} u_2 \underline{v_2^i} w_2 \underline{s_2^i} u_3 \in L$$



Pumping lemma

$$\underline{u_1 v_1^i} \underline{w_1 s_1^i} \underline{u_2 v_2^i} \underline{w_2 s_2^i} u_3 \in L ??$$



Hierarchy on dimensions

rank $r \geq 1$

dim
 $q = 1$

CFL	$\{ a_1^n a_2^n \mid n \geq 0 \}$
-----	-----------------------------------

2

$\{ a_1^n a_2^n a_3^n a_4^n \mid n \geq 0 \}$

3

$\{ a_1^n a_2^n a_3^n a_4^n a_5^n a_6^n \mid n \geq 0 \}$

4

...

5

--

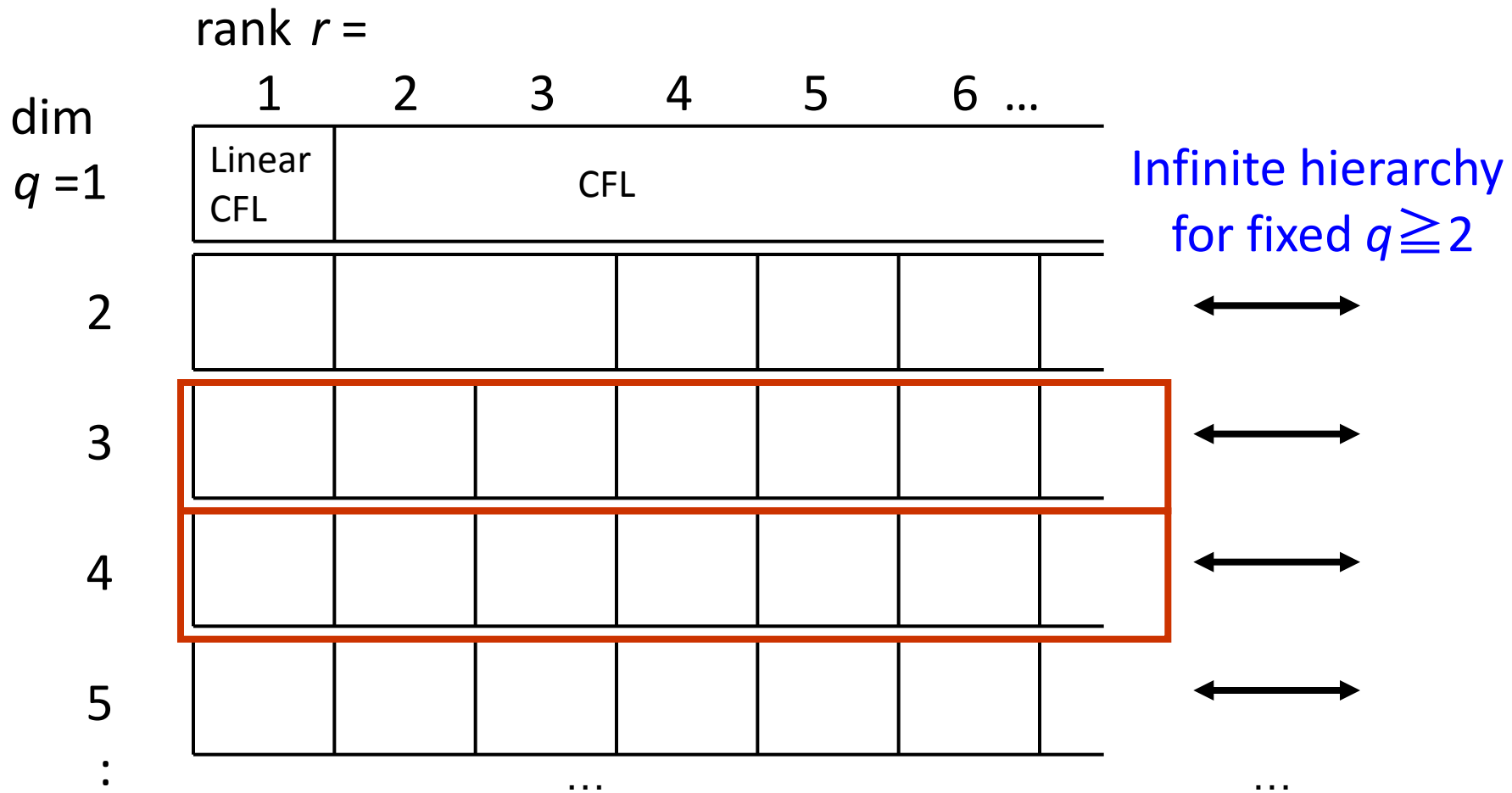
:

--

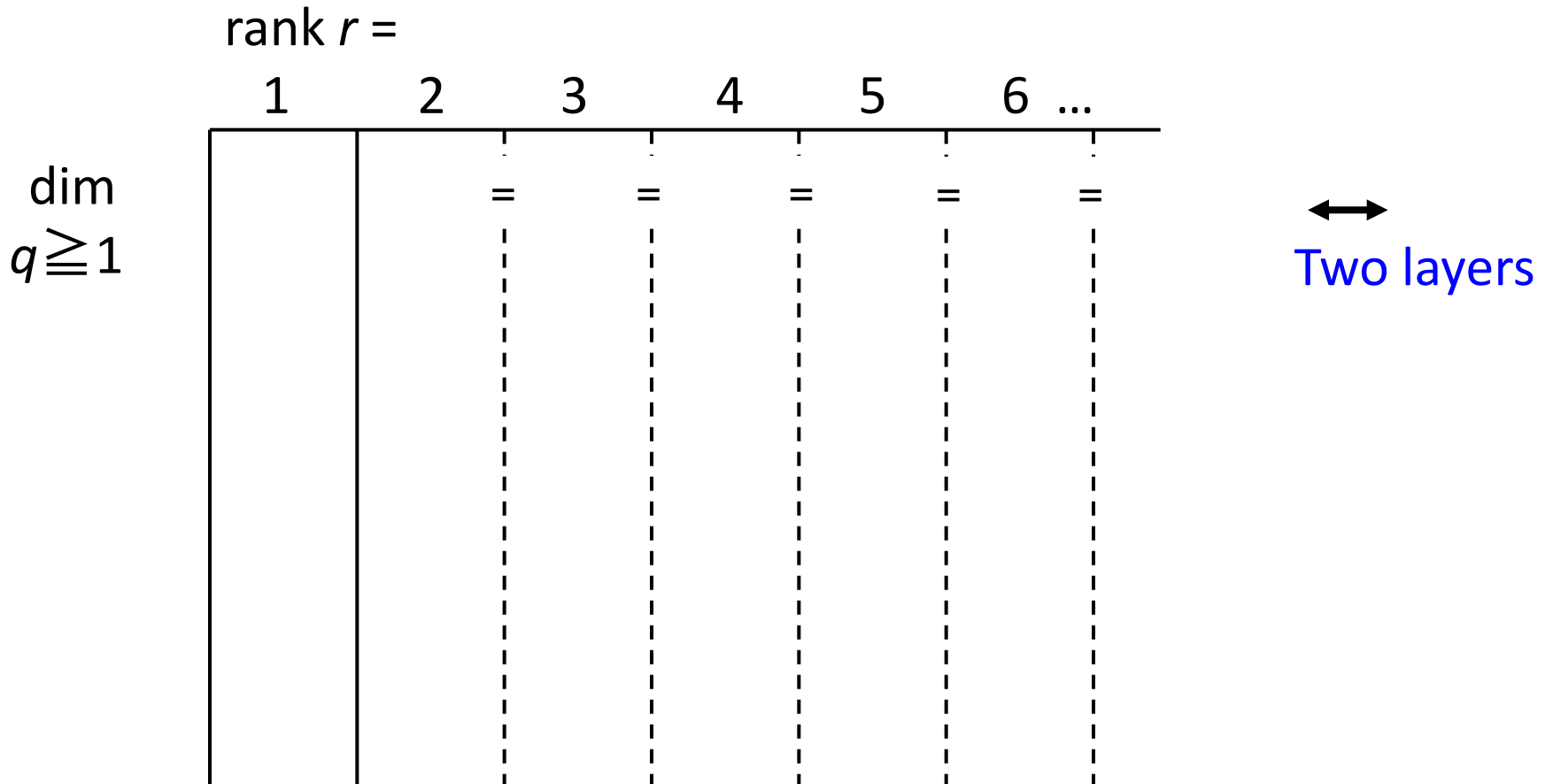
Infinite
hierarchy



Hierarchy on ranks for fixed dim [RS99]

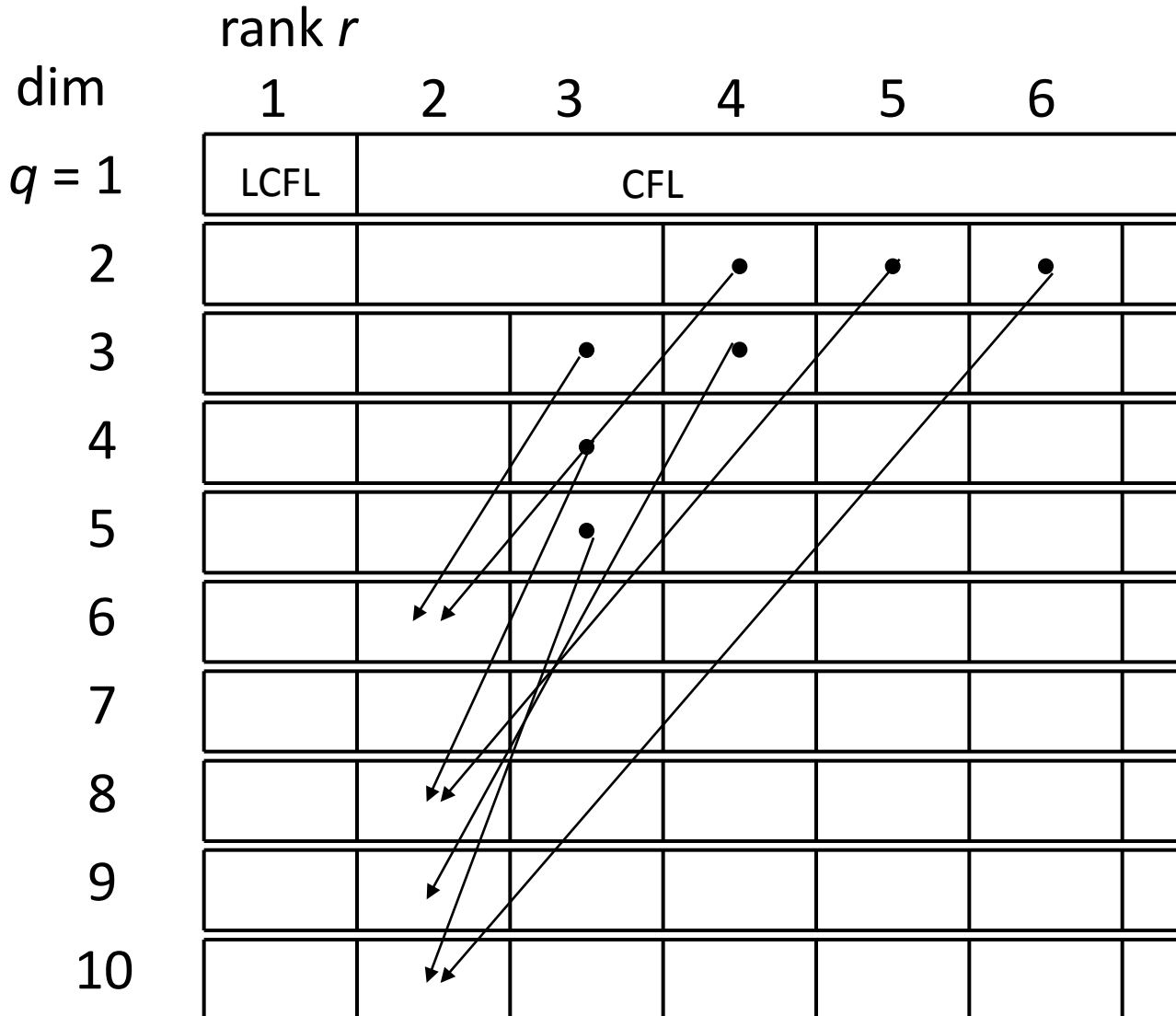


Hierarchy on ranks [RS99][Matsumura89]



MCFL=MCFL(2) \Rightarrow It suffices to consider
MCFL(1) and MCFL(2) as far as ranks are concerned.

Trade-offs [RS99]



$$q\text{-MCFL}(r) \subseteq (k+1)q\text{-MCFL}(r-k)$$

Corollary:

$$q\text{-MCFL}(r) \subseteq (r-1)q\text{-MCFL}(2)$$

Beyond MCFG

$A(x, x) :- B(x), C(x).$

$A(x) :- B(x), C(x).$ Closed under intersection

$A(x, ax) :- B(x).$

$B(x_1y_1, x_2y_2) :- C(x_1, x_2), D(y_1, y_2).$

} MCFG

} PMCFG

} simpleLMG
[Groenink95]

$MCFL \subsetneq PMCFL \subsetneq simpleLML = P$

$\{a^{2^n} \mid n \geq 0\} \in PMCFL - MCFL$

$S(xx) :- S(x), S(a).$

Other topics on MCFG

(Also refer to the other talks.)

- Parsing
 - Earley type parser (Matsumura+89), (Kanazawa 08)
 - ‘Unambiguous’ MCFG ($O(n^2)$ time recognizable) (Nakanishi92)
 - LL(k) MCFG (Nii96)
 - Stochastic CYK parser (Kato+06)
- An extension of Chomsky-Schutzenberger theorem for CFL (Kaji+91, Yoshinaka+10)

Fixed vs. universal recognition

Recognition problem for fixed language L (FRP)

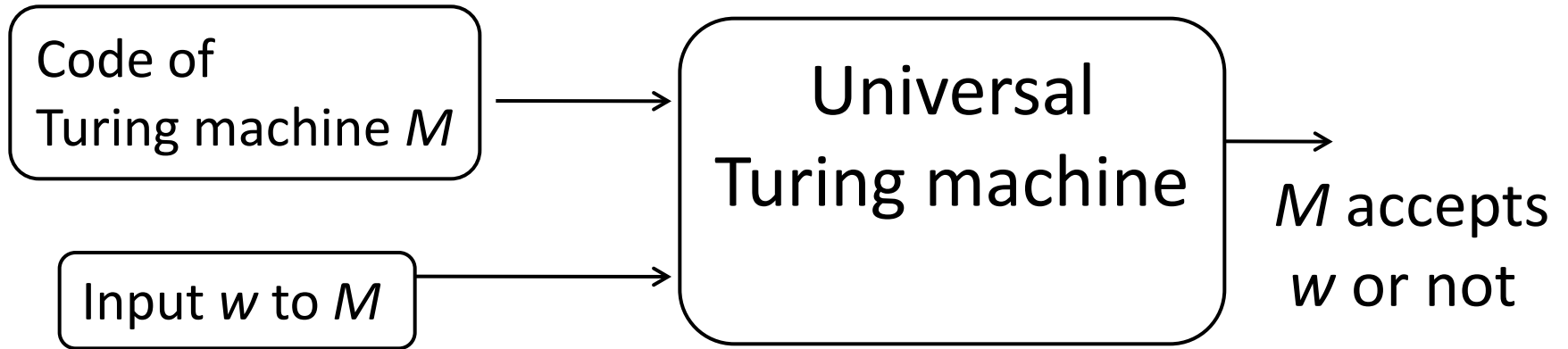
Input: string w

Problem: $w \in L$?

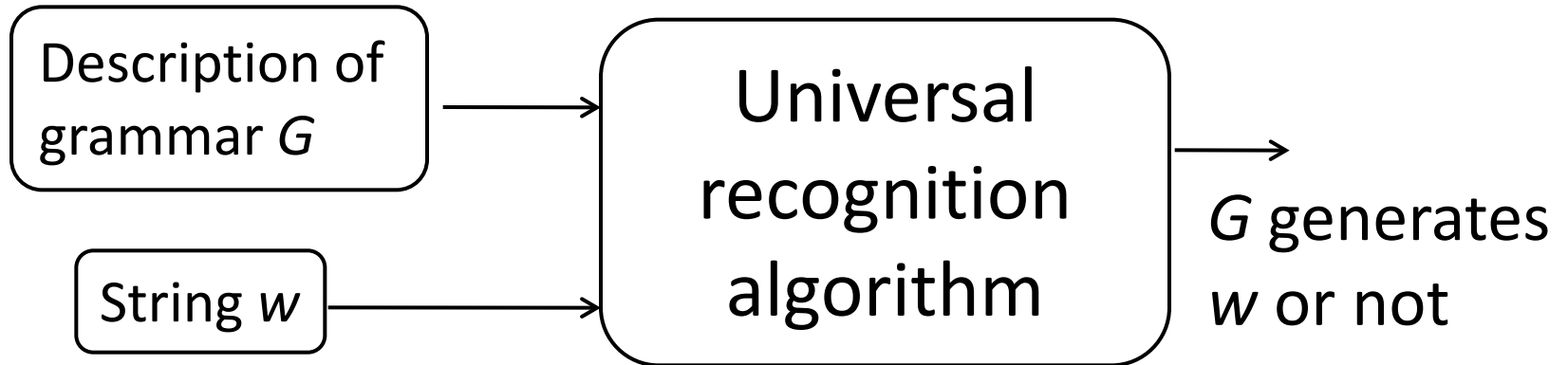
Universal recognition problem (URP)

Input: grammar G , string w

Problem: $w \in L(G)$?



URP



Example

Recognition problem for fixed language L :

For a string w , decide $w \in L$.

$O(|w|^3)$ time when L is CFL.

Universal recognition problem (URP):

For a grammar G & string w , decide $w \in L(G)$.

$O(|G| |w|^3)$ time when G is CFG.

DEXP time-hard when G is GPSG[†].

[†] GPSG (Generalized Phrase Structure Grammar^[GKPS85])

GPSG = CFG in generative power

⇒ Generative power ≠ Complexity of URP

FRP vs. URP

- Generative power can be measured by
Complexity of FRP
- Generative power \neq Complexity of URP
Rather, complexity of URP for grammar class C
 \sim Succinctness of C

Grammar (especially MCFG) as computation device
with bounded resource

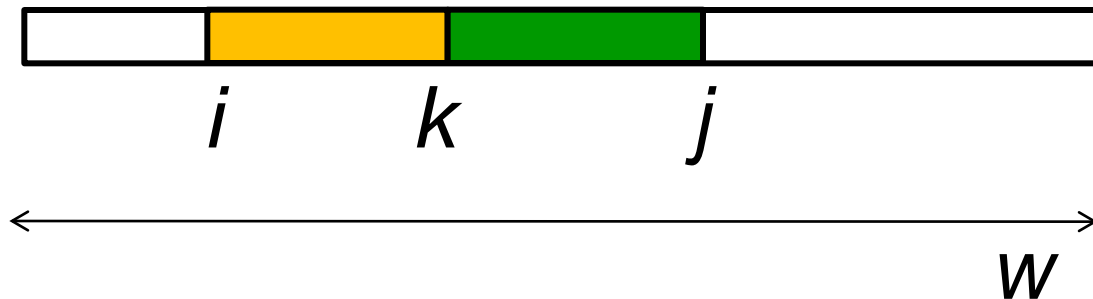
Let's start with FRP.

CYK algorithm for CFG

If rule: $A(xy) :- B(x), C(y)$ ($A \rightarrow BC$)

$B \in \text{table}(i, k) \wedge C \in \text{table}(k, j)$ for $\exists k (i \leq k < j)$

then add A to $\text{table}(i, j)$;



Time complexity: $O(n^3)$ $n = \text{length of input } w$

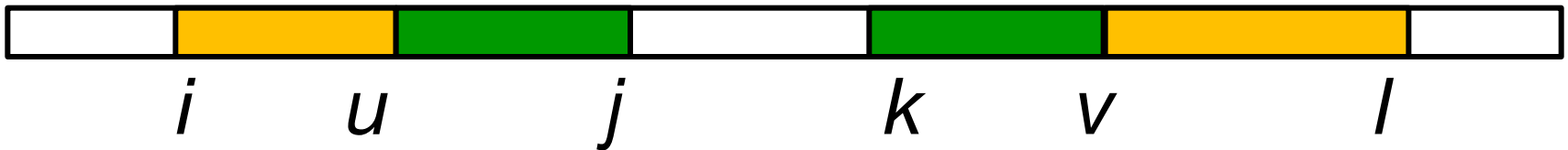
CYK algorithm for MCFG

If rule: $A(y_1 z_1, z_2 y_2) :- B(y_1, y_2), C(z_1, z_2),$

$B \in \text{table}(i, u, v+1, l) \wedge C \in \text{table}(u+1, j, k, v)$

for $\exists u, v (i \leq u < j, k \leq v < l)$

then add A to $\text{table}(i, j, k, l)$;



- Time complexity: $O(n^{\deg(G)})$ ($\deg(G) \leq (r+1)q$) ^{W}
($\deg(G)=6$ in this example)

URP for MCFG

- General case DEXP-complete
- Non-deleting MCFG PSPACE-complete
- q -MCFG with fixed q NP-complete
- q -MCFG(r) with fixed q, r P-complete

General case

- DEXP-complete
 - Upper bound: Simulation of MCFG by DEXP time bounded Turing machine.
 - Lower bound: Simulation of PSPACE bounded *alternating* Turing machine (APSPACE) by MCFG.
Since general MCFG can use deleting rules, it can simulate computation for a universal state (of alternating TM) with poly-length ‘sentential form.’

Alternating Turing machine

ATM is $M = (Q, \Sigma, \Gamma, B, \delta, q_s, Q_F, Q_U, Q_E)$

$Q = Q_U \cup Q_E \cup Q_F \cup \{q_s\}$: state set

$\Sigma (\subseteq \Gamma - \{B\})$: input symbols

Γ : tape symbols B : blank symbol

$\delta \subseteq (Q \times \Gamma) \times (Q \times \Gamma \times \{\rightarrow, \leftarrow\})$: transition relation

q_s : initial state Q_F : final (accepting) states

Q_U : set of universal states

Q_E : set of existential states

Alternating Turing machine

ATM $M = (Q, \Sigma, \Gamma, B, \delta, q_s, Q_F, Q_U, Q_E)$

ID (instantaneous description) is $\alpha q \beta$ where $q \in Q$ & $\alpha, \beta \in \Gamma^*$.

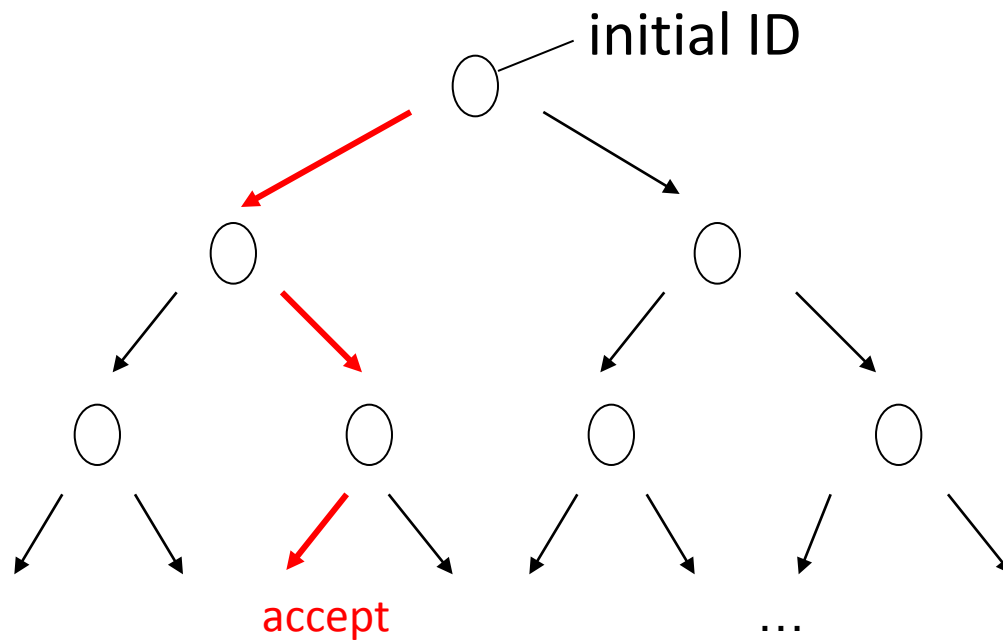
Move relation $\alpha q \beta \rightarrow^*_M \alpha' q' \beta'$ over IDs is defined in a usual way.

For an input w , ID $q_s w$ is the initial ID for w .

ID $\alpha q \beta$ is accepting if $q \in Q_F$.

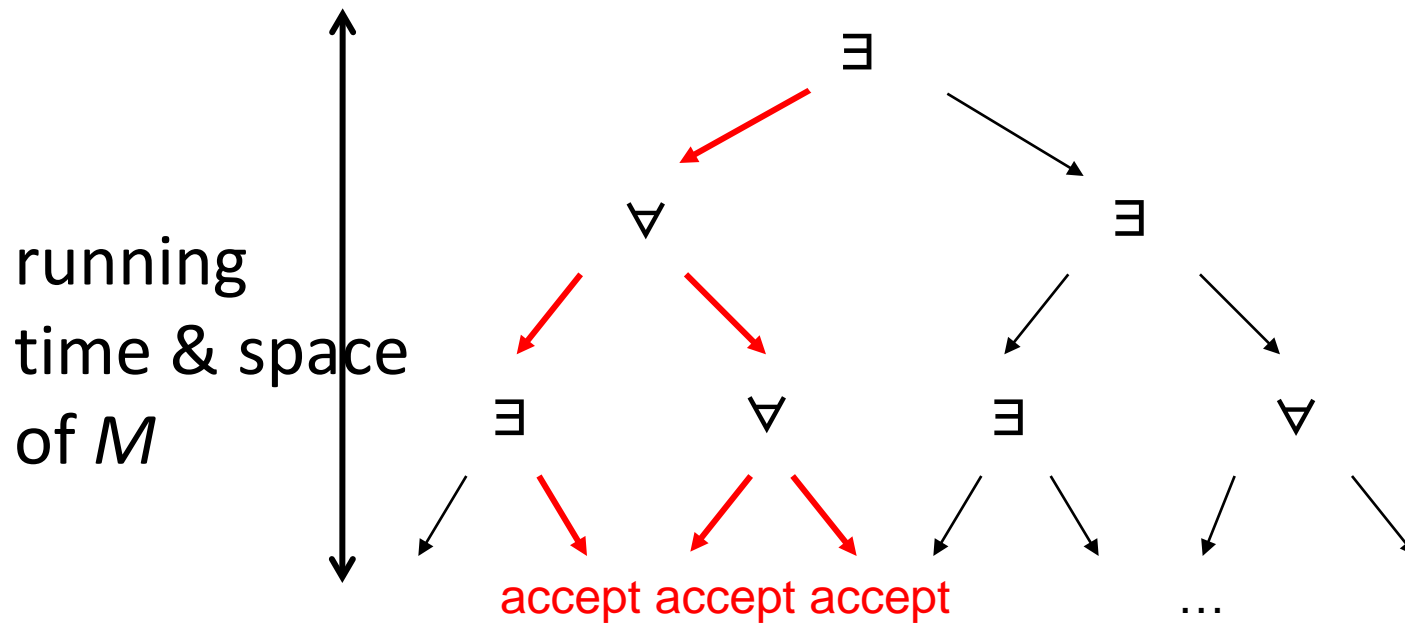
Nondeterministic Turing machine

Input w is accepted if there is a run that reaches an accepting ID from the initial ID for w .



Accepting language

$$L(M) = \{ w \mid \text{Initial ID } q_s w \text{ is accepting} \}$$



TM and ATM

Proposition^[CS80]:

PSPACE=APTIME, DEXPTIME=APSPACE

APTIME := class of problems solvable by poly time-bounded ATM

APSPACE := class of problems solvable by poly space-bounded ATM

Proof of DEXP completeness

For given $p(n)$ -space bounded ATM $M = (Q, \Sigma, \Gamma, B, \delta, q_s, Q_F, Q_U, Q_E)$ and $w \in \Sigma^*$,
construct MCFG $G (N, \{1\}, V, P, S)$ such that
 M accepts w if and only if $\varepsilon \in L(G)$.

APSPACE \leq URP MCFG

- Nonterminals $A[q, k]$ ($q \in Q, k \in [1..p(n)]$)

$$\dim(\cdot) = p(n) \mid \Gamma$$

$$A[q, k] \left[\begin{array}{cccc} 1 & \dots & j & \dots & p(n) \\ & & \varepsilon & & \end{array} \right] c$$

derivable intuitively means that

$\exists \alpha, \beta: \alpha q \beta$ is accepting,

$k = |\alpha| + 1$, the j -th symbol of $\alpha \beta$ is c .

Construction

(0 step acceptance)

$$\forall q_f \in Q_F, \forall k \in [1..p(n)]$$

$$A[q_f, k](\varepsilon, \dots, \varepsilon) .$$

(Initial ID)

$$\forall q \in Q, \forall k \in [1..p(n)]$$

$$S(x_{\langle 1, a_1 \rangle} x_{\langle 2, a_2 \rangle} \dots x_{\langle n, a_n \rangle} x_{\langle n+1, B \rangle} \dots x_{\langle p(n), B \rangle})$$

$$:- A[q, k](\dots)$$

for given input $w = a_1 a_2 \dots a_n$.

Construction (existential state)

- $q \in Q_E$, $(b, p, \rightarrow) \in \delta(q, a)$ (move right)

$A[q, 2] (x_{\langle 1, a \rangle}, x_{\langle 2, b \rangle}, x_{\langle 3, a \rangle},$

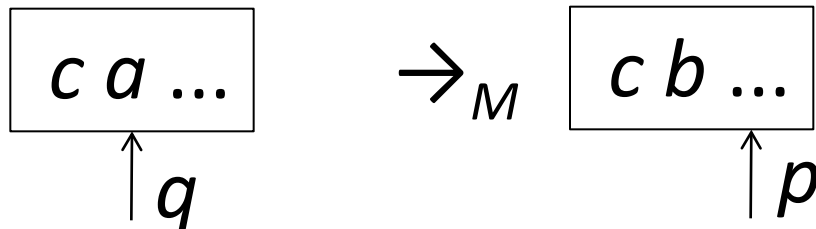
$x_{\langle 1, b \rangle}, \mathbf{1}, x_{\langle 3, b \rangle},$

$x_{\langle 1, c \rangle}, \mathbf{1}, x_{\langle 3, c \rangle}) :-$

$A[p, 3] (x_{\langle 1, a \rangle}, x_{\langle 2, a \rangle}, x_{\langle 3, a \rangle},$

$x_{\langle 1, b \rangle}, x_{\langle 2, b \rangle}, x_{\langle 3, b \rangle},$

$x_{\langle 1, c \rangle}, x_{\langle 2, c \rangle}, x_{\langle 3, c \rangle})$



Construction (universal state)

- $q \in Q_U$, , $(b, p, \rightarrow) \in \delta(q, a)$, $(c, r, \leftarrow) \in \delta(q, a)$

$$A[q, 2] (x_{\langle 1, a \rangle} y_{\langle 1, a \rangle}, x_{\langle 2, b \rangle} y_{\langle 2, c \rangle}, x_{\langle 3, a \rangle} y_{\langle 3, a \rangle}, \\ x_{\langle 1, b \rangle} y_{\langle 1, b \rangle}, \mathbf{1}, x_{\langle 3, b \rangle} y_{\langle 3, b \rangle}, \\ x_{\langle 1, c \rangle} y_{\langle 1, c \rangle}, \mathbf{1}, x_{\langle 3, c \rangle} y_{\langle 3, c \rangle}) :-$$

$$A[p, 3] (x_{\langle 1, a \rangle}, x_{\langle 2, a \rangle}, x_{\langle 3, a \rangle}, \\ x_{\langle 1, b \rangle}, x_{\langle 2, b \rangle}, x_{\langle 3, b \rangle}, \\ x_{\langle 1, c \rangle}, x_{\langle 2, c \rangle}, x_{\langle 3, c \rangle}),$$

$$A[r, 3] (y_{\langle 1, a \rangle}, y_{\langle 2, a \rangle}, y_{\langle 3, a \rangle}, \\ y_{\langle 1, b \rangle}, y_{\langle 2, b \rangle}, y_{\langle 3, b \rangle}, \\ y_{\langle 1, c \rangle}, y_{\langle 2, c \rangle}, y_{\langle 3, c \rangle}),$$

Recent result

Theorem: URP for poly depth-bounded MCFG is PSPACE-complete.

Corollary: URP for poly depth-bounded MCFG(1) is NP-complete.

Summary

non-deletion	dimension	depth	Universal recognition	
-	-	-	DEXP-complete	[KNSK94]
required	-	-	PSPACE-complete	
-	bounded	-	NP-complete	
-	-	poly	PSPACE-complete	[Seki10]
-	-	poly, ($i-1$)-bounded alternation	Σ_i/π_i -complete	
-	-	poly, rank 1	NP-complete	

Appendix

Tadao Kasami Wins the 1999 Claude E. Shannon Award

The Information Theory Society's highest honor, the Claude E. Shannon Award, is awarded annually to an individual who has achieved consistent and profound contributions to the field of information theory. The recipient is chosen by a selection committee consisting of Society officers and two former Shannon Award recipients.

Prof. Tadao Kasami of Hiroshima City University, Japan, has been selected as the 1999 Claude E. Shannon Award recipient. The award was announced at the 1998 International Symposium on Information Theory and will be presented to Prof. Kasami at the 2000 International Symposium on Information Theory.

Hideki Imai held the following interview with Tadao Kasami in honor of his receipt of the Award.

Interview with Tadao Kasami

Hideki: Congratulations on winning 1999 Claude E. Shannon Award. You certainly deserve this award on the basis of your outstanding contributions in coding theory. I know that you started your career as a re-



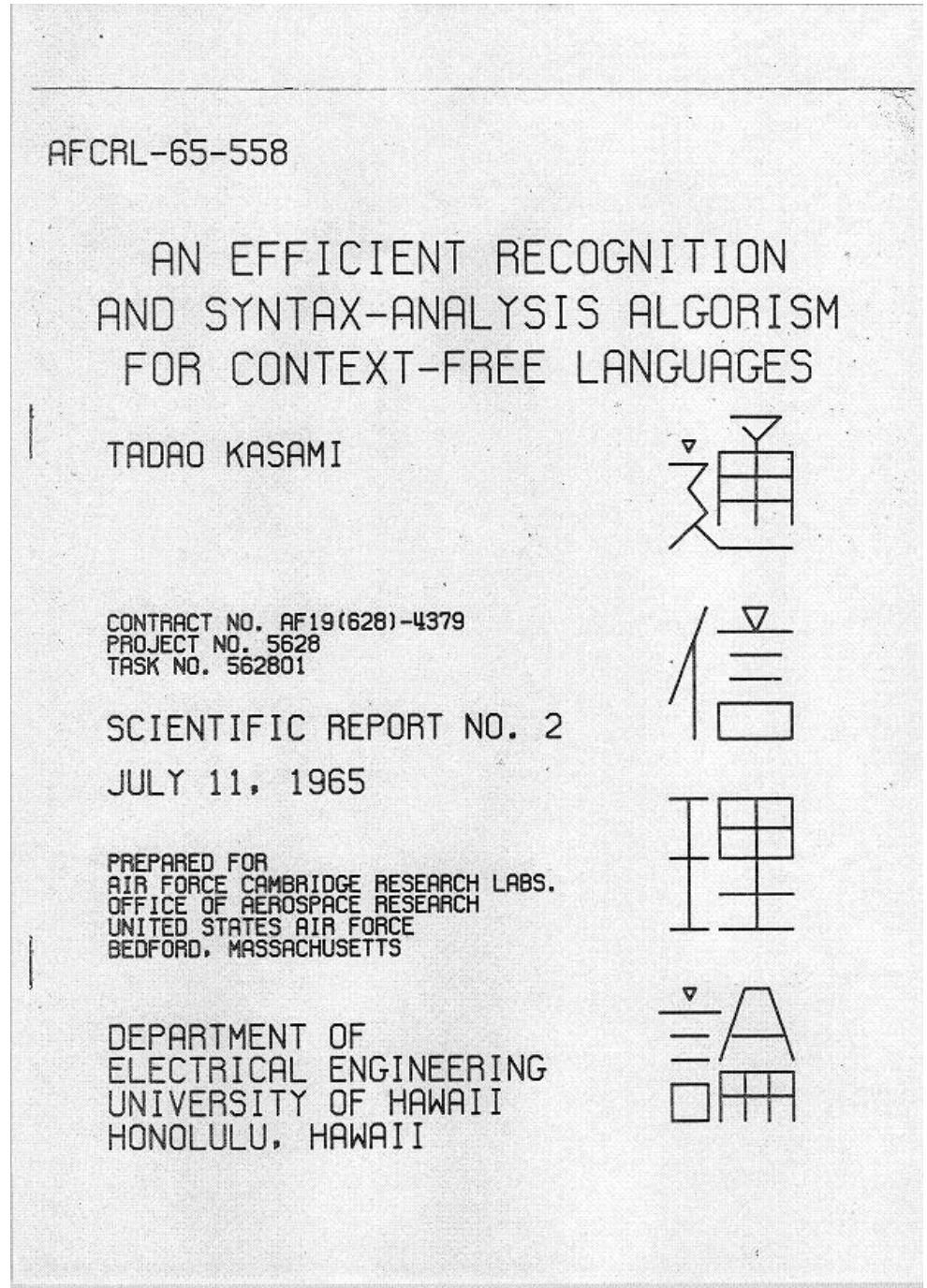
Tadao Kasami

theory in which there remained only very hard long-standing problems. Almost everyday, I scanned papers in IRE Transactions and so on at the library. For instance, I found a paper [2] by N. Abramson and read it with a strong interest. There continued a rush of epoch-making papers by Bose-Chaudhuri, Peterson and so on. I had to study hard in order to catch up with the progress of coding theory. Finally, I finished my Ph. D. thesis (1963) based on published papers [3-8].

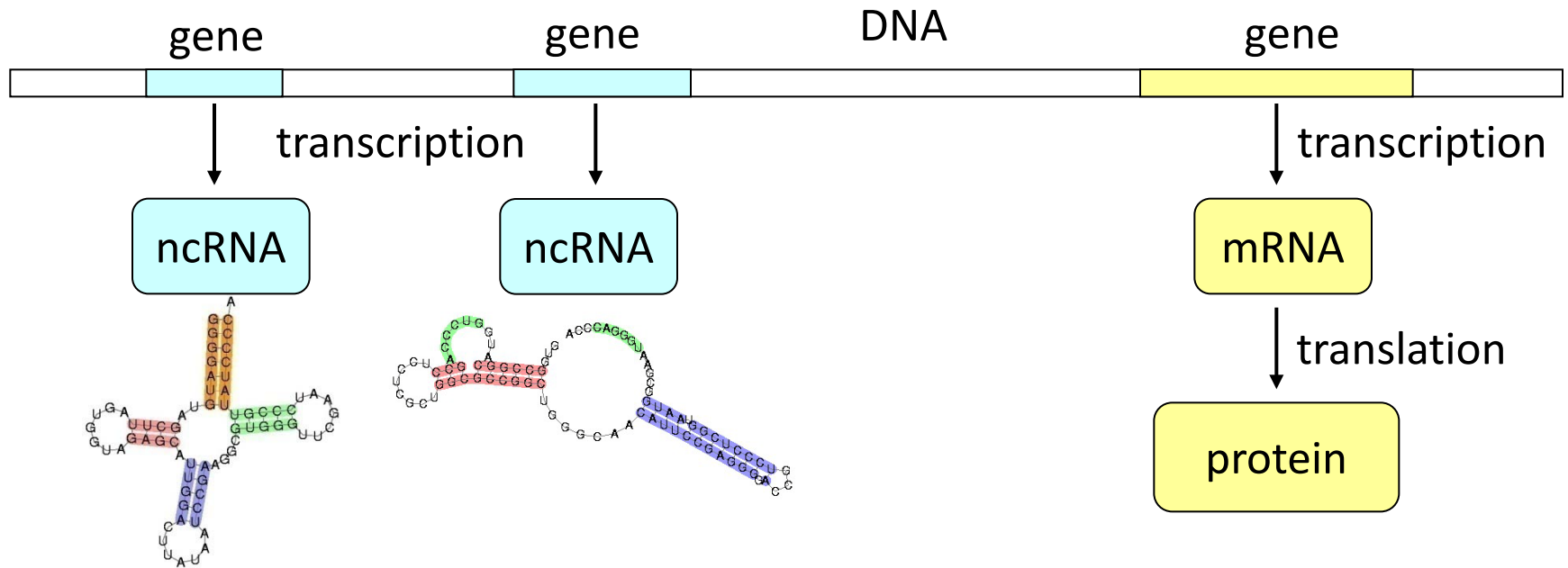
Hideki: How did you meet Prof. W. Wesley Peterson and what did you learn from him?

Tadao: I exchanged letters with Prof. Peterson (Wes) who had been an associate editor of IEEE Trans. on Information Theory. I first had the chance to meet him in 1963 when he visited Japan to attend a URSI meeting in Tokyo. In November 1964, my sincere hope to study under Wes came true by his kind arrangements. I was deeply impressed with his profound knowledge which ran extensively from quantum communication to software engineering. He encouraged me to continue my study [9] on a syntax analysis algorithm for context-free grammar. Wes advised me to select a research subject of main interest in the field. When

- Original paper on CYK algorithm (July, 1965).
- Dr. S. Eddy, a system biologist, Washington U. asked Dr. Kasami to send a copy (Apr. 2002).
- Dr. Kasami searched his house for a week, found one at last and sent it to Dr. Eddy.



noncoding RNA



- **ncRNA (noncoding RNA)**

- RNA not translated to protein, some of them play important roles in bio-chemical reactions (functional RNA).

Structures of ncRNA

- Primary structure: sequence of bases A, U, C, G

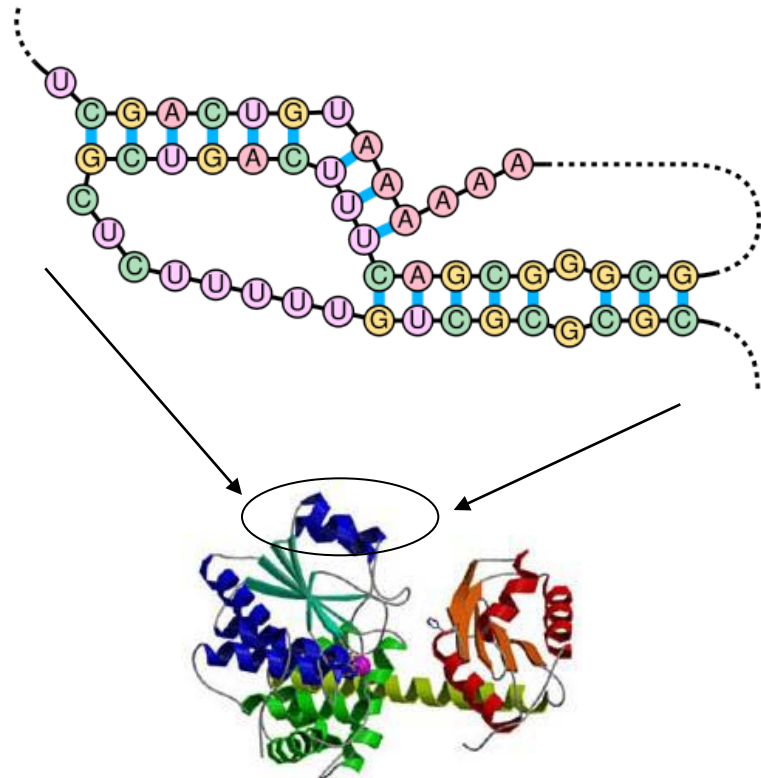
C U U C A U C A G A A A U G A C

(easy to obtain by next-generation sequencer)

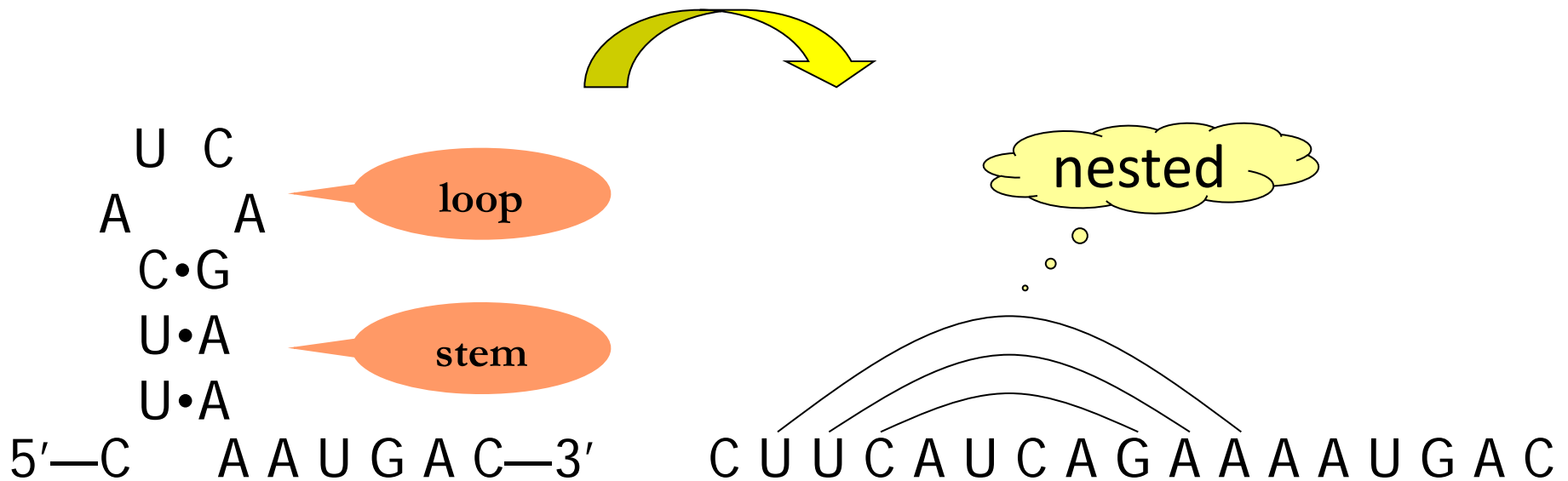
- **Secondary structure:**

folding structure by
hydrogen bonding
between bases

- Tertiary (3D) structure



Stem-loop structure

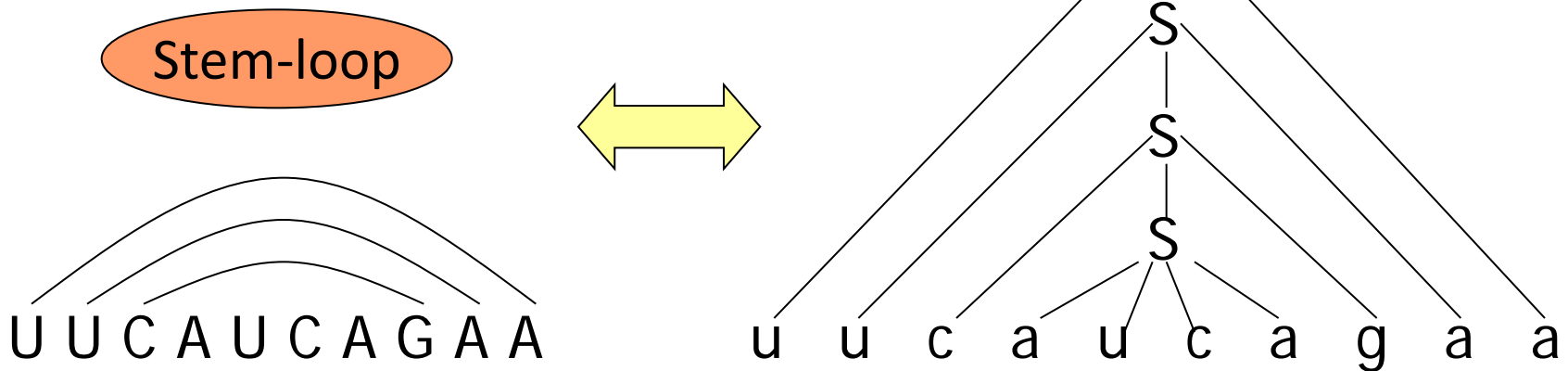


Secondary structure and grammar

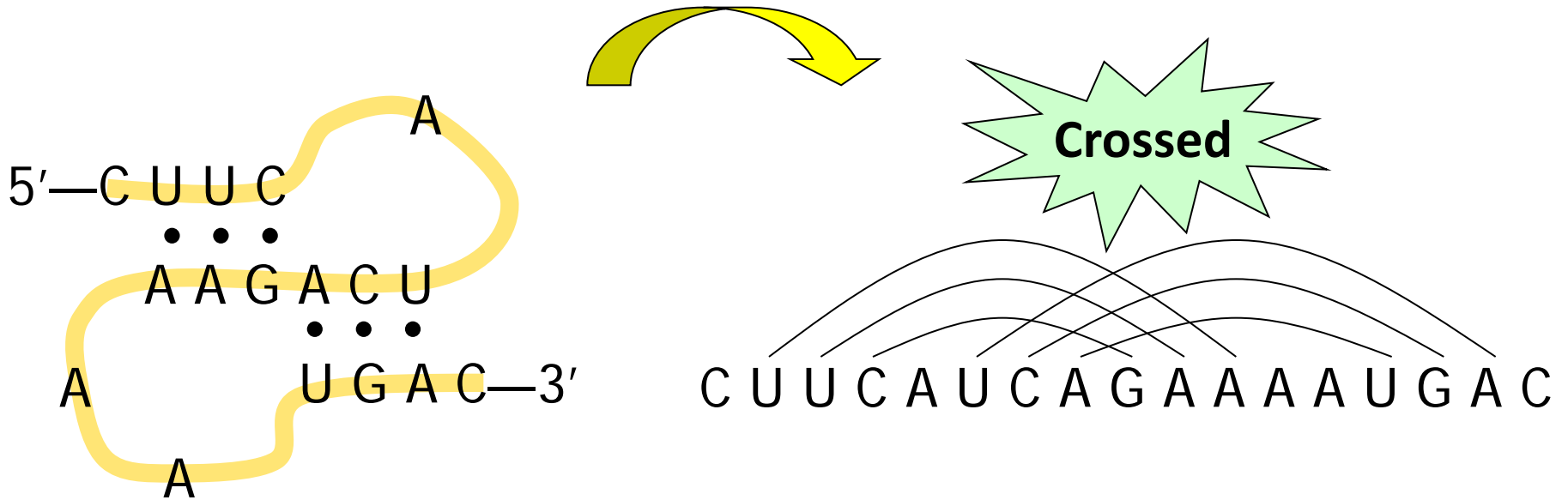
- Modeling biological regularity by a grammar G
- Secondary structure prediction
= **Parsing** primary sequence with G

$$S \rightarrow aSu \mid uSa \mid cSg \mid gSc$$

$$S \rightarrow auca$$



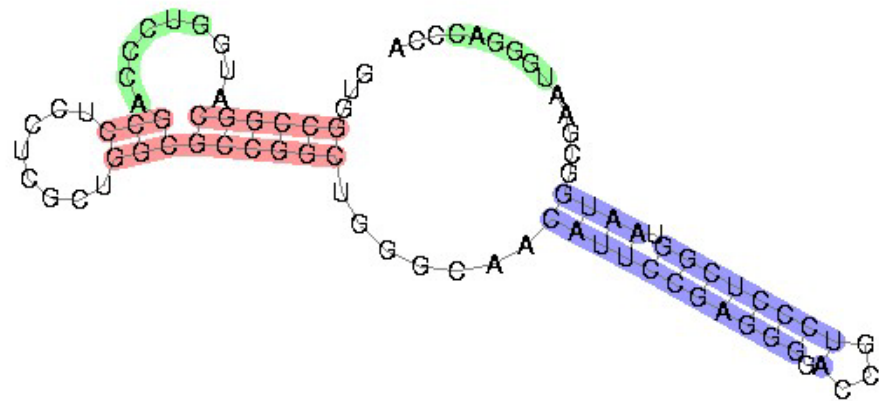
Pseudoknot



- CFG cannot represent **pseudoknot**.

HDV_ribozyme

- Hepatitis delta virus ribozyme
- length: 87–91



Existing methods using comparative approach for pseudoknots

- **hxmatch** [Witwer+04]
 - **ILM** [Ruan+04]
 - **Simulfold** [Meyer&Miklós07]
 - **Pair-SMCFG** [Mizoguchi+09]
 - **Proposed method** [Mizoguchi+11]
- } Not grammar based

[Witwer+04] Witwer , Hofacker & Stadler, IEEE Trans. Computational Biology and Bioinformatics, 2(2), 2004

[Meyer&Miklós07] Meyer & Miklós, PLoS Computational Biology, 3(8), 2007

[Ruan+04] Ruan, Stormo & Zhang, Bioinformatics 20(1), 2004.

[Mizoguchi+09] Mizoguchi, Kato & Seki, 20th International Conference on Genome Informatics poster, 2009

[Mizoguchi+11] Mizoguchi, Kato & Seki, A grammar-based approach to RNA pseudoknotted structure prediction for aligned sequences, *IEEE 2011 ICCABS*.

Prediction accuracy

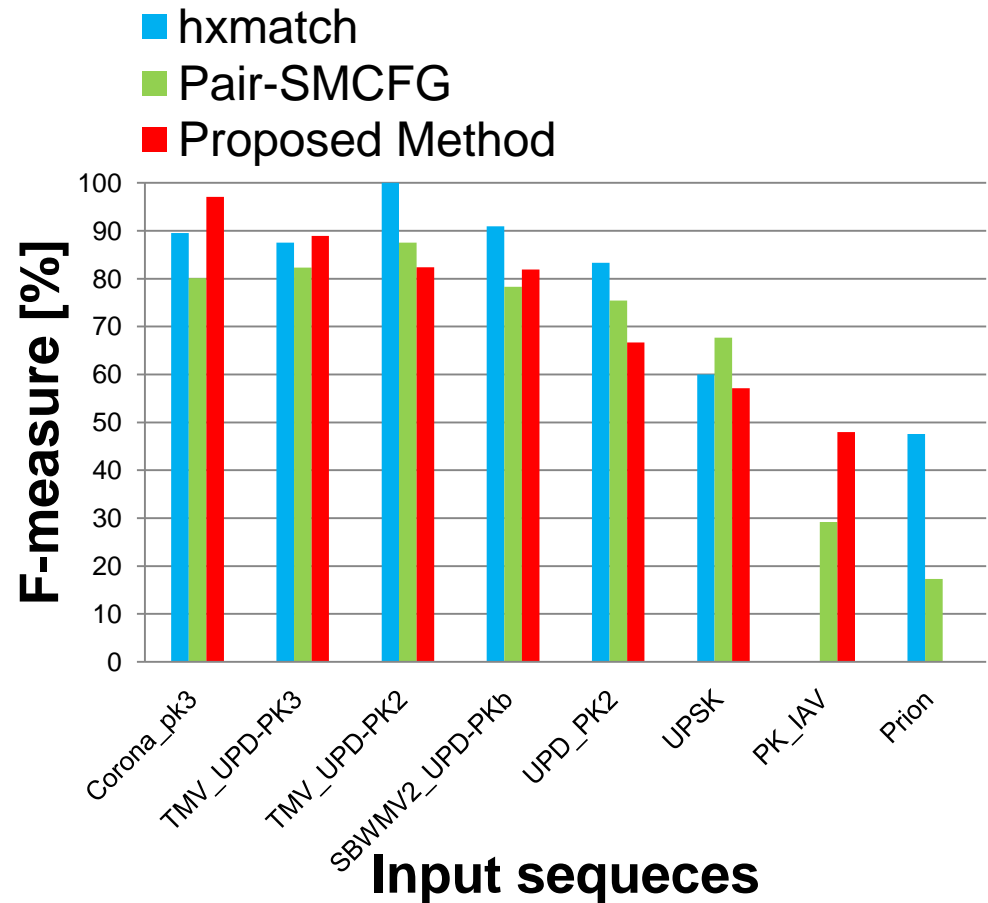
Data set:

8 families from Rfam

database [Griffiths-Jones+05].

Average F-measure:

- ◆ **hxmatch**: 69.85 %
- ◆ **Pair-SMCFG**: 64.73 %
- ◆ **Proposed Method**: 65.26 %



Thank you!