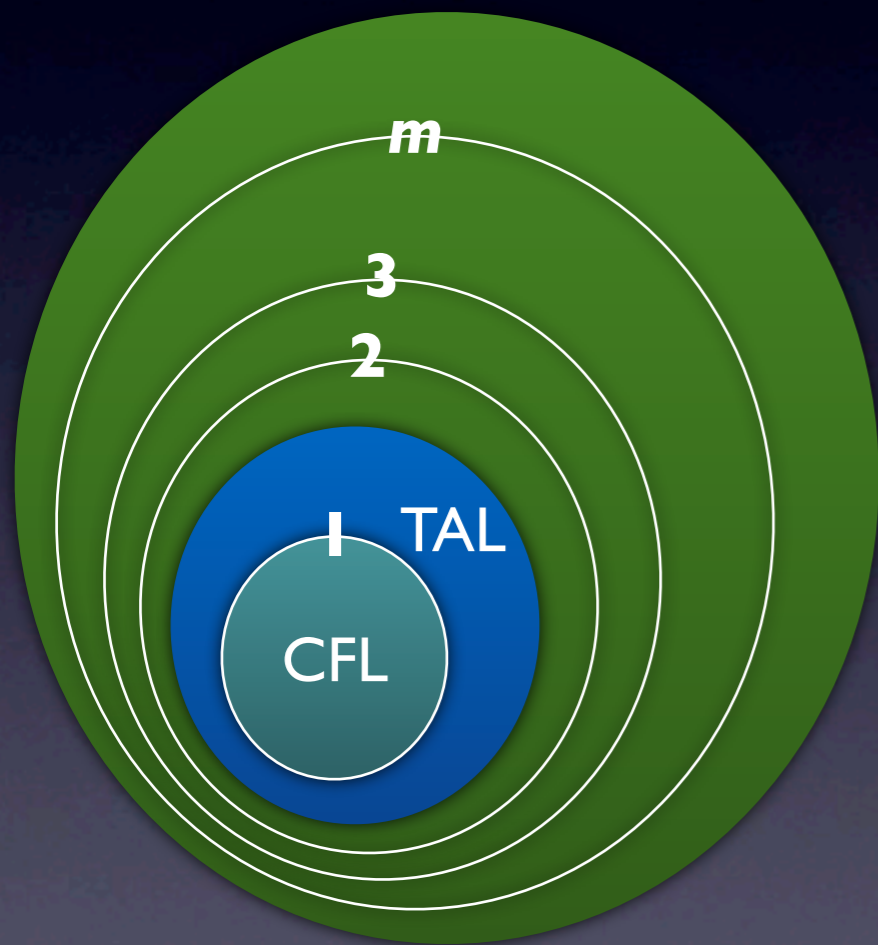


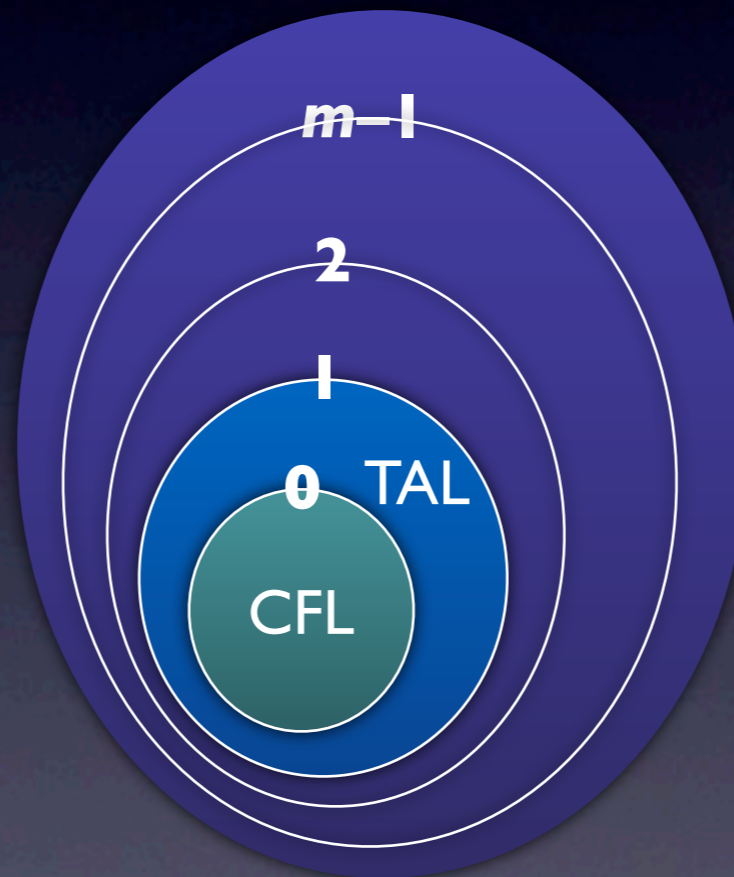
Two Subhierarchies Inside the MCFLs

Makoto Kanazawa
National Institute of Informatics

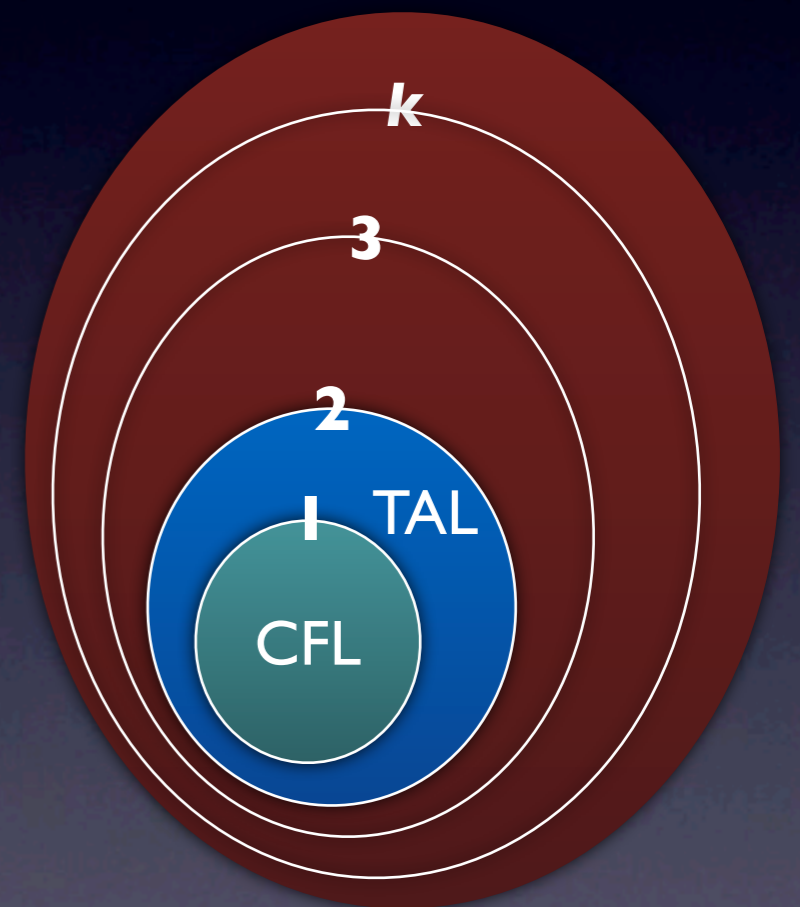
Three Infinite Hierarchies



$$\text{MCFL} = \bigcup_{m \geq 1} m\text{-MCFL}$$



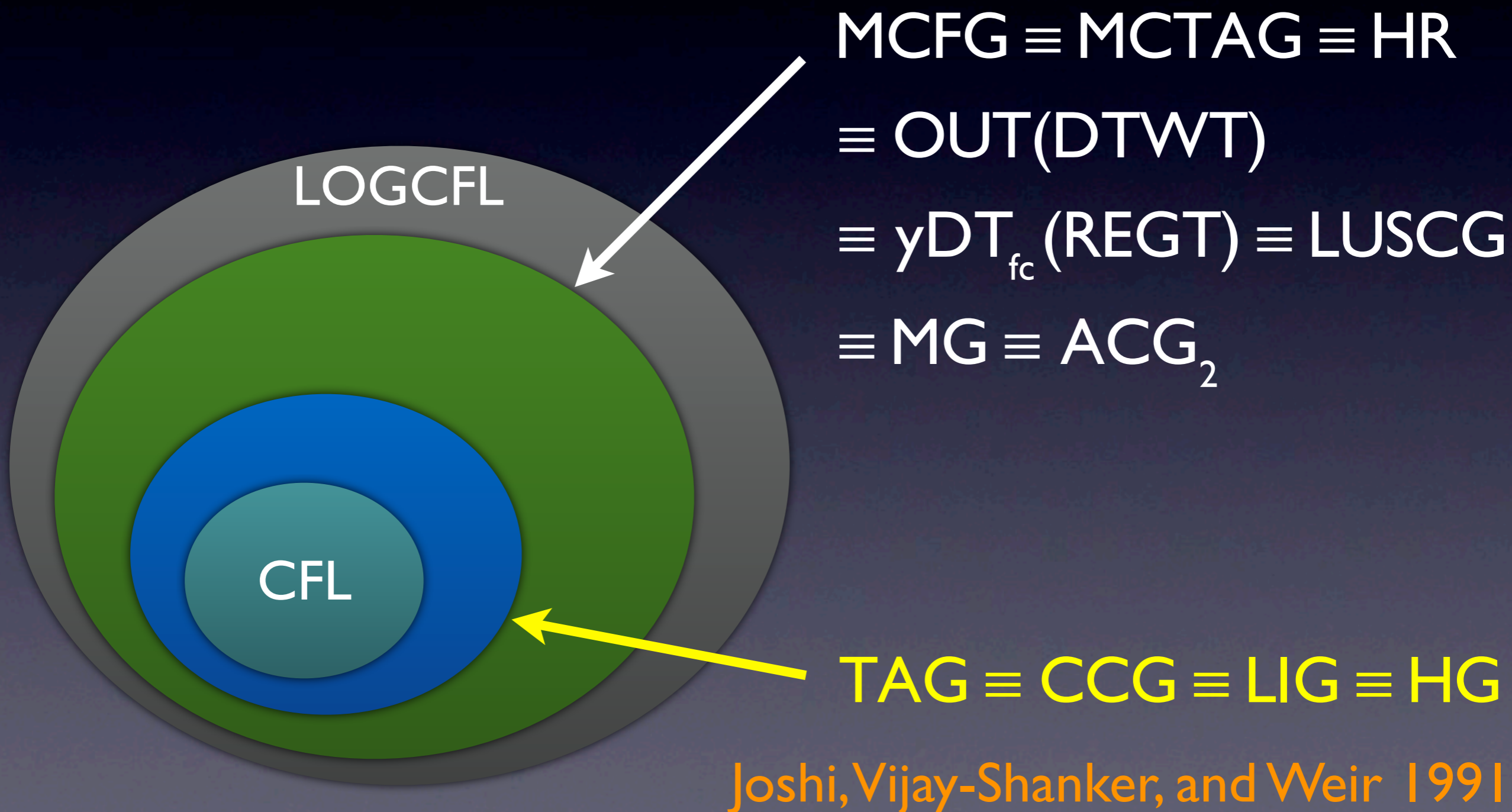
$$y\text{CFT}_{\text{sp}} = \bigcup_{m \geq 1} y\text{CFT}_{\text{sp}}(m-1)$$



$$\mathbf{C} = \bigcup_{k \geq 1} \mathbf{C}_k$$

$y\text{CFT}_{\text{sp}}(1)$ and \mathbf{C}_2 coincide with TAL, the class of languages generated by Tree Adjoining Grammars.

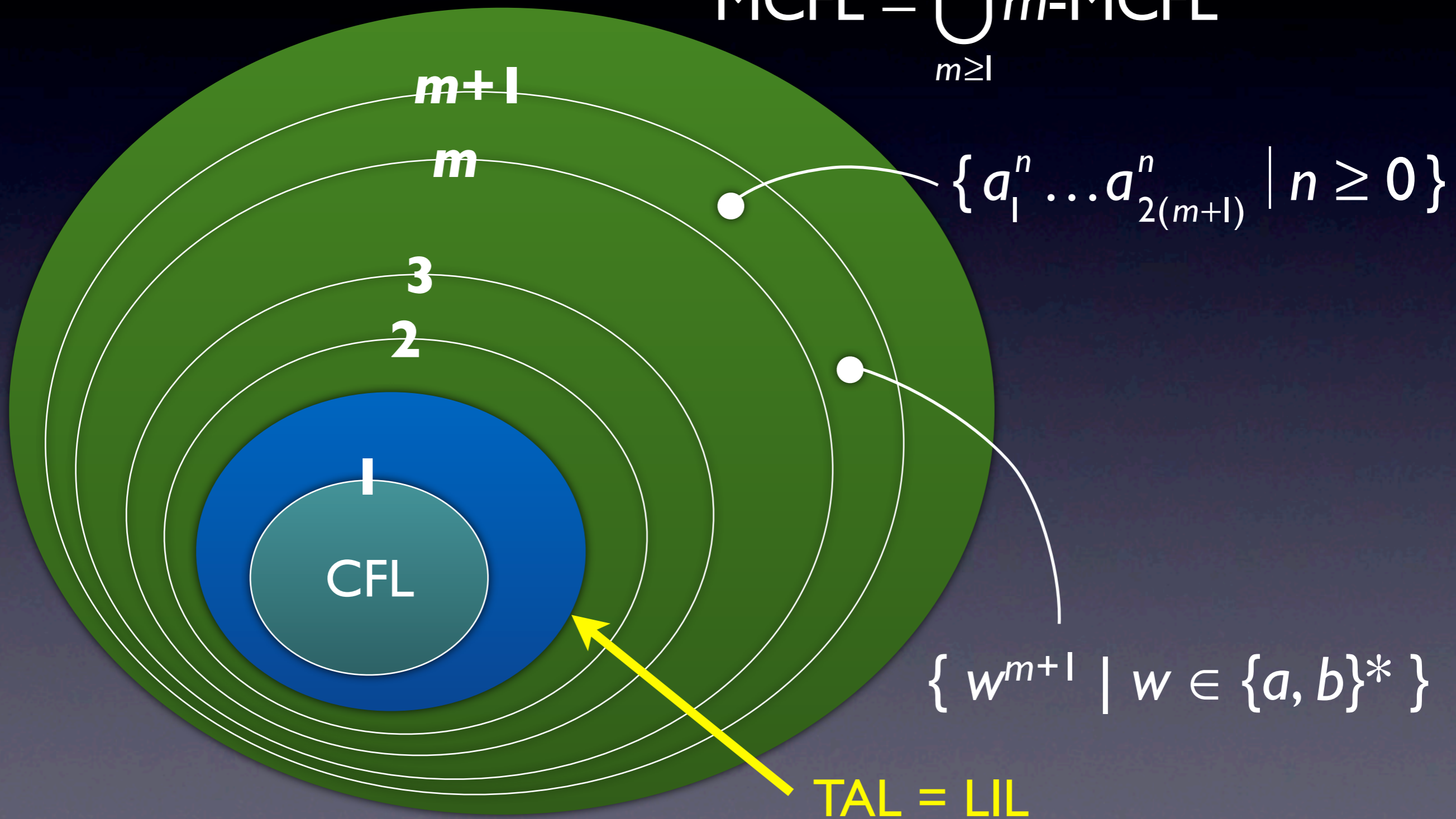
Convergence of Mildly Context-Sensitive Grammar Formalisms



The “convergence of mildly context-sensitive ...” originally referred to TAG, CCG, LIG, HG. CCG is rather ad hoc. TAG and HG are not very different.

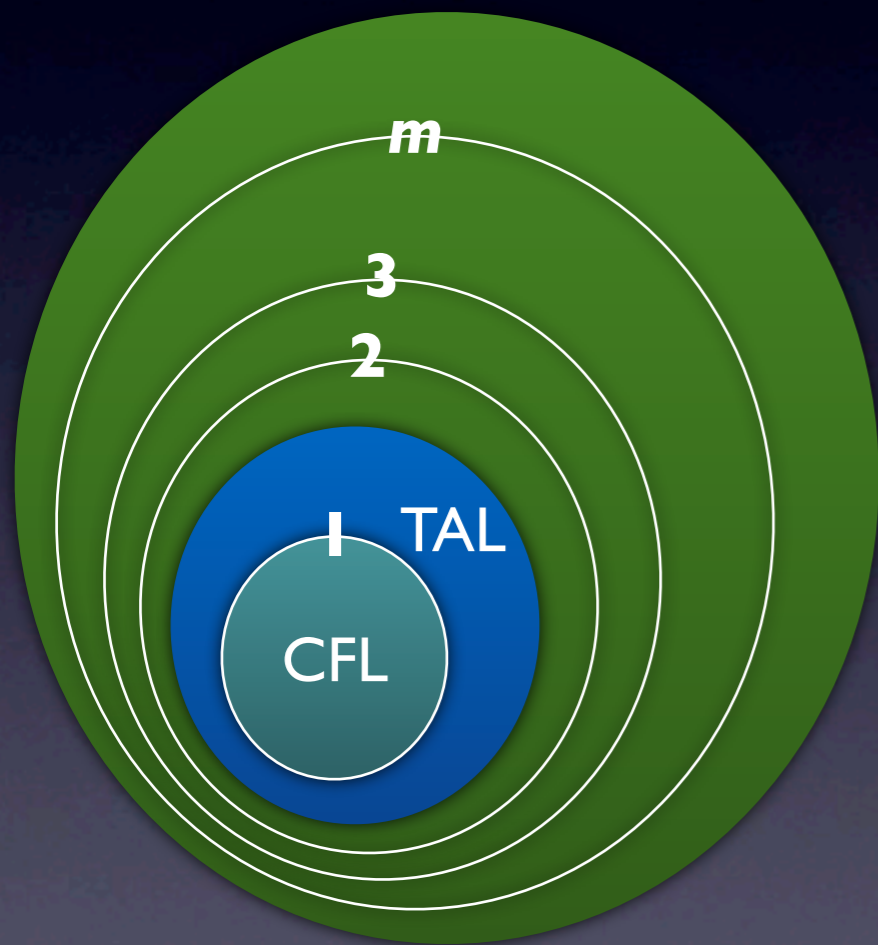
MCFL Hierarchy

$$\text{MCFL} = \bigcup_{m \geq 1} m\text{-MCFL}$$

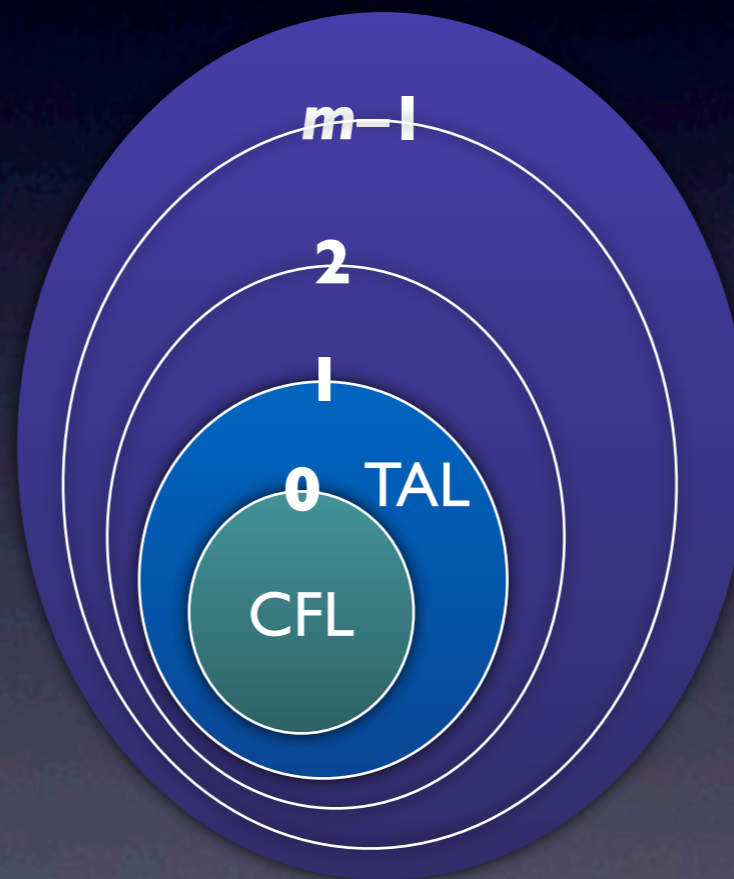


This is an infinite hierarchy.

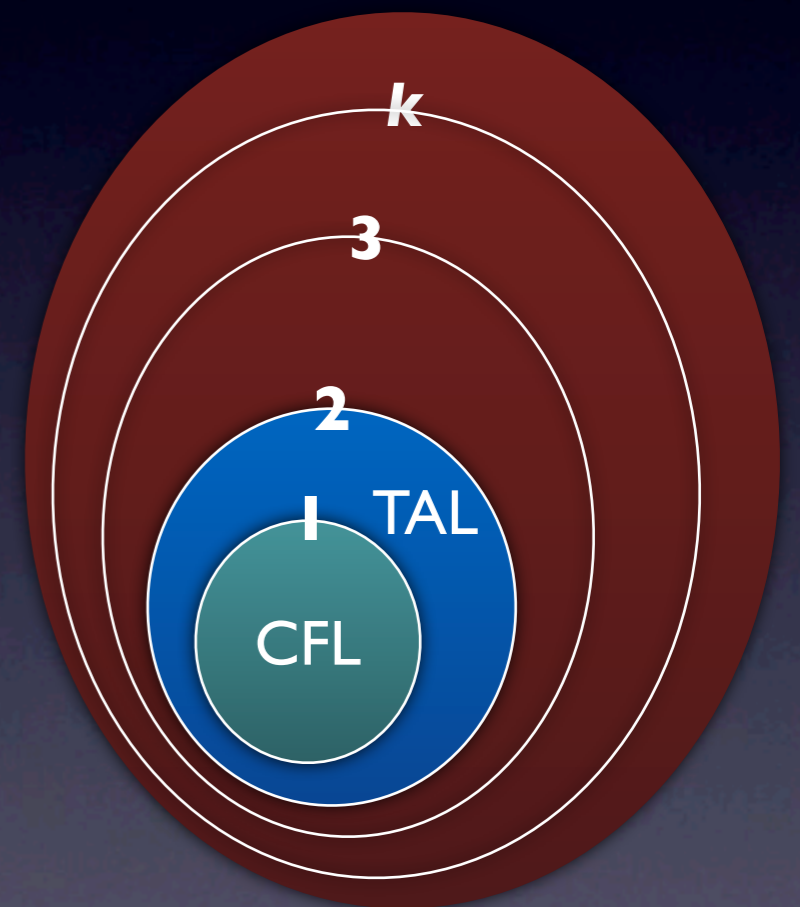
Three Infinite Hierarchies



$$\text{MCFL} = \bigcup_{m \geq 1} m\text{-MCFL}$$



$$y\text{CFT}_{\text{sp}} = \bigcup_{m \geq 1} y\text{CFT}_{\text{sp}}(m-1)$$



$$\mathbf{c} = \bigcup_{k \geq 1} \mathbf{c}_k$$

$$\text{TAG} \equiv \text{CFT}_{\text{sp}}(\text{I})$$

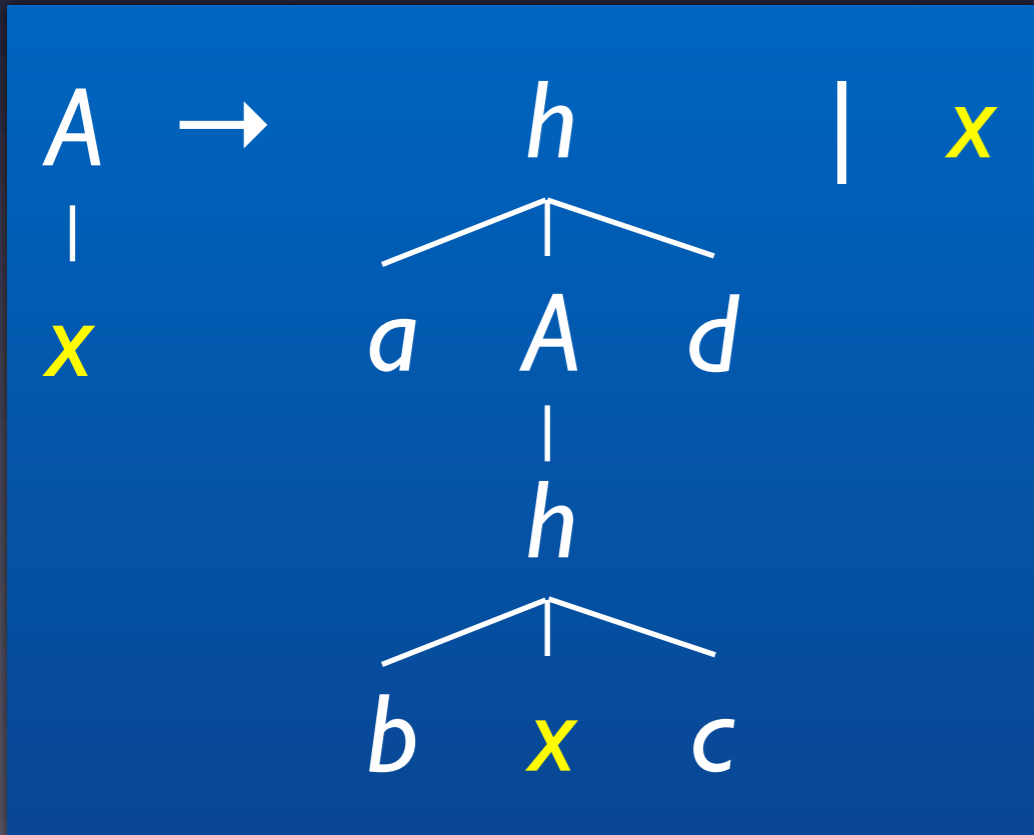
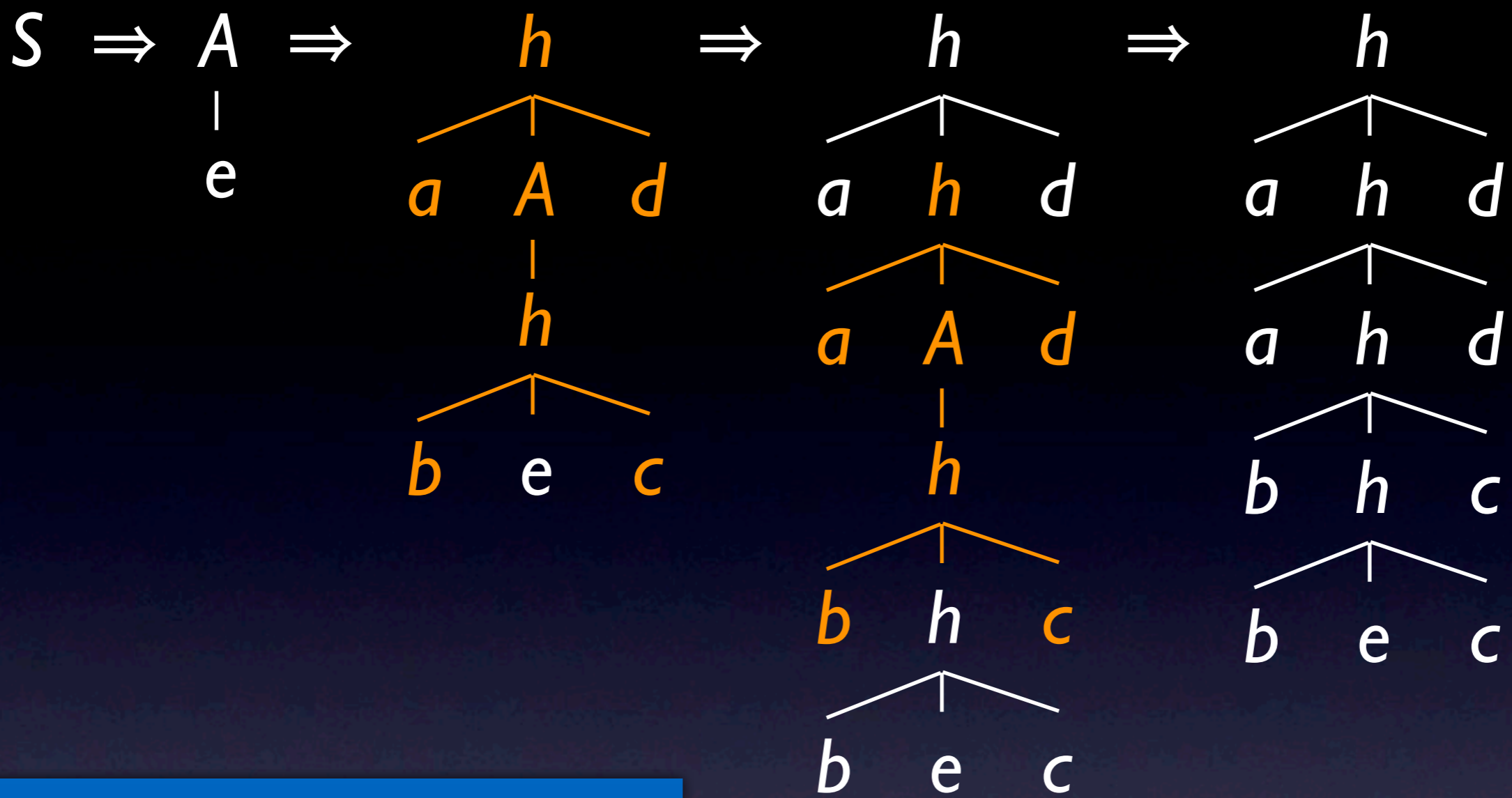


$$N = N^{(0)} \cup N^{(1)} = \{S\} \cup \{A\}$$

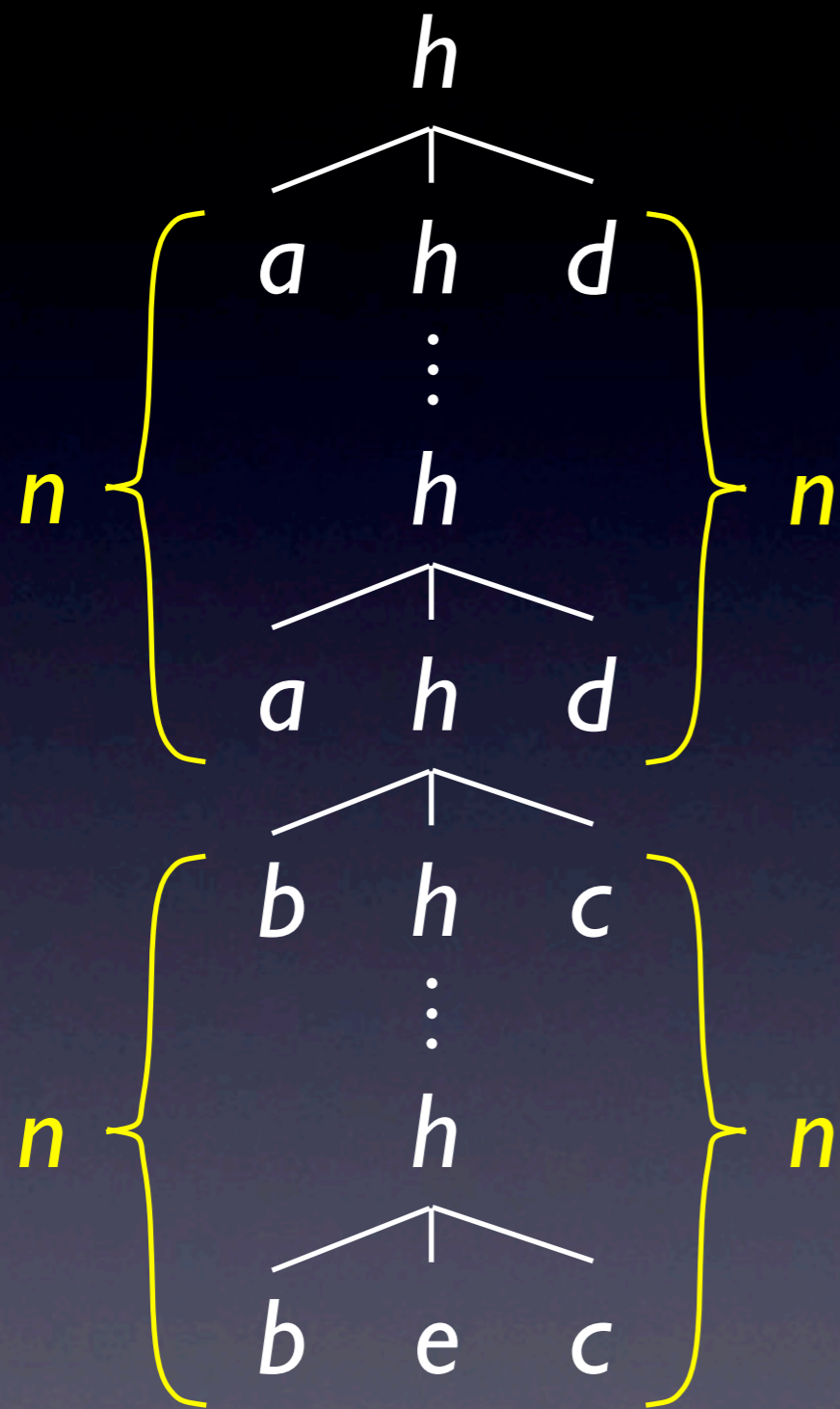
$$\Sigma = \Sigma^{(0)} \cup \Sigma^{(1)} = \{a, b, c, d, e\} \cup \{h\}$$

A CFT grammar generates a ranked tree language. Nonterminals and terminals come with ranks.

This grammar is “monadic” because the maximal rank of a nonterminal is 1.



An example of a derivation.

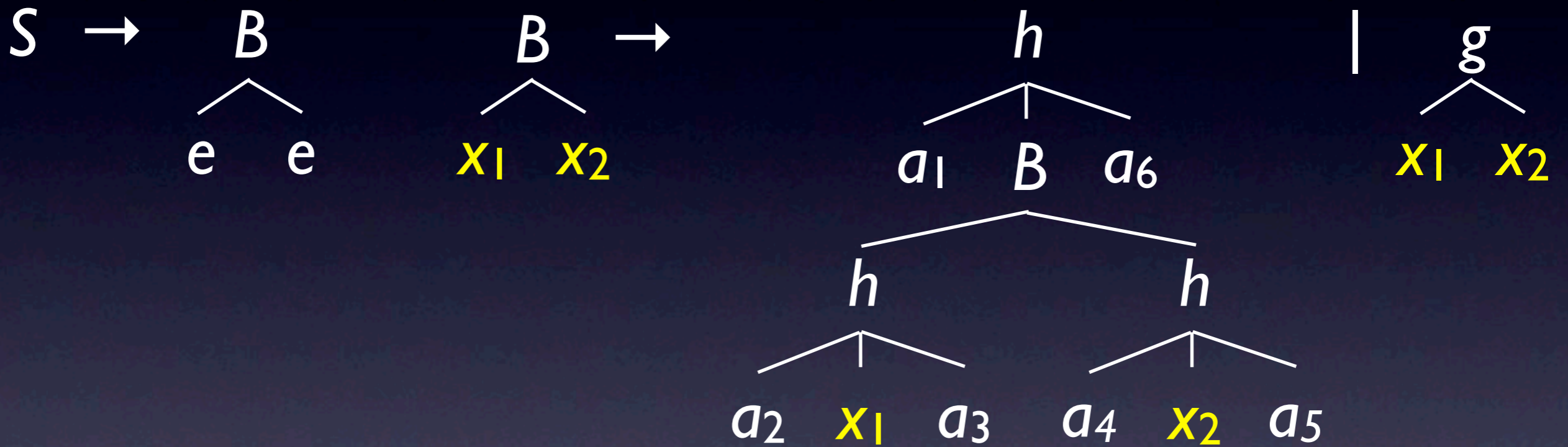


$a^n b^n e c^n d^n$

yield

This example grammar generates all trees of this form.

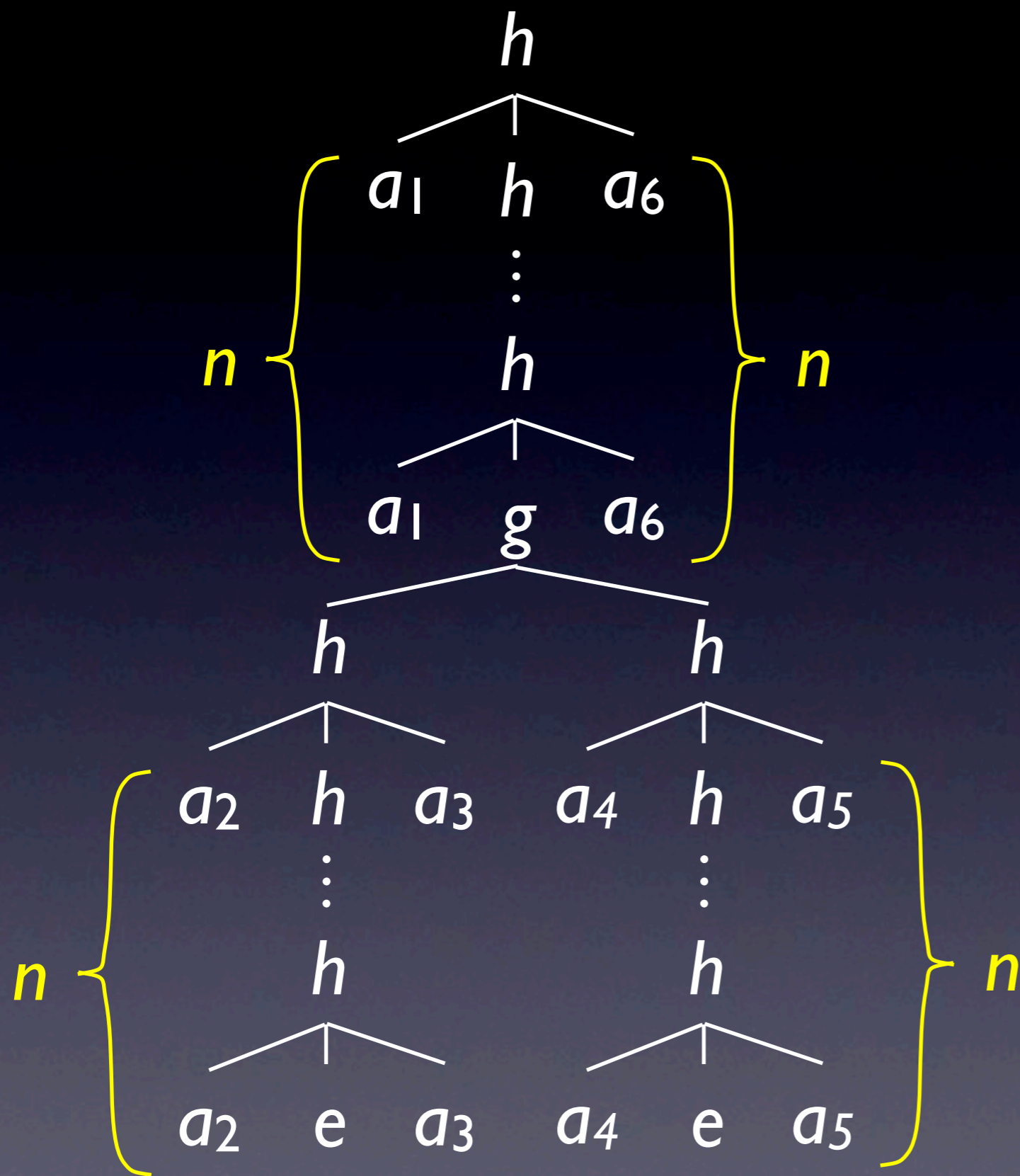
CFT_{sp}(2)



$$N = N^{(0)} \cup N^{(2)} = \{S\} \cup \{B\}$$

$$\Sigma = \Sigma^{(0)} \cup \Sigma^{(2)} \cup \Sigma^{(3)} = \{a_1, a_2, a_3, a_4, a_5, a_6, e\} \cup \{g\} \cup \{h\}$$

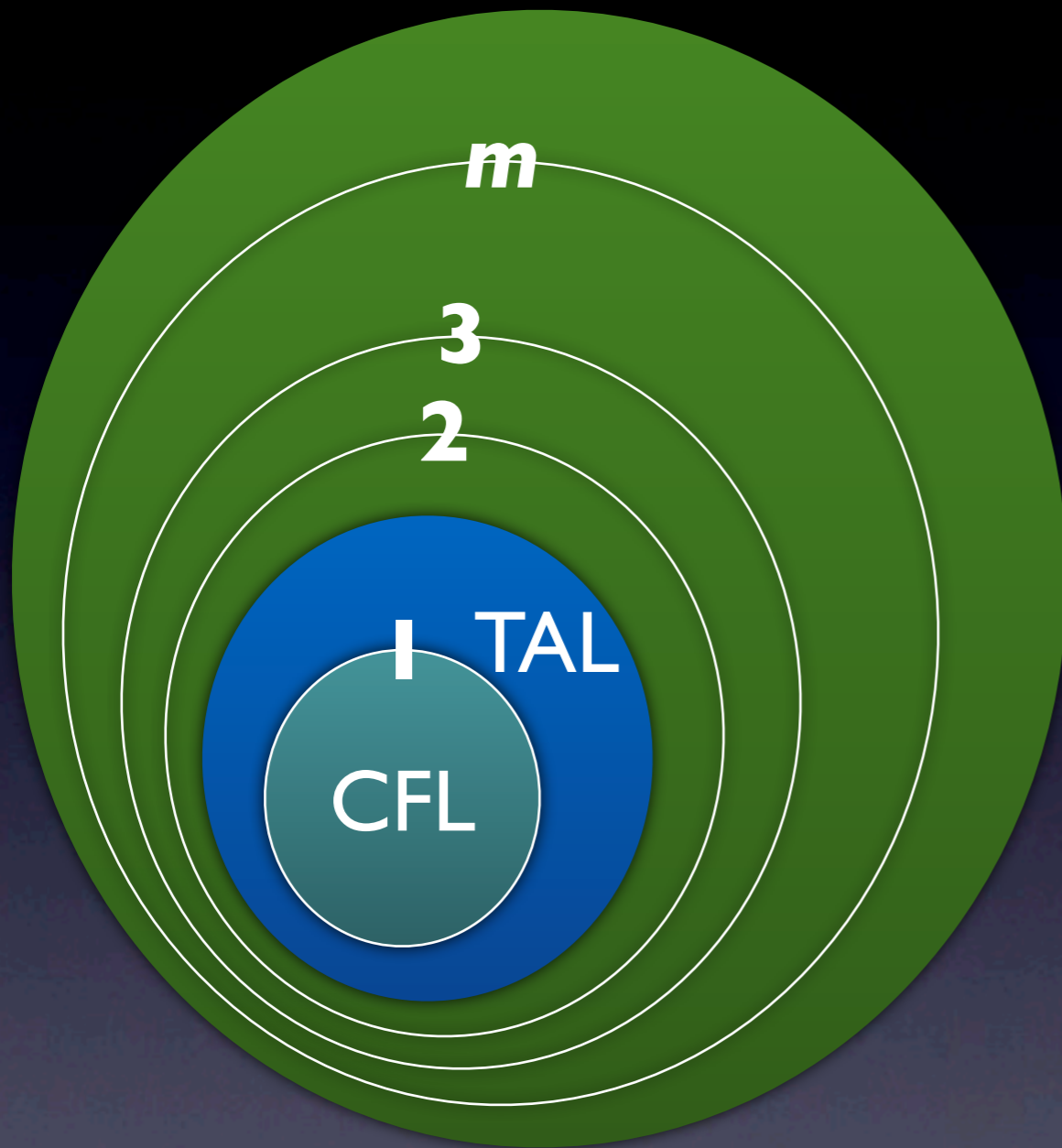
This grammar has a rank 2 nonterminal.



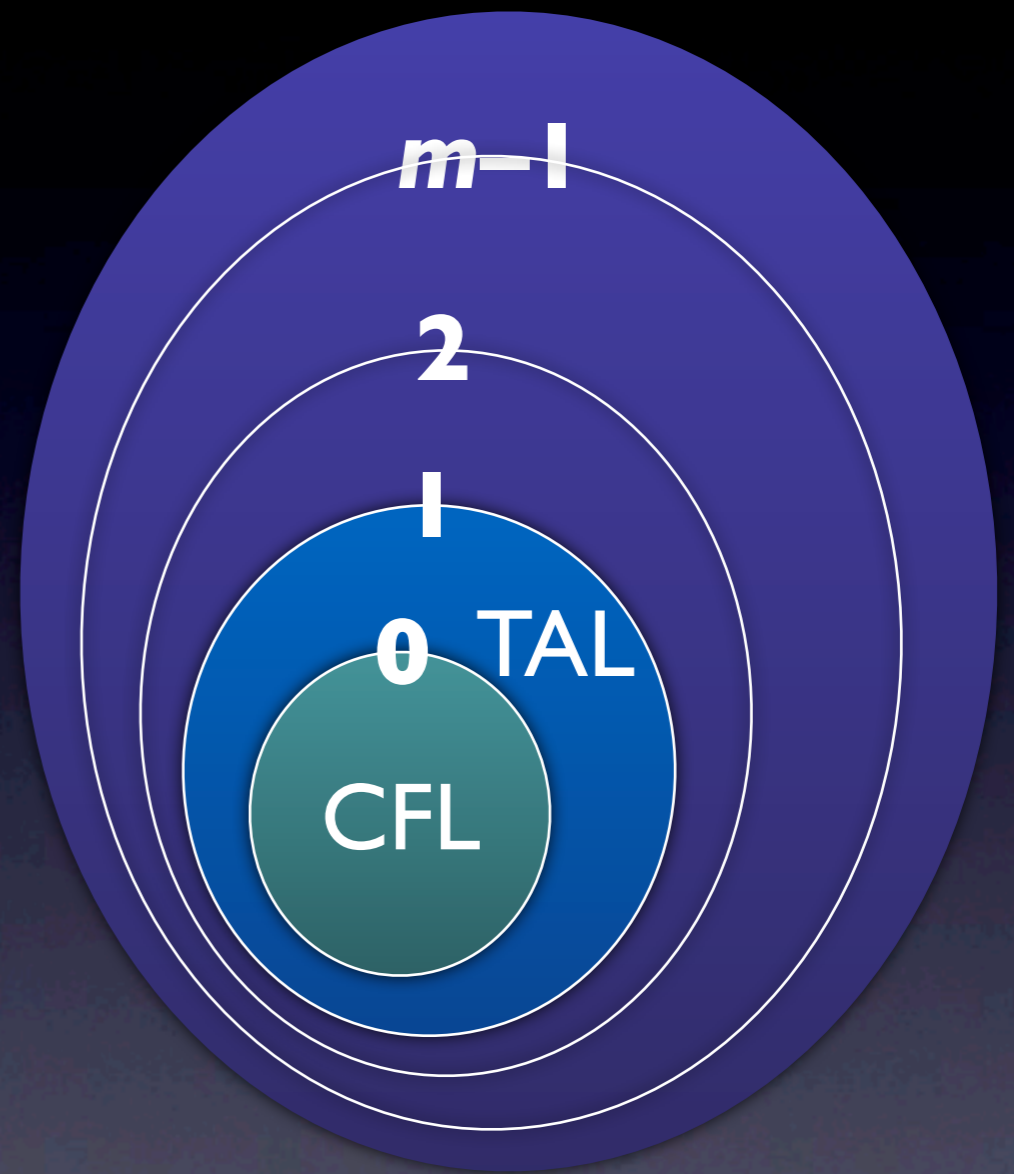
$a_1^n a_2^n e a_3^n a_4^n e a_5^n a_6^n$

yield

The yield image of the language of this grammar is not a TAL.



$$\text{MCFL} = \bigcup_{m \geq 1} m\text{-MCFL}$$



$$y\text{CFT}_{\text{sp}} = \bigcup_{m \geq 1} y\text{CFT}_{\text{sp}}(m-1)$$

Let's look at how corresponding levels of these two hierarchies compare with each other.

$$y\text{CFT}_{\text{sp}}(m-1) \subseteq m\text{-MCFL}$$

Seki and Kato 2008

de Groote and Pogodalla 2004, Salvati 2007

“Context-Free” Grammar Formalisms

$$S \rightarrow \begin{array}{|c|c|} \hline B & S \\ \hline B & \\ \hline \end{array}$$

top-down rewriting of
sentential forms

$$S \left(\begin{array}{|c|c|} \hline X & Z \\ \hline Y & \\ \hline \end{array} \right) \leftarrow B(X), B(Y), S(Z)$$

bottom-up construction of
derived objects

In order to convert a CFT_{sp} to an MCFG, we take a bottom-up view of the former, because the latter is a bottom-up formalism.

Top-down vs. bottom-up

Type 0

EFS

Smullyan 1961

Type 1 = CSG

I.b. EFS \equiv CSG

Arikawa et al. 1989

simple LMG \equiv P

Groenink 1997

PMCFG

Seki et al. 1991, Groenink 1997

MCFG

Seki et al. 1991, Groenink 1997

Type 2 = CFG

simple EFS \equiv CFG

Arikawa 1970

Type 3

rewriting systems

logic programs on strings

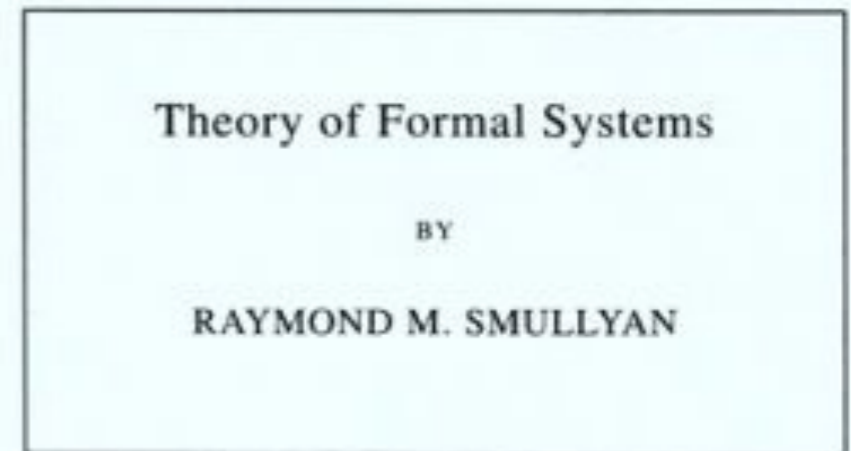
The bottom-up formalisms give you a richer picture. Note that almost all formalisms that were not defined by Chomsky do not fit within the Chomsky Hierarchy. For example, an indexed grammar is *not* an instance of a Type 0 grammar. An MCFG is an instance of an Elementary Formal System of Smullyan, which was rediscovered by Groenink.

1956

Thue → Post → Chomsky →

Smullyan

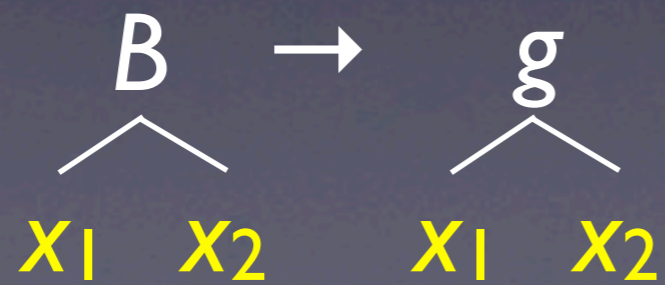
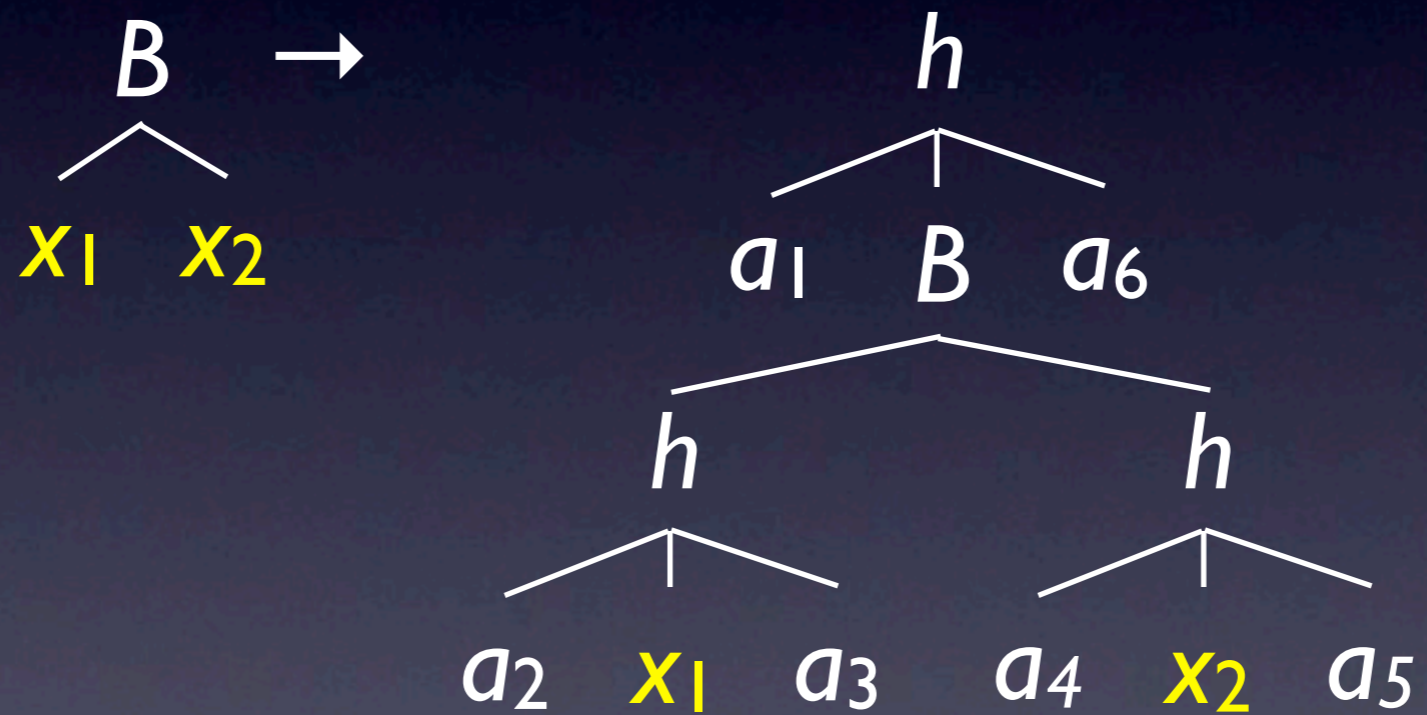
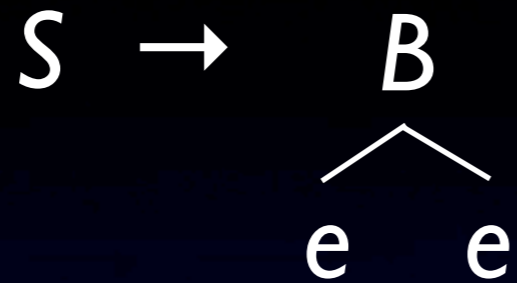
1961



ANNALS OF MATHEMATICS STUDIES
PRINCETON UNIVERSITY PRESS

It's too bad that Smullyan's work came a little too late for Chomsky to take notice. The world would have been a better place if Chomsky had based his theory of formal grammar on Smullyan's work, rather than the work of Thue and Post.

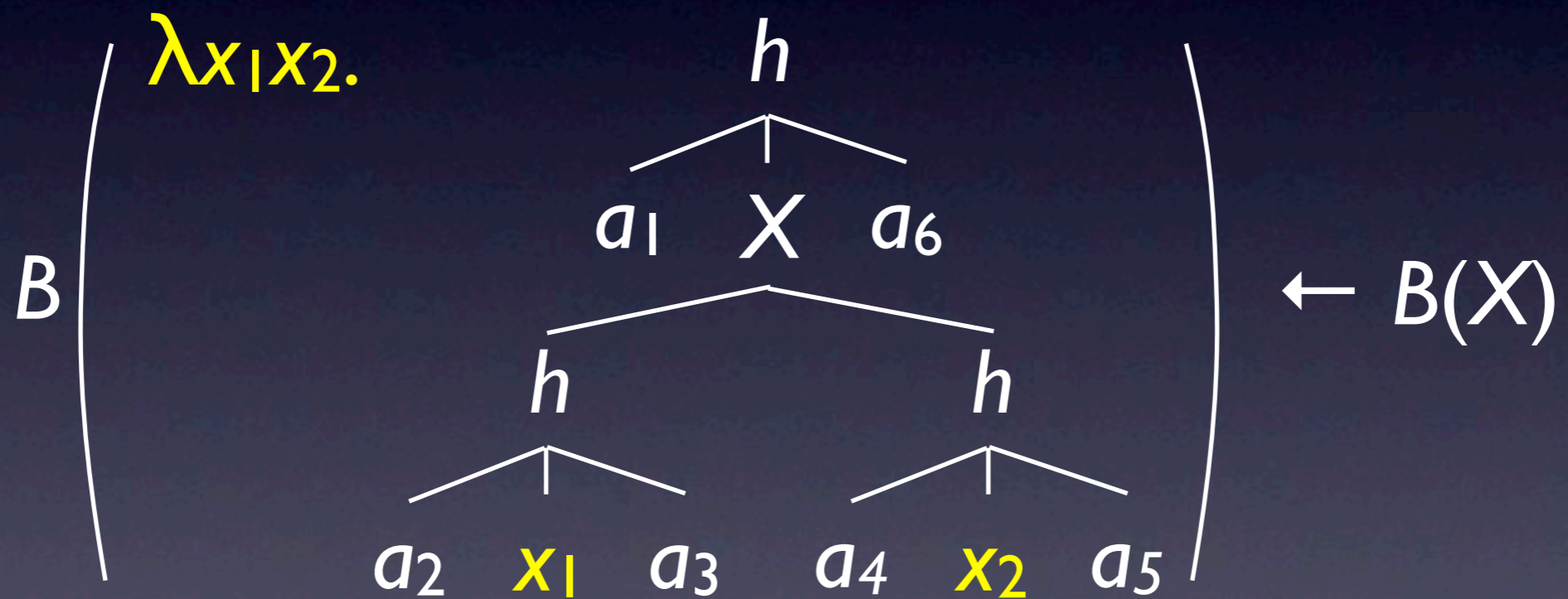
Top-down View



Standard, top-down view of CFT grammar rules.

Bottom-up View

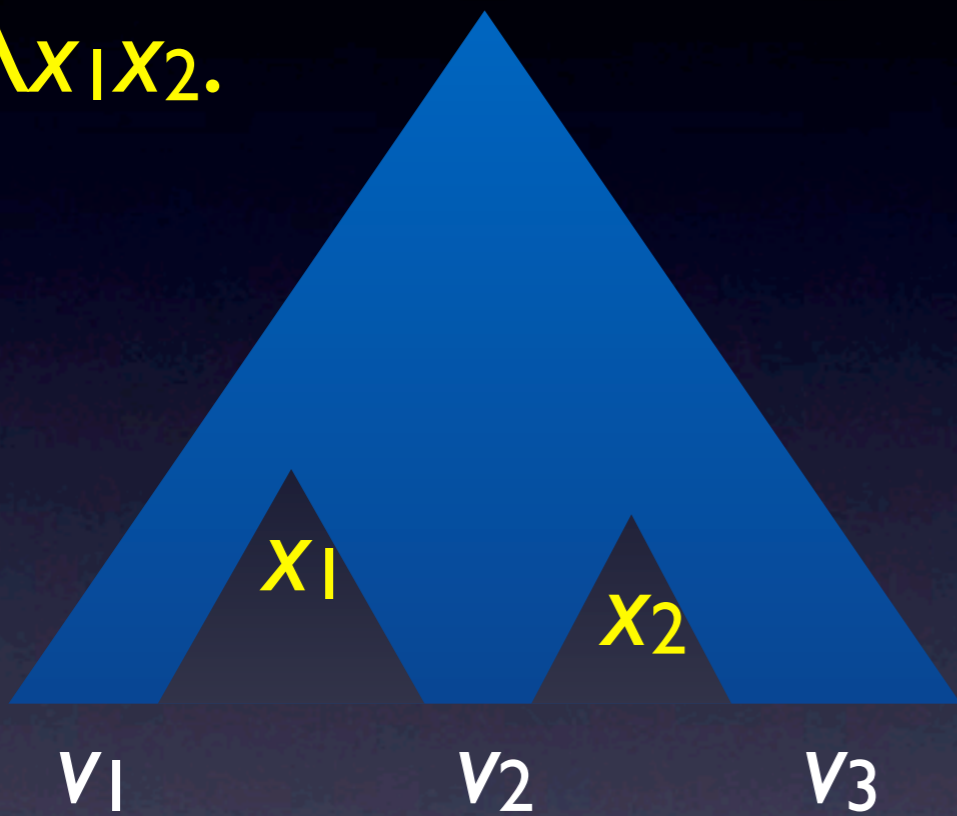
$$S \left(\begin{array}{c} X \\ \wedge \\ e \quad e \end{array} \right) \leftarrow B(X)$$



$$B \left(\begin{array}{c} \lambda x_1 x_2. \\ \wedge \\ g \\ \wedge \\ x_1 \quad x_2 \end{array} \right) \leftarrow$$

An alternative bottom-up view.

$\lambda x_1 x_2.$

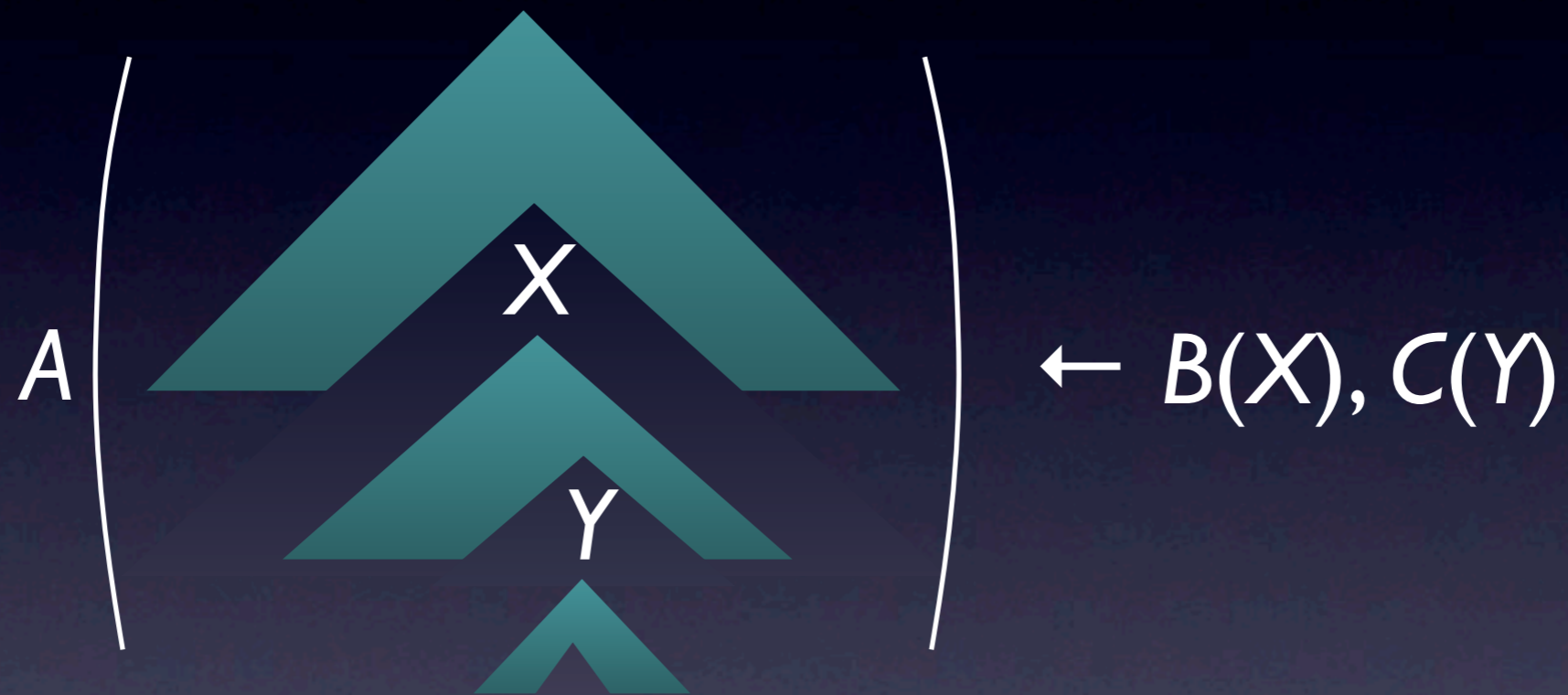


(v_1, v_2, v_3)

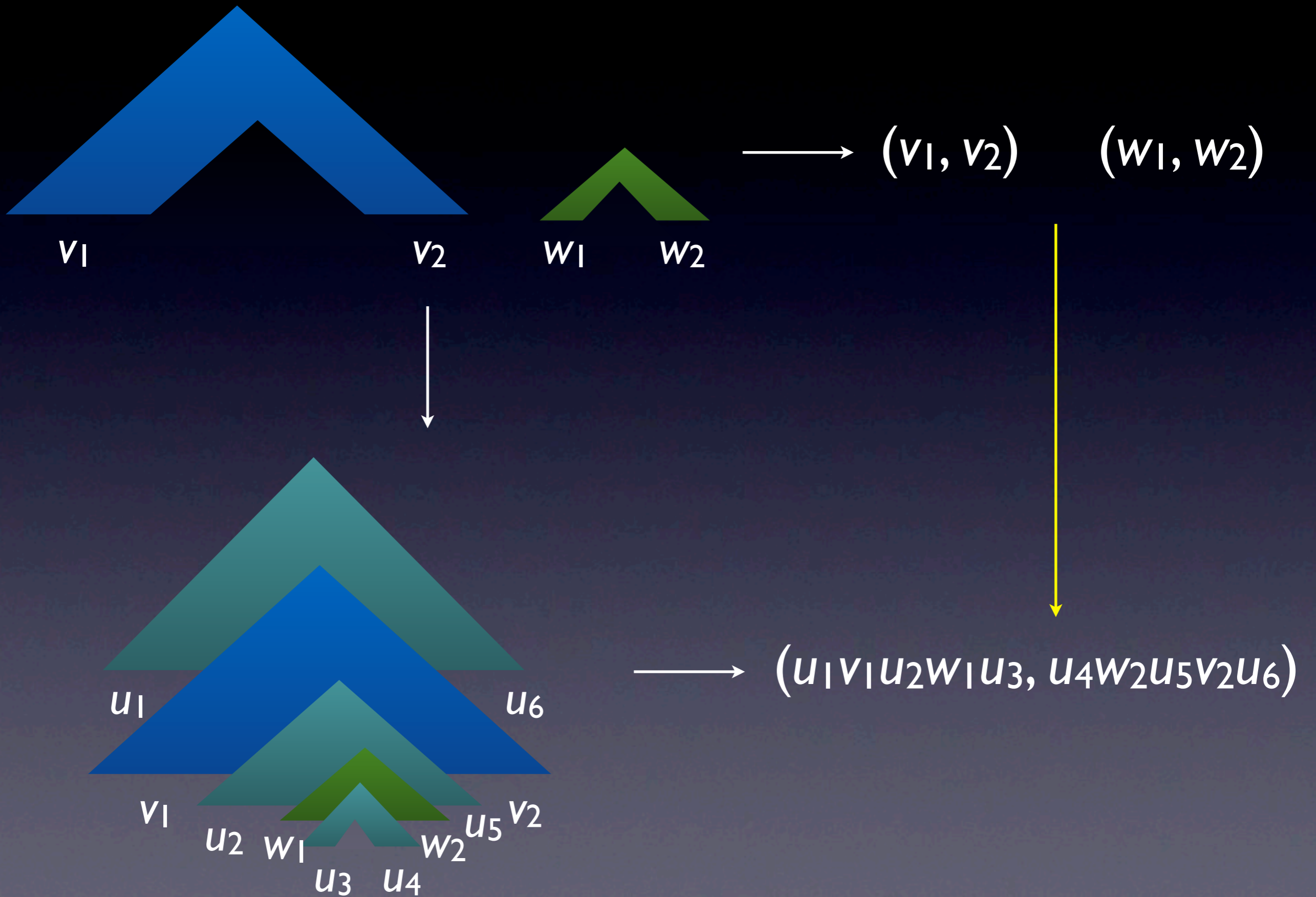
n-ary tree context

$(n+1)$ -tuple of strings

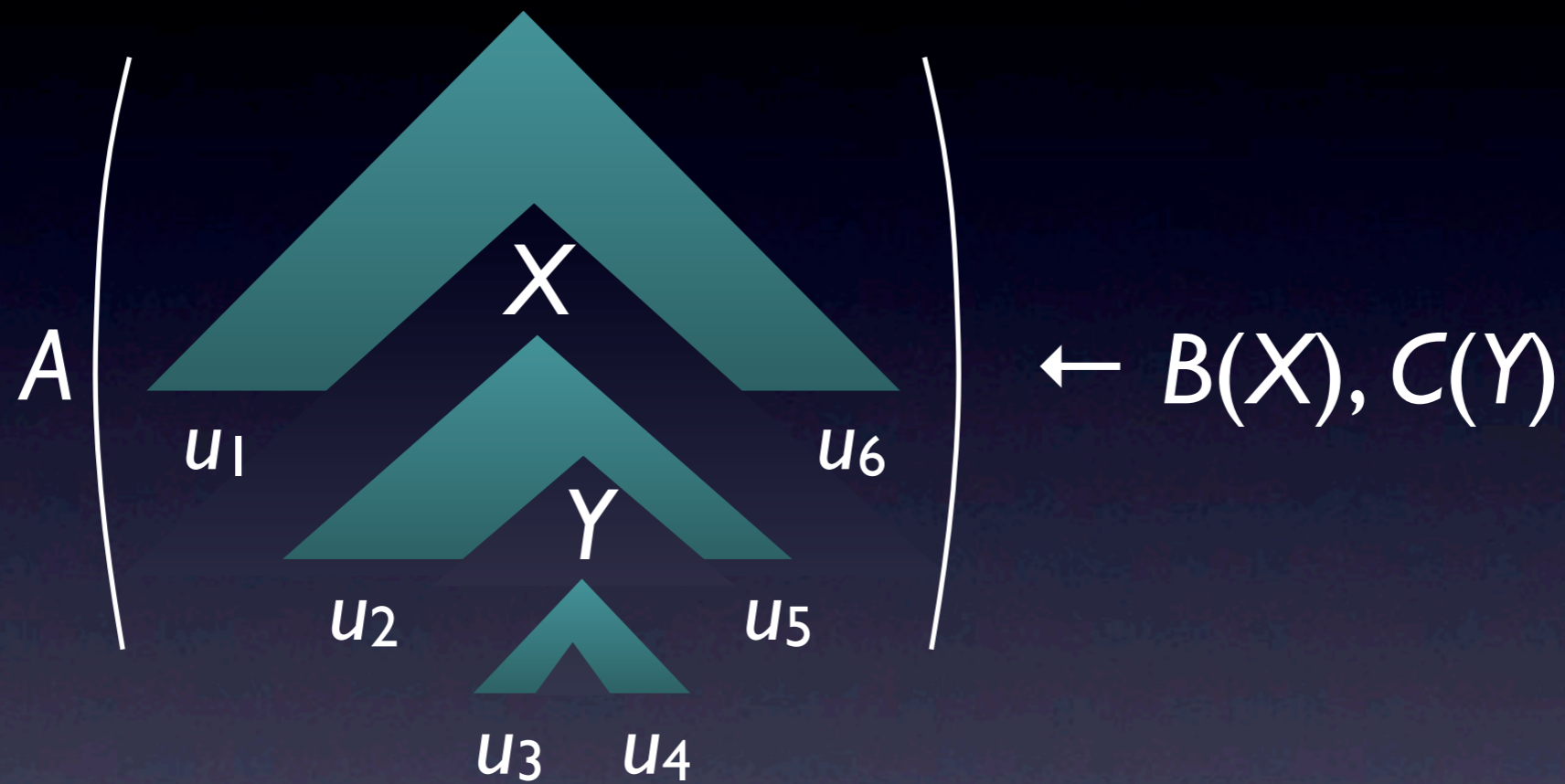
There's a natural mapping from n-ary tree contexts to $(n+1)$ -tuples of strings.



What about the operation in a CFT_{sp} rule? Does it naturally correspond to an operation on tuples of strings?



Yes, it does. This is a commutative diagram!



$$A(u_1 \mathbf{x}_1 u_2 \mathbf{y}_1 u_3, u_4 \mathbf{y}_2 u_5 \mathbf{x}_2 u_6) \leftarrow B(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2)$$

This is how to translate a $\text{CFT}_{\text{sp}}(1)$ rule to a 2-MCFG rule.

$$y\text{CFT}_{\text{sp}}(m-1) \subseteq m\text{-MCFL}$$


Seki and Kato 2008

de Groote and Pogodalla 2004, Salvati 2007


The translation establishes this inclusion.

Well-nested MCFGs

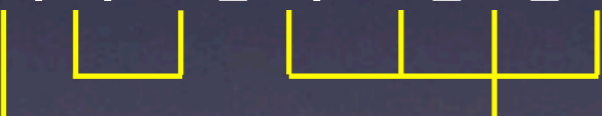
✗ $S(x_1 y_1 x_2 y_2) \leftarrow A(x_1, x_2), B(y_1, y_2)$




✓ $S(x_1 y_1 y_2 x_2) \leftarrow A(x_1, x_2), B(y_1, y_2)$



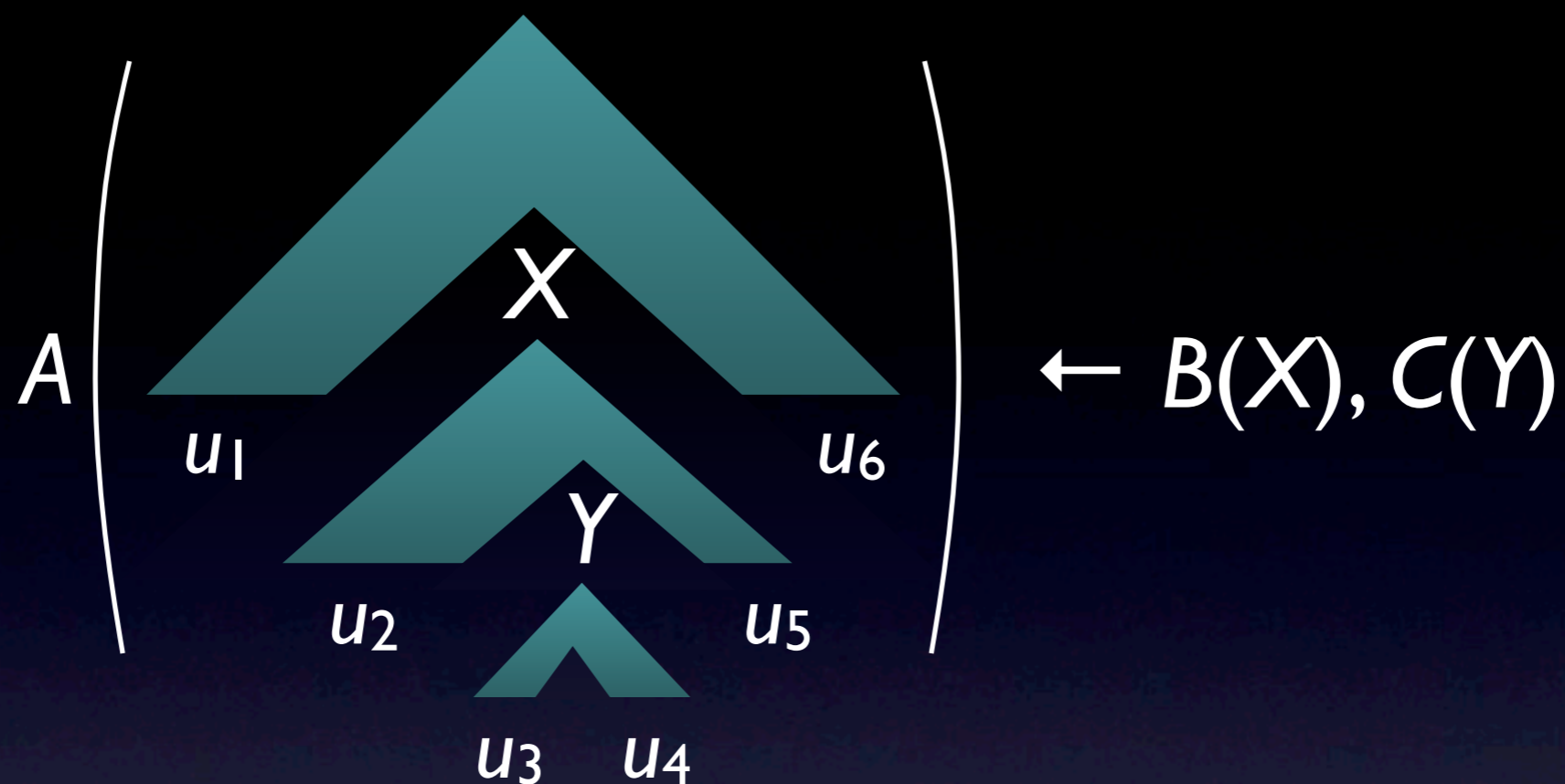
✗ $C(x_1 y_1 y_2 z_1 z_2 x_2 z_3) \leftarrow A(x_1, x_2), B(y_1, y_2), C(z_1, z_2, z_3)$



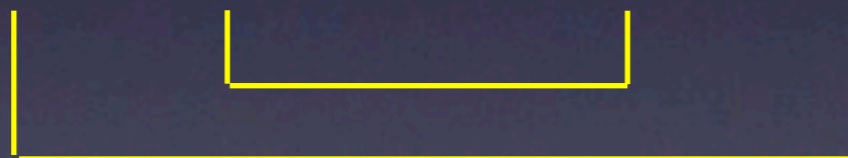
✓ $C(z_1 x_1 x_2 z_2 y_1 y_2 z_3) \leftarrow A(x_1, x_2), B(y_1, y_2), C(z_1, z_2, z_3)$



Cf. Kuhlmann 2007



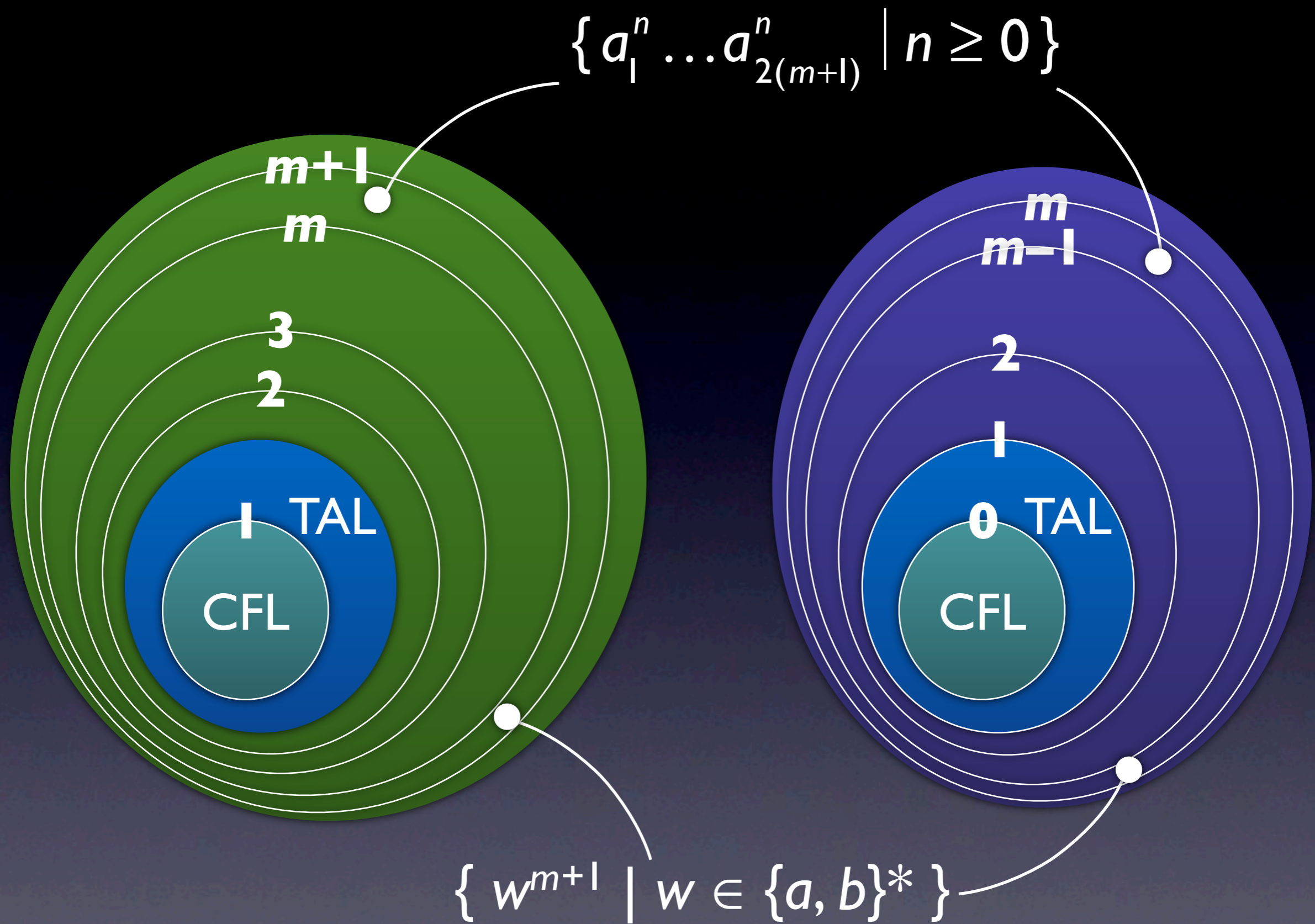
$$A(u_1 \mathbf{x}_1 u_2 \mathbf{y}_1 u_3, u_4 \mathbf{y}_2 u_5 \mathbf{x}_2 u_6) \leftarrow B(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2)$$



$$y\text{CFT}_{\text{sp}}(m-1) \subseteq m\text{-MCFL}_{\text{wn}}$$

$$y\text{CFT}_{\text{sp}}(m-1) \supseteq m\text{-MCFL}_{\text{wn}}$$

The translation of a CFT_{sp} rule gives you a well-nested MCFG rule.



$$\text{MCFL} = \bigcup_{m \geq 1} m\text{-MCFL}$$

$$y\text{CFT}_{\text{sp}} = \bigcup_{m \geq 1} y\text{CFT}_{\text{sp}}(m-1)$$

What about separation of the corresponding levels of the two hierarchies? These languages do not separate them.

m -MCFL vs. $yCFT_{sp}(m-1)$

$$RESP_2 = \{ a_1^i a_2^i b_1^j b_2^j a_3^i a_4^i b_3^j b_4^j \mid i, j \geq 0 \} \quad \text{Weir 1989}$$

$$RESP_2 \in 2\text{-MCFL} - yCFT_{sp}(1) \quad \text{Seki et al. 1991}$$

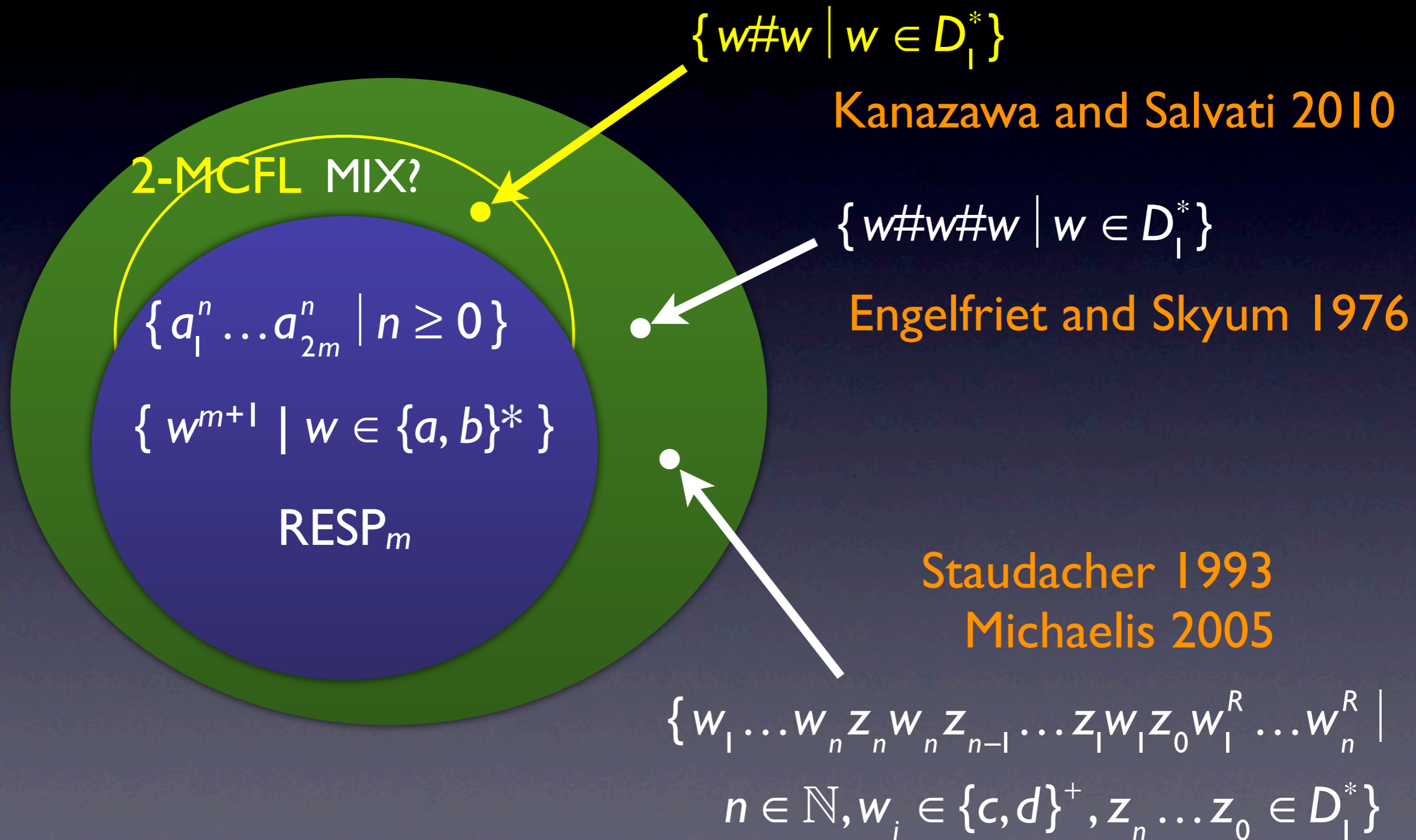
$$RESP_m = \{ a_1^i a_2^i b_1^j b_2^j \dots a_{2m-1}^i a_{2m}^i b_{2m-1}^j b_{2m}^j \mid i, j \geq 0 \}$$

$$RESP_m \in m\text{-MCFL} - yCFT_{sp}(m-1) \quad \text{for } m \geq 2$$

Seki and Kato 2008

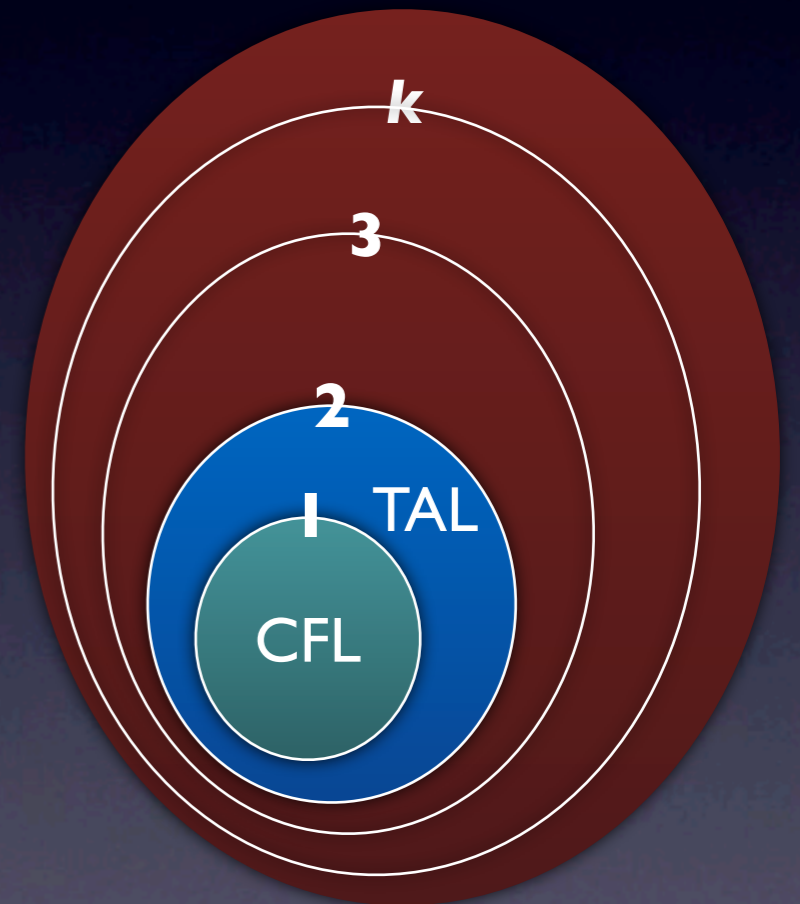
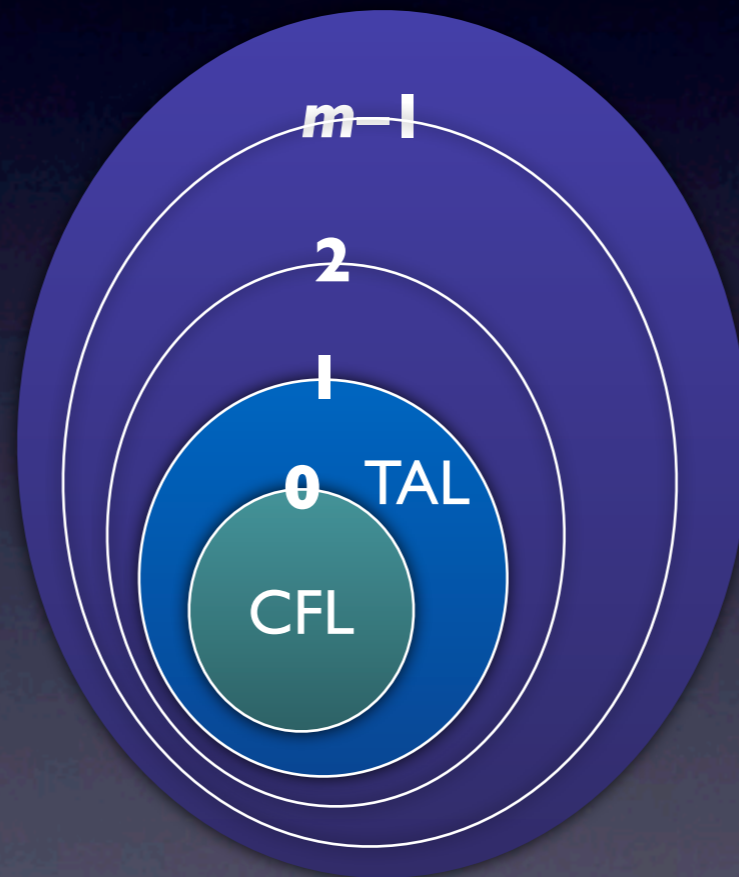
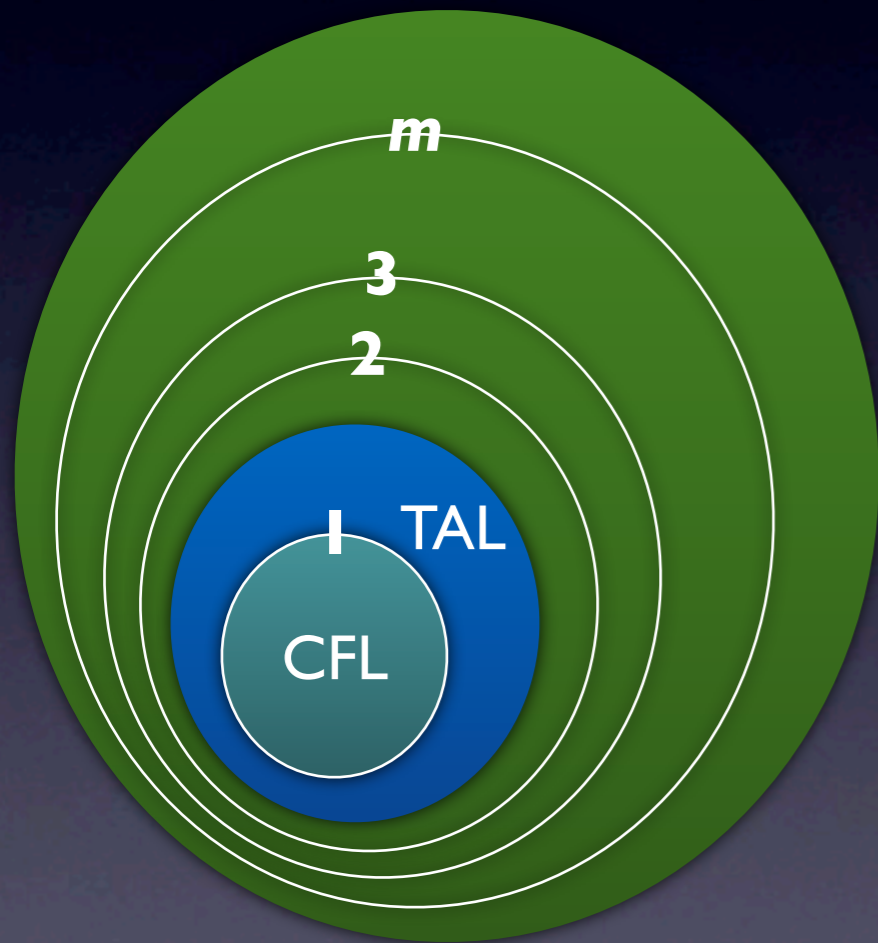
$$RESP_m \in yCFT_{sp}(2m-1)$$

MCFL vs. $yCFT_{sp}$



With Kanazawa and Salvati's (2010) theorem, we can see $2\text{-MCFL} - \text{MCFL}_{wn} \neq \emptyset$.
Improves known results.

Three Infinite Hierarchies



$$MCFL = \bigcup_{m \geq 1} m\text{-MCFL}$$

$$yCFT_{sp} = \bigcup_{m \geq 1} yCFT_{sp}(m-1)$$

$$\mathbf{c} = \bigcup_{k \geq 1} \mathbf{c}_k$$

Let's look at the third hierarchy.

Controlled Tree Languages

$$\text{CT}(L_1, L_2) = L_1 \cap (\widehat{L_2})^{*,c}$$

$$L_1 \subseteq T_\Delta$$

$$L_2 \subseteq \Delta^*$$

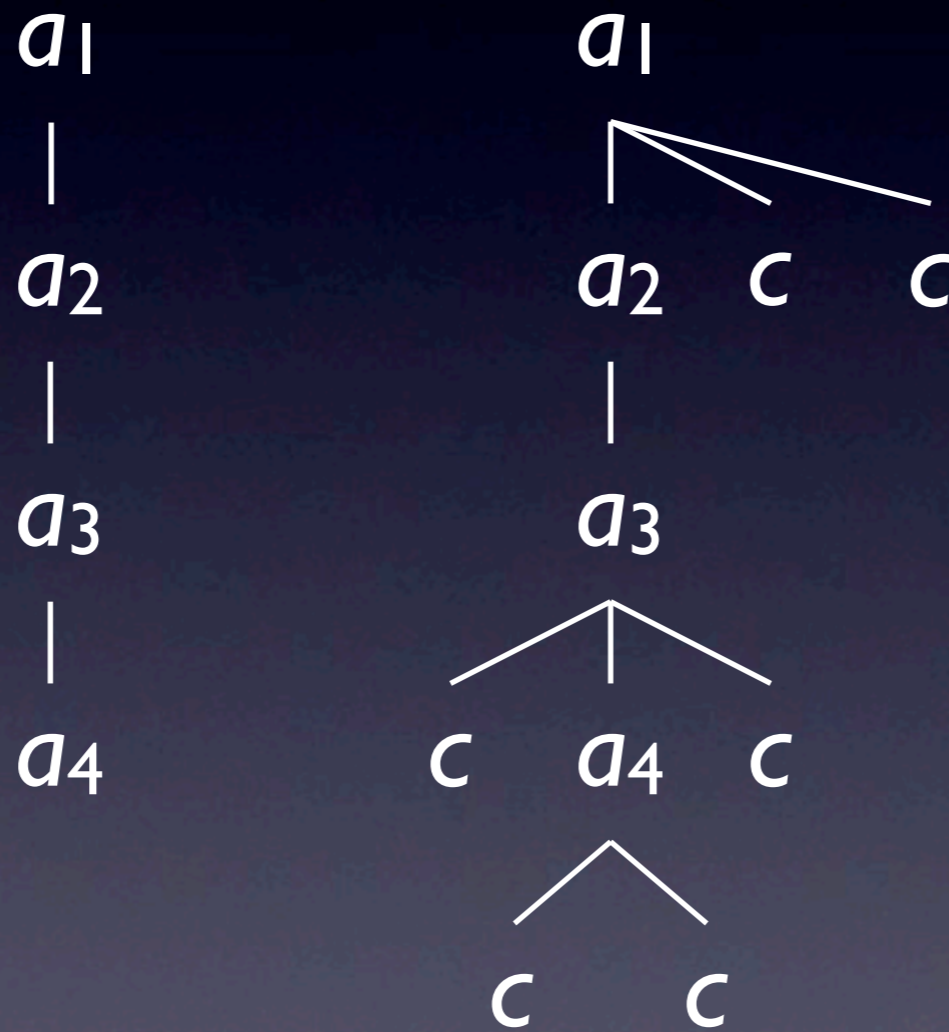
$$h : \Delta \rightarrow \mathbb{N}$$

$$d \in \Delta^{(r)} \Rightarrow h(d) \in \{0, \dots, r\}$$

Let's define an operation CT, which takes a tree language L_1 and a string language L_2 and returns a subset of L_1 . This operation is parametric on a function h .
Defined in terms of three operations: the hat operation, the tree analogue of the Kleene star operation, and intersection.

$$L \mapsto \hat{L}$$

$a_1 a_2 a_3 a_4$



$$a_1 \in \Delta^{(3)}, h(a_1) = 1$$

$$a_2 \in \Delta^{(1)}, h(a_2) = 1$$

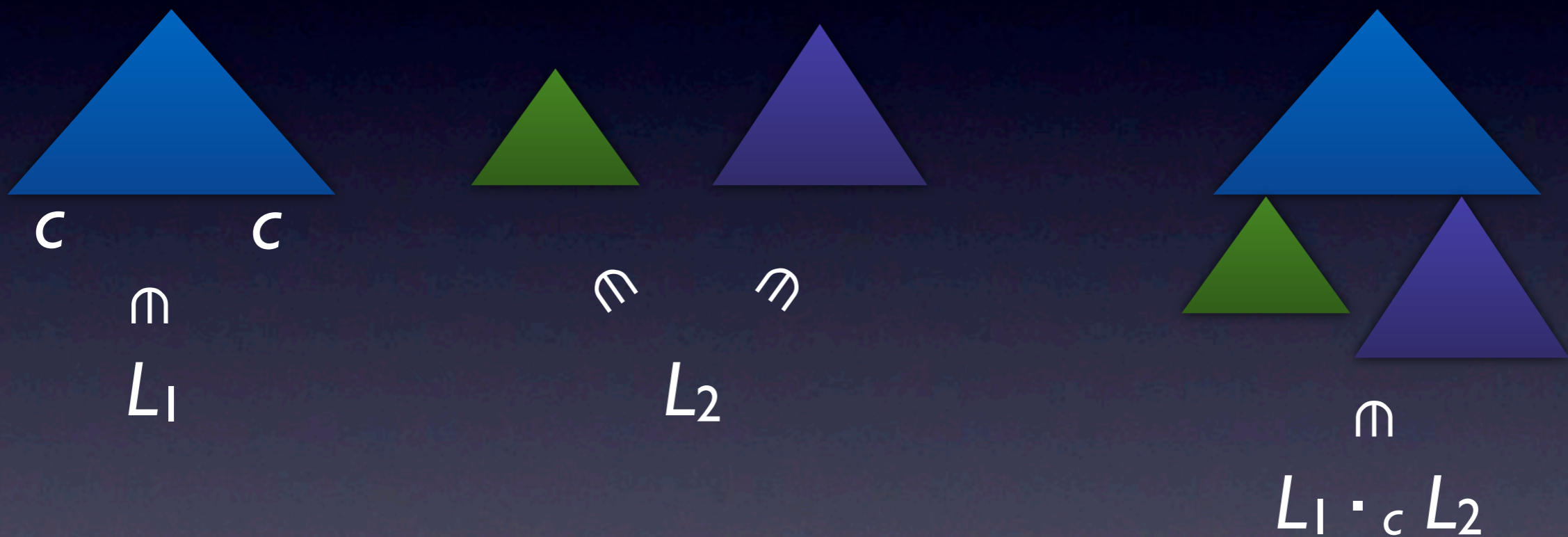
$$a_3 \in \Delta^{(3)}, h(a_3) = 2$$

$$a_4 \in \Delta^{(2)}, h(a_4) = 0$$

The hat operation turns a string into a tree.

Regular Tree Operations

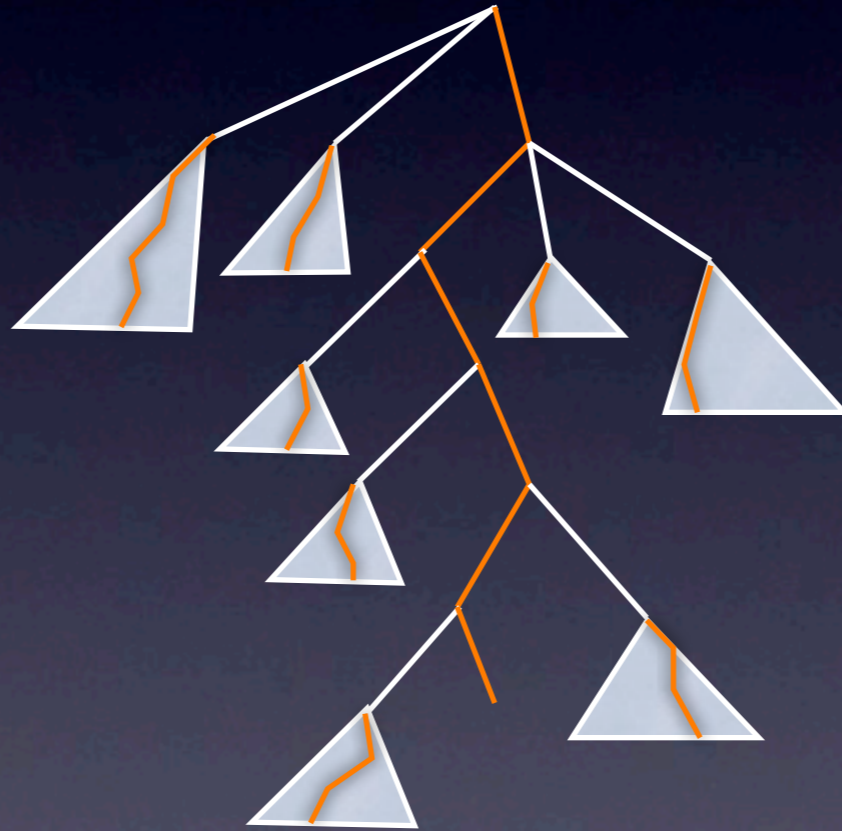
Concatenation



Kleene star

$$L^{*,c} = \{c\} \cup L \cdot_c L^{*,c}$$

$$\text{CT}(L_1, L_2) = L_1 \cap (\widehat{L_2})^{*,c}$$



This depicts the CT operation. You can observe elements of L_2 along some paths in a tree in $\text{CT}(L_1, L_2)$. These paths are determined by the h function.

$$\mathbf{CT}_1 = \text{LOC}$$

$$\mathbf{CT}_{k+1} = \text{CT}(\text{LOC}, y\mathbf{CT}_k)$$

$$\mathbf{C}_k = y\mathbf{CT}_k$$

$$\mathbf{C}_1 = \text{CFL}$$

$$\mathbf{CT}_2 \subseteq \text{CFT}_{\text{sp}}(1)$$

$$\mathbf{C}_2 = y\text{CFT}_{\text{sp}}(1)$$

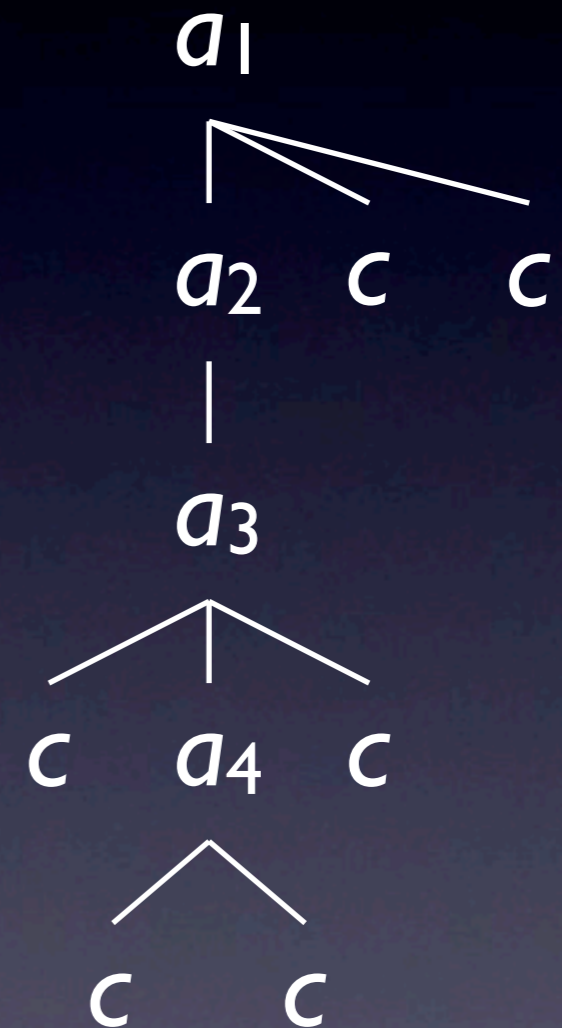
Define the Control Tree Language Hierarchy in terms of CT, starting from the class of local tree languages.

Weir's Control Language Hierarchy is the yield image of this hierarchy.

Let's show \mathbf{CT}_2 is included in $\text{CFT}_{\text{sp}}(1)$.

$$L \mapsto \hat{L}$$

$a_1 a_2 a_3 a_4$



substring

unary tree context

linear non-deleting
homomorphism

Going from strings to monadic trees, substrings are mapped to unary tree contexts, so a CFG is mapped to a $CFT_{sp}(1)$. Addition of nodes labeled by c is a simple case of a linear non-deleting homomorphism.

$$\mathbf{CT}_2 = \{ L_1 \cap (\widehat{L_2})^{*,c} \mid L_1 \in \text{LOC}, L_2 \in \text{CFL} \}$$

$$\subseteq \text{CFT}_{\text{sp}}(1)$$

$$\mathbf{C}_2 \subseteq \gamma\text{CFT}_{\text{sp}}(1)$$

$$\subseteq \text{2-MCFL}$$

$\text{CFT}_{\text{sp}}(1)$ is closed under linear non-deleting homomorphism, Kleene star, and intersection with regular sets. This shows that CT_2 is included in $\text{CFT}_{\text{sp}}(1)$.

$$\mathbf{C}_k \subseteq 2^{k-1}\text{-MCFL}$$

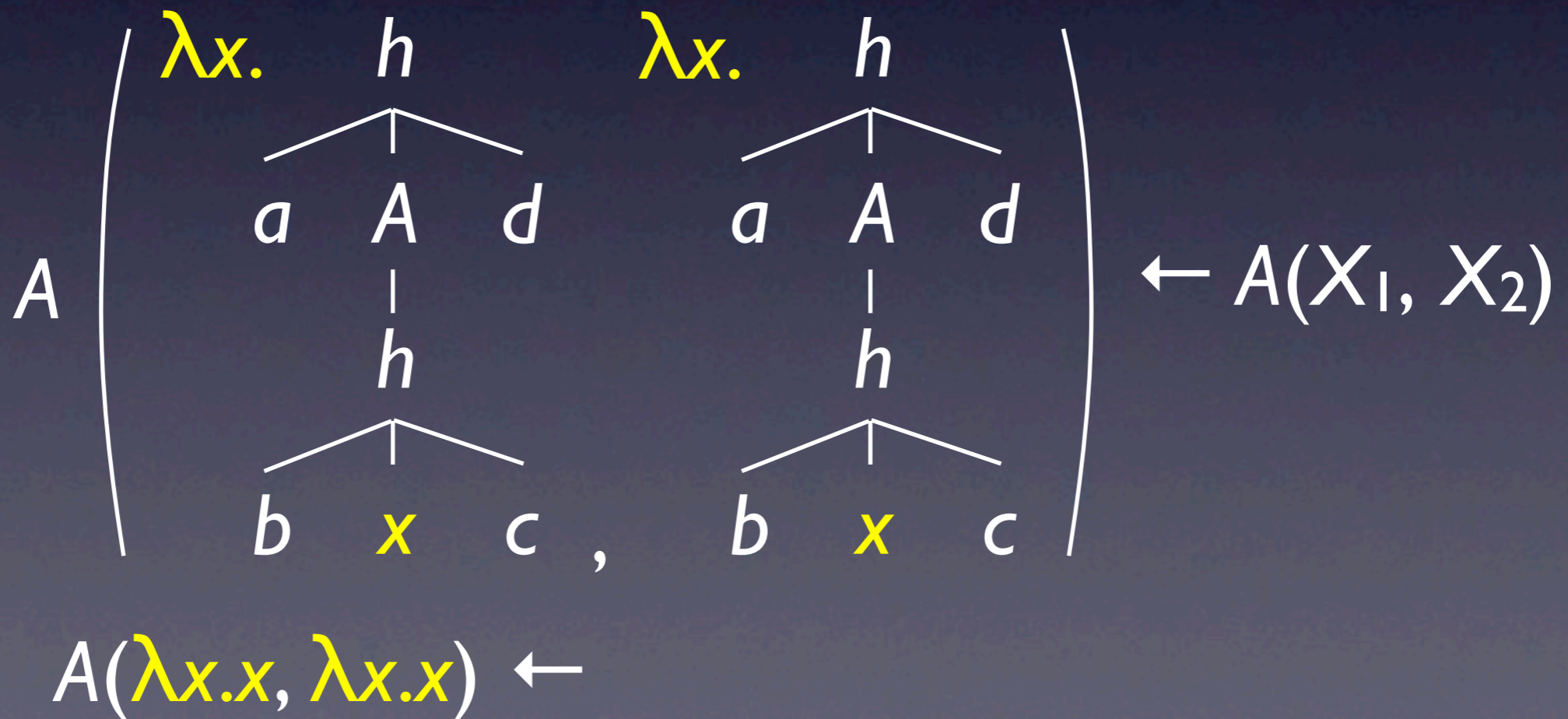
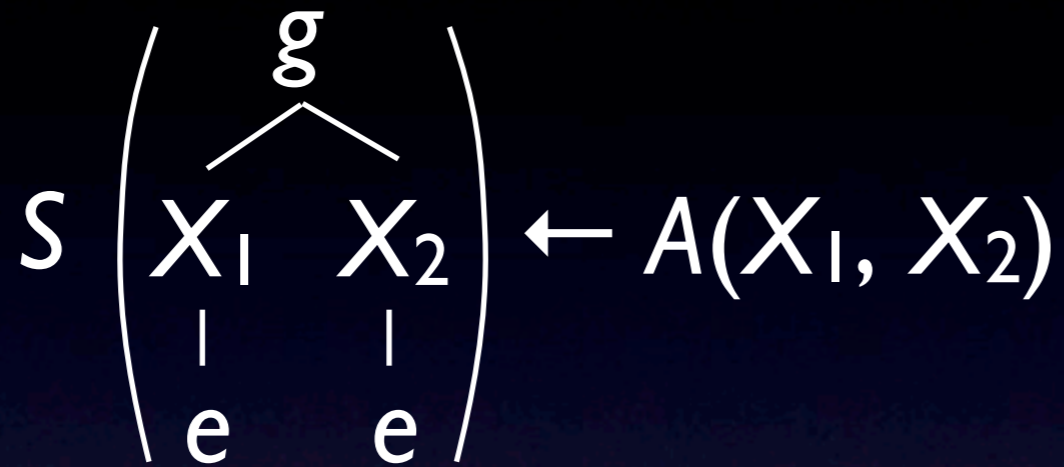
Kanazawa and Salvati 2007

$$\mathbf{CT}_k \subseteq 2^{k-2}\text{-MCFT}_{\text{sp}}(1) \quad (k \geq 2)$$

\approx MCTAG (Weir 1988)

The induction step is similar. Going from strings to monadic trees, substrings are mapped to unary tree contexts, and an MCFG is mapped to a “multiple monadic simple context-free tree grammar”. The latter is basically the same as a multicomponent TAG.

2-MCFT_{sp}(1)



Here's an example of a 2-MCFT_{sp}(1).

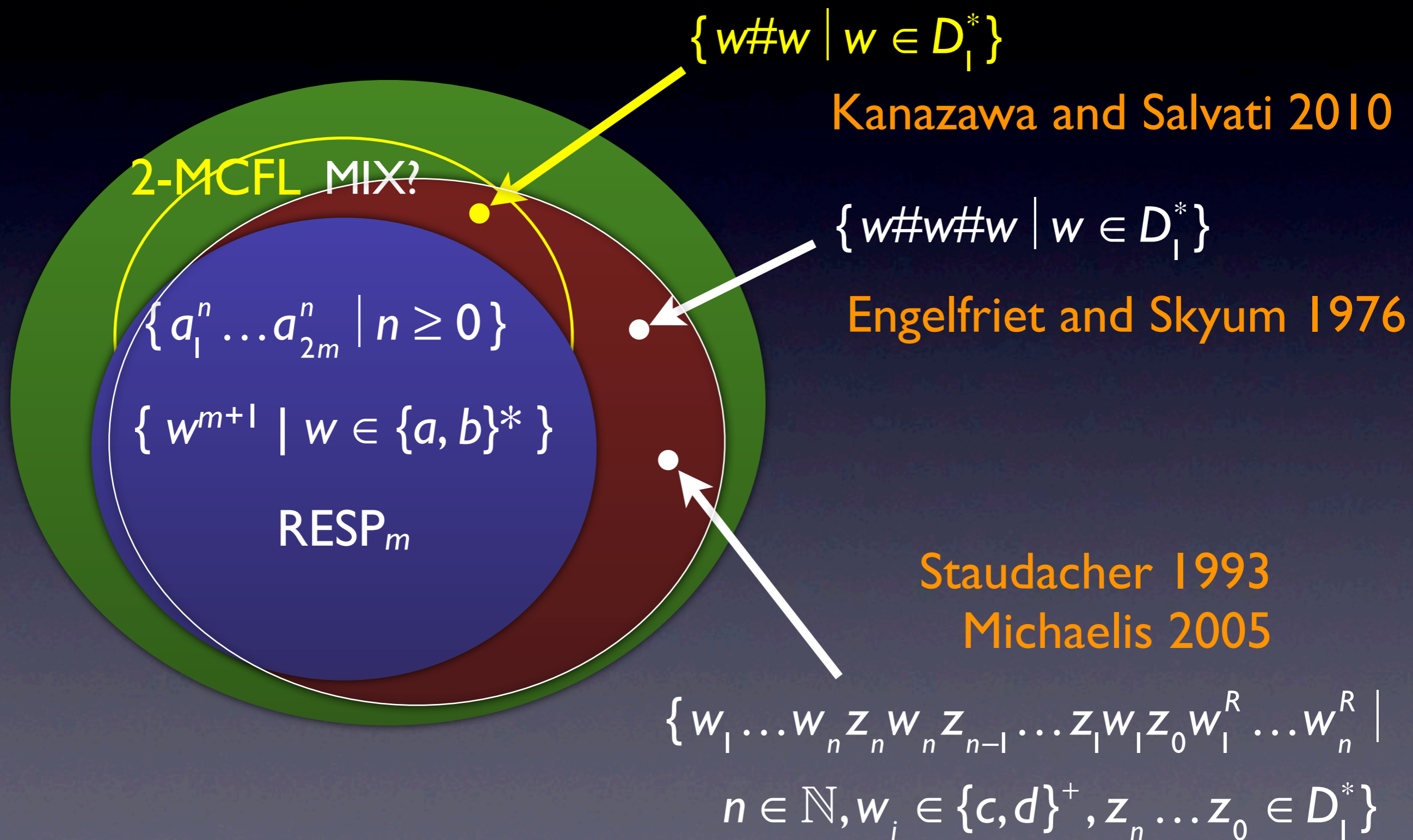
$$\begin{aligned}
\mathbf{CT}_{k+1} &= \{L_1 \cap (\widehat{L_2})^{*,c} \mid L_1 \in \text{LOC}, L_2 \in \mathbf{C}_k\} \\
&\subseteq \{L_1 \cap (\widehat{L_2})^{*,c} \mid L_1 \in \text{LOC}, L_2 \in \gamma(2^{k-2}\text{-MCFT}_{\text{sp}}(\mathbb{I}))\} \\
&\subseteq \{L_1 \cap (\widehat{L_2})^{*,c} \mid L_1 \in \text{LOC}, L_2 \in 2^{k-1}\text{-MCFL}\} \\
&\subseteq \{L_1 \cap L_2^{*,c} \mid L_1 \in \text{LOC}, L_2 \in 2^{k-1}\text{-MCFT}_{\text{sp}}(\mathbb{I})\} \\
&\subseteq 2^{k-1}\text{-MCFT}_{\text{sp}}(\mathbb{I})
\end{aligned}$$

$$\begin{aligned}
\mathbf{C}_k &= \gamma \mathbf{CT}_k \\
&\subseteq \gamma(2^{k-2}\text{-MCFT}_{\text{sp}}(\mathbb{I})) \\
&\subseteq 2^{k-1}\text{-MCFL}
\end{aligned}$$

The induction step goes like this.

The yield image of the tree language of an m -MCFT_{sp}(1) is the language of a $2m$ -MCFL.

MCFL vs. C



C is closed under copying.
 It's harder to separate MCFL and C.