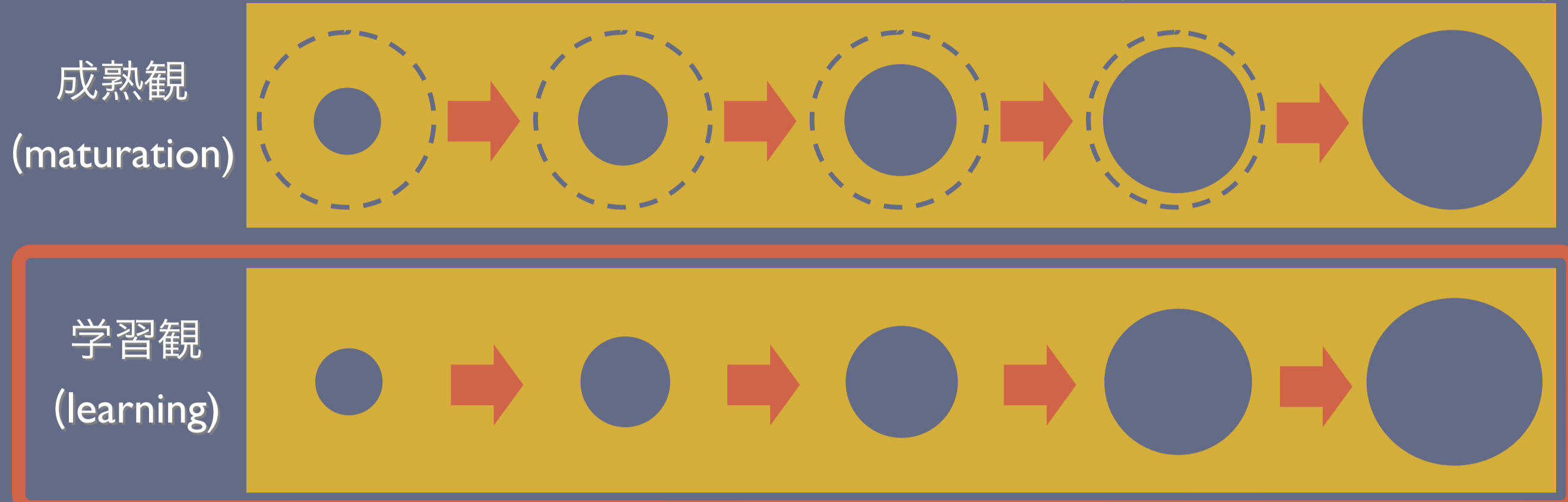


## <研究の背景と目的>

二つの言語習得観 (Bates & Elman 2002)



課題

刻一刻と変化する知識状態をどう記述し段階的な変化をどう計測するか

(憂うべき)現状

ほとんどこの検証は行われていない! (例外: Borenstajn et al. 2009)

そこで...

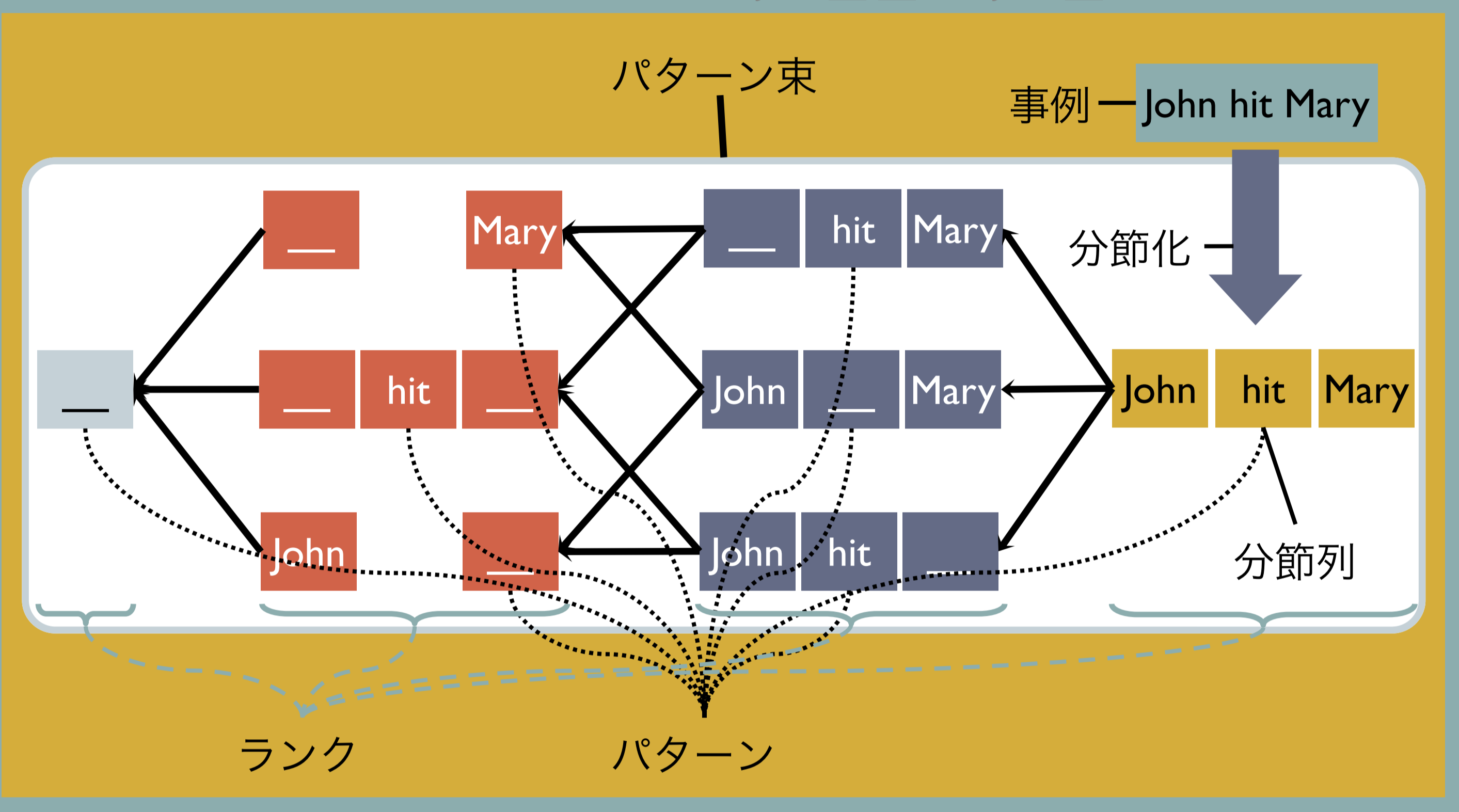
ならばやってみよう!

概要

- パターン束モデル (Pattern Lattice Model: e.g., 黒田・長谷部 2009)に基づき
- Brown コーパス(Brown 1973) in CHILDES (MacWhinney 2000)を使用し
- 幼児の発話から「パターン」= 統語知識の候補を生成し
- その「生産性 (Productivity)」をシャノンのエントロピーで算定し
- 年齢を経る毎に生産性が高まっていく段階的な発達プロセスを検証する

## <パターン束モデル (Pattern Lattice Model)>

- 黒田・長谷部 (2009)の提案したヒトの言語知識のモデル
  - ◎ (現時点では)主に言語形式の組織化のモデル
  - ◎ ヒトの統語知識は(ツリーではなく)「パターン」の集合であると考え
- パターン束
  - ◎ 事例 e の任意の分節モデル T による分節化 T(e) に対する再帰的変項化の産物
  - ◎ 内実はパターンの「集合」= パターン集合 P(e)
  - ◎ 部分一致に基づく「継承関係」(is-a)の規定された半順序集合 = パターン束 L(e)
- 簡単な例
  - ◎ e = John hit Mary. T = 単語分節 T(e) = [ John, hit, Mary ]
  - ◎ P(e) = { (John, hit, \_), (John, \_, Mary), (\_, hit, Mary), (John, \_), (\_, hit, \_), (\_, \_, Mary), (\_ , \_ , \_) }
  - ◎ 補足: 連続する変項を単一の変項に縮約 ((John, \_, \_) → (John, \_))

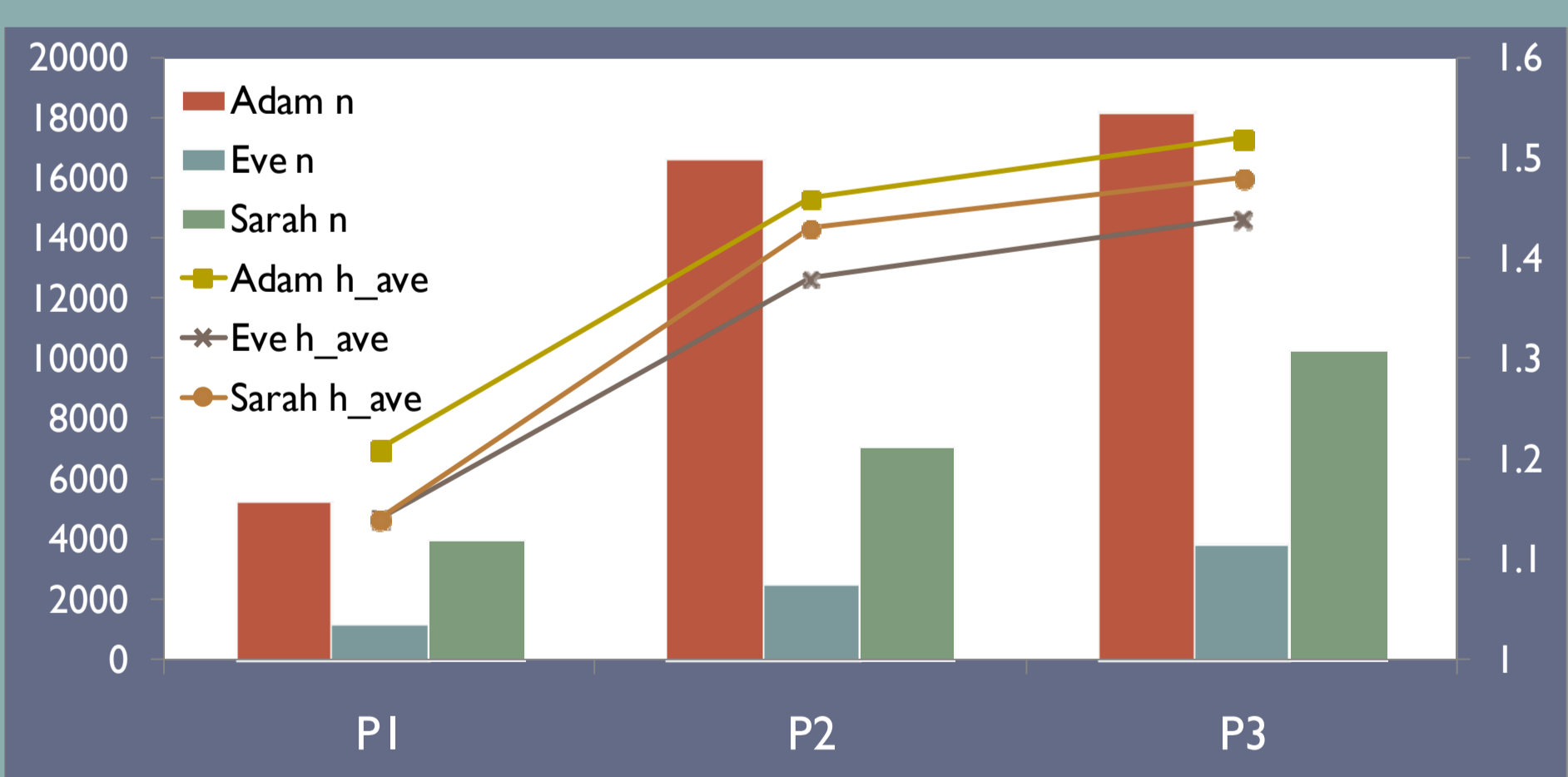


## CHILDES & Brown's Corpus

Files	Age	#sent.	vocab.	t/t	
Adam					
P1	1;16	2;3-2;11	11,184	1,407	0.056
P2	17-32	2;11-3;6	11,578	2,010	0.053
P3	33-48	3;6-4;5	9,071	2,006	0.055
Eve					
P1	1;7	1;6-1;9	3,485	669	0.102
P2	8-14	1;9-2;0	3,395	785	0.083
P3	15-20	2;1-2;3	3,535	958	0.087
Sarah					
P1	1;45	2;3-3;2	11,693	1,389	0.063
P2	46-90	3;2-4;1	8,384	1,706	0.075
P3	91-135	4;1-5;0	8,525	1,944	0.071

## <結果と考察>

● パターン数・生産性の平均値: 右下図・下図の通り



得られたパターンの例 (Adam, 頻度上位20位)

id	pattern	freq	h	rank
p3748	(what, _)	609	4.75	1
p258	(_, dat)	581	3.35	1
p60	(_, it)	413	6.94	1
p211	(_, a, _)	360	1.21	1
p83	(where, _)	333	6.39	1
p51	(I, _)	312	7.59	1
p47	(yeah, _)	304	0.00	1
p5398	(what, dat)	260	0.00	2
p81	(_, go)	239	6.09	1
p37	(_, dere)	223	5.35	1
p71	(no, _)	223	0.00	1
p29	(put, _)	216	6.99	1
p10	(who, _)	201	1.78	1
p119	(_, in, _)	200	1.59	1
p42	(Adam, _)	186	6.92	1
p591	(_, there)	185	6.35	1
p2798	(dat, _)	170	6.79	1
p2144	(_, dat, _)	161	1.51	1
p594	(_, in, there)	153	6.35	2
p1008	(_, it, _)	150	1.52	1

● 生産性の差の検定結果 ⇒

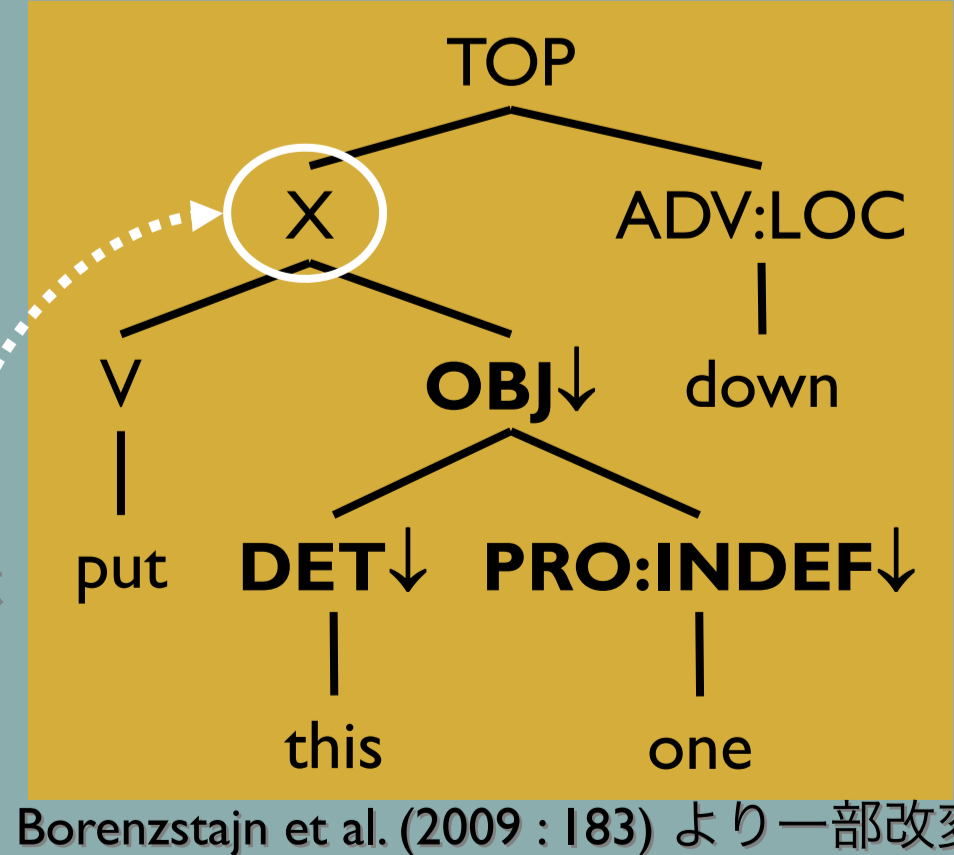
	Adam	Eve	Sarah
P1-P2	p-value 2.57E-117 W 34759997	p-value 7.74E-20 W 1212357	p-value 8.95E-83 W 11034700
P2-P3	p-value 4.16E-22 W 1.43E+08	p-value 5.55E-06 W 4515610	p-value 9.06E-09 W 34798423
P1-P3	p-value 2.36E-187 W 35466067	p-value 9.32E-38 W 1746899	p-value 2.02E-129 W 15237634

● エントロピー上昇は何の証左か?

- ◎ 幼児の発話における「体系的ばらつき」の増加 = 適度な部分一致を含むバラエティ豊かな発話への変化

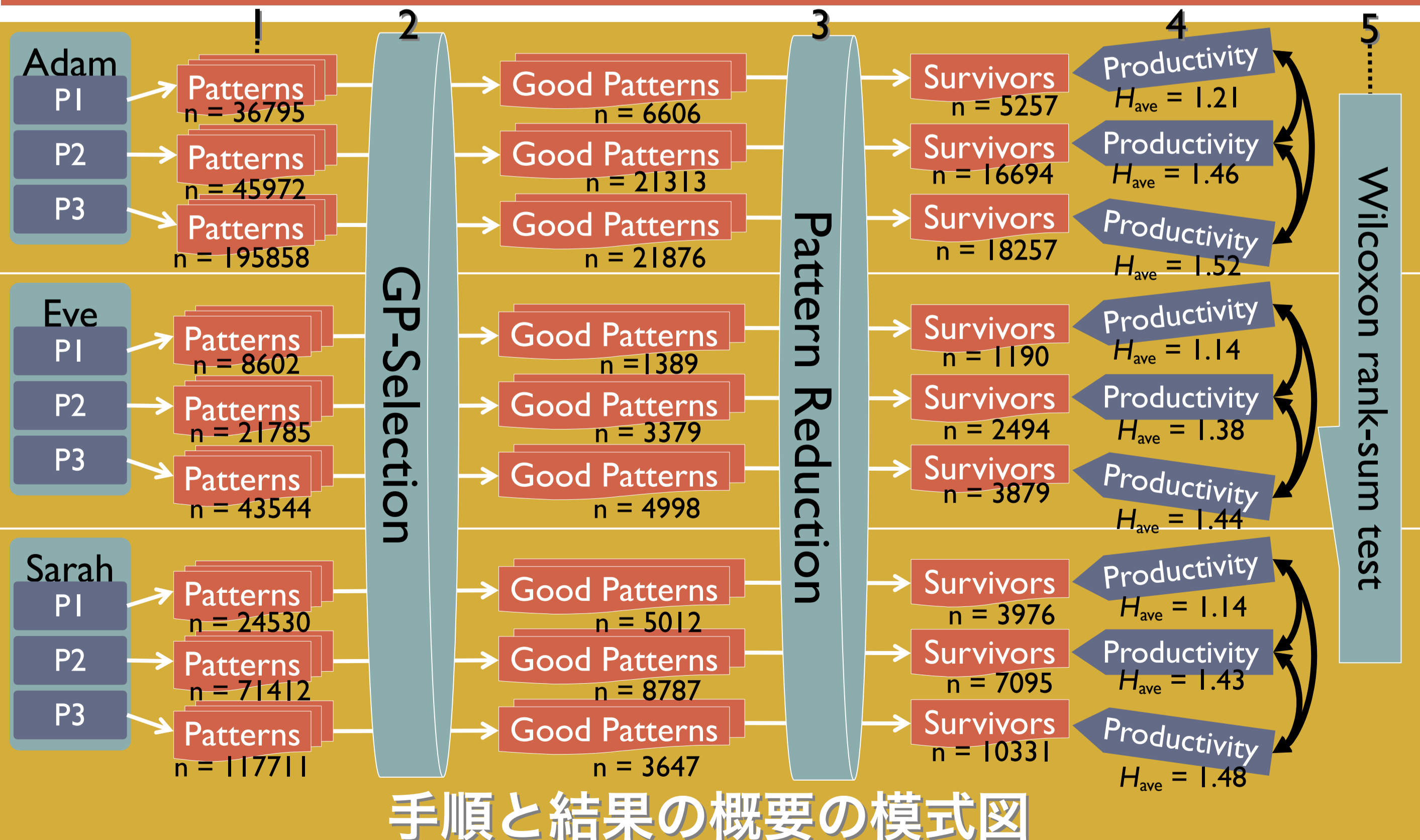
● 先行研究 (Borenstajn et al. 2009)との簡単な比較

- ◎ "Right representation"? : Tree Substitution Grammar = 木表示
  - ▶ 「木」より「束」の方が「弱い」かつ「網羅的」
  - ▶ 「抽象度」を過大評価している節あり



## <調査>

- データ
  - ◎ CHILDES (MacWhinney 2000) 内のBrownコーパス(Brown 1973)
    - ▶ 幼児 = {Adam, Eve, Sarah} の発話のみを抜き出し
    - ▶ 重複・言いさし、ポーズの含まれる発話を除外
  - ◎ 3幼児それぞれのデータをデータ量が揃うように3等分 = {P1, P2, P3}
- 方法
  - ◎ 3幼児 × 3データに対しそれぞれPLMのアルゴリズムによるパターン生成 ... 1
  - ◎ 頻度 ≥ 2 のパターンを選択 (「良いパターン(Good Patterns)」の選定) ... 2
  - ◎ 良いパターンから「パターン削減 (Pattern Reduction)」によりさらに選抜 ... 3
    - ▶ 削減1: ランクが一つ上 = 階層が一つ下のパターンの異なりが一つなら削除
    - ▶ 削減2: ランクが一つ上 = 階層が一つ下のパターンが全て良いパターンなら削除
  - ◎ 選抜パターンに対しシャノンのエントロピーを用いて生産性を算定 ... 4
- 生産性の算定
  - ◎ 算定方法は別紙 (Appendix 1) 参照
  - ◎ P1-P2間, P2-P3間, P1-P3間で生産性の平均の差の検定 (Wilcoxonの順位和検定) ... 5



### <謝辞>

黒田 航氏 (情報通信研究機構)  
 長谷部 陽一郎氏 (同志社大学)  
 井上 逸兵教授 (慶應義塾大学)  
 中村 文紀氏 (慶應義塾大学大学院)

### <主要参考文献>

- Brown, R. 1973. *A first language: The early stages*. Cambridge, MA.: Harvard University Press.
- 黒田航・長谷部陽一郎. 2009. Pattern Lattice を使った(ヒトの) 言語知識と処理のモデル化. 言語処理学会第15 回大会発表論文集(pp.670-673).
- MacWhinney, B. 2000. *The CHILDES project: Tools for analyzing talk*. Mahwah, New Jersey: Lawrence Erlbaum Associates.
- Bates, E., & Elman, J. 2002. Connectionism and the study of change. In Johnson, M., Munakata, Y., & Gilmore, R. O. (eds.) *Brain development and cognition: A reader* (2nd ed., pp.420-440) Oxford: Blackwell Publishing.
- Borenstajn, G., Zuidema, W., & Bod, R. 2009. Children's grammars grow more abstract with age: Evidence from an automatic procedure for identifying the productive units of language. *Topics in Cognitive Science, 1 (1)*, 175-188.

