



Comprehension outcores production in language acquisition: Implications for Theories of Vocabulary Learning

Dominic W. Massaro¹

University of California, Santa Cruz

Bill Rowe²

University of California, Santa Cruz

Received : 02.04.2015
Accepted : 08.09.2015
Published : 30.09.2015

Abstract

It is well-documented that children comprehend many more words than they are able to produce. Without exception, a child appears to understand various words that they do not use in their own speech. These results are used to test three different theories of speech perception. Motor theory assumes that motor processes are necessarily recruited for speech perception. A common representation theory claims that the same representation exists for both comprehension and production. Comprehension can be viewed as a recognition task whereas production can be considered a recall task, and recall is necessarily more difficult than recognition. The third theory assumes that speech perception and speech production and their acquisition cannot be based on the same underlying representation. Speech perception follows prototypical pattern recognition processes whereas speech production involves intricate motor processes that attempt to match a speech target. We analyze data from the MacArthur-Bates Communicative Development Inventories (CDI) to determine 1) if difficulty of articulation and parental input frequency influence perception and production equivalently and 2) assess whether equivalent representations can account for perception and production. The results falsify motor theory and common representation theory and support the pattern recognition account of speech perception.

Keywords receptive language, expressive language, vocabulary learning, speech perception, speech production

1. Introduction

A persistent question of how the child seamlessly learns language is being informed by a variety of empirical and theoretical research findings (e.g., Werker, Yeung, & Yoshida, 2012). All researchers acknowledge the seemingly impossible challenge that the infant faces in accomplishing this feat. A prototypical example is a contemporary version of Gavagai (Quine, 1960, 1990/1992). A mother tells her child that “the cat is on the mat”. How does a child make sense of the continuous speech stream to arrive at some veridical understanding? What are the words, their meanings and what does their combination signify? For the purposes of the present paper, it is important to realize that this illustration is framed in terms of the understanding of

¹ Department of Psychology, Santa Cruz, CA 95060 USA. Corresponding author massaro@ucsc.edu. Language acquisition of both speech and reading. Synergetic applications of behavioral science and technology.

² Santa Cruz Institute for Particle Physics, Santa Cruz, CA 95060 USA. wrowe@ucsc.edu. Consciousness, and the origin of language.

language as opposed to its production. But surely an equally perplexing challenge is how a child would produce the proposition, even when it is signified by a single-word utterance such as *cat*. Phrasing the anecdote in terms of understanding rather than production reflects the intuitive belief that language understanding paves the way for language production. We expect that a child would necessarily have to understand an utterance or a similar utterance before being able to produce it. The goal of this paper is to assess predictions of theories of speech perception and production based on the well-known asymmetry of the acquisition of perceiving and producing words. We now describe three extant theories of speech perception.

1.1 Motor Theory for Speech Production

Speech provides a natural domain for a motor-theory account because speech perception and speech production appear to be so tightly linked. A motor theory of speech perception contains different assumptions for different theorists but they all tightly link speech perception with speech production. Our interpretation of motor theory is illustrated in Figure 1. As seen in the figure, speech production processes somehow intervene in speech understanding. In one review, the theory claims that speech production processes are necessarily recruited for speech perception and comprehension (Galantucci, Fowler, & Turvey, 2006). Associated with this motor theory, it is also sometimes assumed that an interlocutor perceives the intended phonetic gestures of the speaker rather than the more variable auditory speech. Arbib (2012), a long-time advocate of the motor theory of speech perception, assumes that we recognize the sounds of speech by creating a motor representation of how those sounds would be produced (Moulin-Frier & Arbib, 2013). Motor theory is favored by many, especially those committed to an embodied cognition framework (Glenberg, 2008; Glenberg & Gallese, 2012).

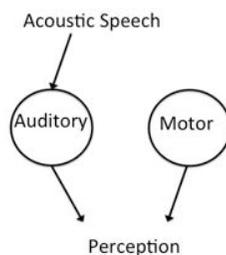


Figure 1. Schematic description of motor theories of speech perception in which a motor representation necessarily mediates speech perception.

We believe that it is fair to assume the engagement of the motor system illustrated in Figure 1 is assumed for all variants of motor theory. Massaro & Chen (2008) and Hickok (2009, 2014) question the viability of a motor theory. Hickok, Houde, and Rong (2011) agree that motor theory fails in its strong form. As evidence, they cite the accuracy of speech perception in individuals who have a variety of speech production deficits or temporally-induced interference with speech production processes. In addition, they emphasize that the putative existence of mirror neurons does not disqualify the existing evidence against motor theory.



1.2 Common Representation for Speech Perception and Production

Figure 2 illustrates a second theory claiming that speech perception and speech production rely on the same underlying representation. Using an information processing analysis, any difference between the two forms of vocabulary knowledge emerges because of the differences in the two tasks required of the child (Huttenlocher, 1974). Comprehension can be viewed as a recognition task whereas production can be considered a recall task, and recall is necessarily more difficult than recognition. Using a signal detection framework, comprehension exists in a context with just a few response alternatives whereas recall occurs in a context of many response alternatives. In a recent formal model of McMurray, Horst, and Samuelson (2012), the child has many more competitor alternatives in production than in comprehension, and these alternatives necessarily lower performance in production relative to perception. The bottleneck in production relative to comprehension putatively occurs because there are many more competing production alternatives than comprehension alternatives. This difference exists even though the quantity and quality of the underlying representation is assumed to be equivalent in the comprehension and production scenarios.

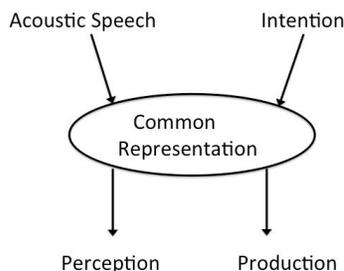


Figure 2. Schematic illustration of the view that speech perception and production use the same common representation.

McMurray et al.'s (2012) simulated results indicated that at a given stage of learning (20,000 trials), a production task was accurate on just 5 words whereas the comprehension tasks with 3, 5, and 10 alternatives were accurate on 32, 28, and 25 words, respectively. In addition, the production task gave roughly the same result as the 35-alternative comprehension task. We agree with the authors that receptive vocabulary will necessarily be larger than expressive vocabulary when the child uses the same representation for both perception and production. We question, however, whether the same representation is in fact used for these two functions. It might be too parsimonious to attribute differences between receptive and expressive vocabulary to simply task differences, and assume that “memory demands, difficulty planning articulation, or the earlier age at which speech perception develops” (McMurray et al., 2012, p. 845) do not play a role. Our analysis, therefore, focuses on whether it is reasonable to assume that equivalent representations are used for perception and production. In our data analysis, we ask whether simply an equivalent representation can account for the discrepancy between perception and production of individual children.

It is possible that an incomplete memory representation would allow accurate comprehension but not accurate production of some words. As illustrated in Figure 2, a receptive language advantage over production putatively occurs because the child has a sketchy representation that allows her to recognize a word but not recall it for accurate production. This test centers on whether there exists a memory representation that allows the child to understand a word but not produce it. The results give for each child the words she produces and the words she understands. Are there possible memory representations of these words that would allow correct understanding of the words but production of only a subset of these words? An ideal memory representation for the description in Figure 2 would be one that allowed perception of all of the words but production of only some of them. If a child comprehended only the words *mommy* and *bye* but did not produce them, a partial memory representation could account for the results. For example, *mommy* could be represented by a nasal sound and *bye* by an /ai/ as in the sound of the word *eye*. These two representations would allow recognition of the words but they would be insufficient for accurate production.

1.3 Pattern Recognition: Different Processing for Perception and Production

Figure 3 illustrates a class of explanations that assume that speech perception and speech production are fundamentally different functions and, therefore, their acquisition cannot be based on the same processing or the same underlying representation. In this view, speech perception follows prototypical pattern recognition processes whereas speech production involves intricate motor processes that attempt to match a speech target (Massaro, 1987, 1998; Guenther, 1995). The pattern-recognition explanation follows from language understanding research, which accounts for comprehension without regard to production processes (Massaro, 1975, 1998; Movellan & McClelland, 2001). An emerging fundamental principle is that there are highly analogous processes involved across speech perception, reading, and higher-order language processing (Massaro, 1975, 1987, 1998; Massaro, 2005). Language processing is a form of pattern recognition, is influenced by multiple cues or sources of information, and is quantitatively described by the Fuzzy Logical Model of Perception (FLMP; Christiansen et al., 1998; Movellan & McClelland, 2001). For example, the processing for perception in face-to-face situations can include information about visible speech, not just auditory speech, and information about many linguistic and contextual cues (Massaro, 1998). Generalizing from this view, the assumption is that comprehension develops somewhat independently of production so that processes involved in speech production cannot account for comprehension ability. This model considers that speech perception and speech production are served by separate mostly independent systems.

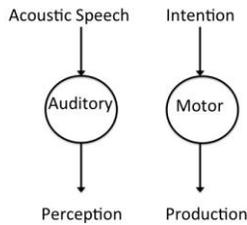


Figure 3. Schematic illustration of the view that speech perception and production use different representations.

These same pattern-recognition processes also occur in language acquisition, not just in accomplished language users (Fennell & Waxman, 2010; Gogate & Hollich, 2010; Hollich et al., 2000; Massaro, 1987, Chapter 8). An Emergentist Coalition Model describes how children rely on multiple cues over development in the mapping of words onto referents (Golinkoff & Hirsh-Pasek, 2008; Hirsh-Pasek et al., 2000). The use of and the weight given to these cues change across development. For example, infants initially rely mostly on perceptual cues and gradually begin to use a speaker's intent and linguistic cues to determine word reference. These cues and constraints are graded (not categorical) in nature, suggesting further that they must be combined to give a more reliable understanding of the input. Evidence to date indicates that this combination process is highly efficient or optimal, as described by Bayes Law (Massaro, 1987, 1998; Massaro, 2008). This type of Bayesian inference is consistent with predictive models, which are gaining in popularity in recent interdisciplinary accounts of behavior (Clark, 2013; Friston, 2010). If correct, a pattern recognition model of word comprehension would challenge the two views that production somehow intervenes in the comprehension of a word or that comprehension and production rely on exactly the same representation.

One formalization of this pattern recognition model at the neural level is the dual-route model that has ventral and dorsal processing streams (Hickok, 2008). The ventral stream processes the speech input for comprehension and utilizes structures in the superior and middle portions of the temporal lobe. The dorsal stream responsible for speech production utilizes the posterior planum temporale region and posterior frontal lobe. This stream maps the acoustic signal into an articulatory representation. It is also assumed that this stream can also generate articulatory representations with other sensory inputs such as visible speech and also without any input at all. Similar to the Bayesian approach just mentioned, the comprehension system is highly predictive in that it automatically exploits multiple sources of information to facilitate speech understanding.

1.4 Existing Evidence

We know that comprehension is accomplished well before production of many words. It is well-documented that children comprehend many more words than they are able to produce (Bates, 1993; Fenson et al., 2000). Without exception, a child appears to understand a variety of words that they do not use or use correctly in their own speech. Although there is no apparent controversy about the vocabulary advantage of perception over

production, its explanation is still somewhat elusive. In this paper, we evaluate several possible explanations against a large database of results. As described, the explanations differ primarily in terms of whether speech production is somehow recruited for speech perception to succeed, whether the same representation is involved in perception and production, or whether speech perception is fundamentally pattern recognition and occurs relatively independent of speech production processes.

The MacArthur-Bates Communicative Development Inventories (CDI) has been used successfully to measure comprehension and production vocabulary. It uses a checklist to ask parents to report their child's word comprehension, word production, and grammar. One version provides norming data on age of acquisition for 396 individual words collected from the parents of children ages 0;8 to 2;6. For each word on the checklist, there are two boxes that can be checked. The check boxes are labeled 1) understands and 2) understands and says, respectively. The parent either does not put a mark in one of the two check boxes or marks just one. This design of the inventory form precludes a parent from indicating that a child produces a word but does not understand it. Thus, the results indicate whether each word on the list is understood by the child or is understood and produced by the child. A produced word requires that the child's utterance could be understood out of context by at least a parent or caregiver. The child is also not given credit for produced words that are simply imitated, as in the case when the parent simply asks the child, "Can you say 'banana'?" These norms are available on the internet (<http://www.sci.sdsu.edu/cdi/cdiwelcome.htm>) and support the well-known finding that vocabulary comprehension far outpaces vocabulary production (Bates & Goodman, 1999; Fenson et al., 2000).

Notwithstanding various criticisms, the validity of these measures has been substantiated in various investigations (see Fenson et al., 1994; for a review; Fenson et al., 2007). One relevant issue is that it would not be surprising to find that parents might overestimate what their children know (Tomasello & Mervis, 1994). Particularly relevant to our inquiry about differences between perception and production, a parent might indeed score a word as produced even though their child mispronounced it. We believe, however, that this would not impact the results because a word produced is also scored as understood in the scoring of the CDI inventory. The principal metric to be used in our analyses is the difference between words understood and words produced. Thus, if anything, a liberal inclusion of words produced would diminish any difference between words produced and words understood.

The validity of the CDI procedure has also been substantiated by positive correlations with experimental tests of comprehension. Bates, Bretherton & Synder (1988) found that the CDI measure of comprehension correlated with an experimental test of comprehension for infants aged 1 year, 1 month. Further, comprehension at 1 year, 8 months correlated with a laboratory test (Dale, Bates, Reznick & Morisset, 1989). The strongest evidence for the validity of the CDI comes from event-related potentials (ERPs) of infants to auditorily-presented words that were either checked as understood or not (Mills, Coffey-Corina & Neville, 1993, 1997). The infants ERPs distinguished between these two classes of words as scored on the CDI.



The average production vocabulary has been analyzed as a function of comprehension vocabulary, as measured by the CDI (Fenson et al., 1994). Although there is a great deal of variability across children in both comprehension and production vocabulary, comprehension of words far exceeds their production throughout the first years of language development. An average child already understands 200 words by the time she or he produces 50 words.

Another study included 659 children whose age in months ranged from 8 to 18 months. For each word on the checklist of 396 items, the parent indicated which of the words on the list are understood by the child or understood and produced by the child (MacArthur-Bates Communicative Development Inventories American Cross-sectional CDI studies: Words & Gestures (American CDI, 2014; Dale & Fenson, 1996). At 8 months of age, 50% of the children in this sample comprehended the words *mommy*, *daddy*, *bye*, *peekaboo*, *bottle*, and *no*. None of these words was uttered by 50% of the children. At age 16 months, 175 words were comprehended and 20 words were produced by 50% of the children in the sample. Thus, there is a dramatic discrepancy between words understood and produced.

The conclusions from the CDI are supported by rigorous laboratory experiments. Similar to the MDCDI inventory, studies of word learning also find a big discrepancy between the learning of receptive and productive language. For example, Hahn and Gershkoff-Stowe (2010) found that two and three year olds were 60 and 70% correct in a four-alternative task in choosing the correct object for a previously paired nonsense word, but only 5 and 12% correct in producing the label. Similar results were found in a second study with eight alternatives. It is possible, however, that the receptive advantage over production is simply a guessing advantage when given a small number of response alternatives, consistent with the Common-Representation theory (Figure 2). A motor theory (Figure 1) might offer a similar explanation.

How might the observed acquisition asymmetry in the perception and production of words inform these theories of speech perception? As already noted, analyses to date have fallen short of any definitive choice among the theories. The current paper extends these analyses in several ways. When possible, the analyses will be carried out on individual children to better measure an individual child's perception and production capabilities at a given stage of language acquisition. An analysis on an individual child's results should be able to test whether the same representation of a word is sufficient to account for both perception and production. For example, showing that a word can be understood even though none of the segments in the word are produced would seem inconsistent with motor theory (Figure 1) in which production processes are assumed to mediate speech perception. Similarly, this same result would question whether perception of a word while failing to produce it can be described by the same memory representation (see Figure 2). On the other hand, this example result would be consistent with a pattern recognition process that assumes that speech perception is not constrained by production processes (Figure 3). We now carry out a more thorough and detailed analysis of existing results of

individual children to evaluate the three theories of speech perception we have described.

2. Methodology

Our caveat is that to be completely valid, data analysis should be applied to an individual child's words rather than to group results. The database consisted of 396 words from the Words & Gestures American, organized in different words categories such as animal sounds and vehicles. There were 1089 unique cases in which the words produced and the words understood were tabulated for a given child at a given age (Fenson et al., 2007). The ages of the children in months ranged from 8 to 18 months. Each case was from a different child and the first time the CDI form was completed by a parent.

Table 1 gives a summary of the number of words understood and produced by the 1089 individuals divided into 6 groups based on the number of words produced. The range and number of words understood far exceeds the comparable results for words produced. As an example, there were 176 cases in which no words were produced, and only 5 of these cases also had no words understood. The remaining 171 cases of children with no expressive words at all had between 1 and 341 words understood with an average of 45.6 words understood. A similar discrepancy is apparent for the other five groups. The question remains whether motor theory or the task difference postulated by Huttenlocher (1974) and McMurray et al. (2012) could account for these huge discrepancies with the assumption that the same representation was used for perception and production.

Table 1. Six partitions of the individual results of 1089 children based on the number of words produced, giving the range and average of the number of words produced, the range and average number of words understood, and the number of cases contributing to each partition of the analysis.

Range Produced	Average Produced	Range Understood	Average Understood	Cases
0	0	0-341	45.63	176
1-3	1.91	1-261	58.70	194
4-7	5.39	7-243	76.74	183
8-17	12.59	41-280	116.21	172
18-49	29.29	45-261	168.44	211
51-376	120.11	112-396	264.69	153
Overall	28.22	Overall	121.74	1089

We counted the number of times each child understood and said the initial and final consonants of words based on the words they understood and said.



There were 1089 children with measures of words understood and words said. We determined the number of times each child understood a given consonant and the number of times the child said the consonant based on his or her words that were checked as 1) understood and 2) understood and said, respectively. Recognizing or producing a word boils down to having a memory representation of at least some of the segments that make it up. As mentioned, there were 171 children who understood some words but did not produce any. How is a child able to accurately distinguish dozens or even hundreds of words without being able to produce any of them? Recognizing (distinguishing among) a large number of words would require a fairly complete memory representation. If a child uses the same memory representation for perception and production, we should have expected that the child should be able produce at least some of the words.

3. Findings

3.1 Individual Subject Analyses.

To give an example of the results for an individual child, Table 2 gives the words understood and produced for a child who understood 121 words and produced only 17. Table 3 gives the number of times speech segments occurred in initial position for understood and said words for this child. As can be seen in the table, this child was able to understand words with all of the 22 possible initial phonemes but could produce words that begin with

Table 2. The words understood and produced (Table 2a) for a child who understood 121 words and produced only 17.

all gone	bubbles	ear	Ice cream	ouch	show	tongue	Woof woof
baa baa	bye	eat	juice	outside	sky	touch	yes
baby	careful	eye	jump	owie	slide	toy	yum yum
bad	car	finger	kiss	park	spaghetti	toothbrush	zipper
balloon	cat	fish	kitty	peas	splash	tummy	
ball	chicken	foot	leg	peekabo	spoon	tv	
banana	cherios	gentle	light	plant	stop	uh oh	
bathtub	cookie	girl	look	play	stroller	up	
bath	cracker	good	meow	please	swing	vroom	
bed	cup	go	milk	pull	teddy bear	walk	
bite	daddy	grandm a	mommy	puppy	show	wanna	
belly button	dance	hair	moon	push	teddy bear	wash	
blow	diaper	hand	mother	run	telephone	watch	
book	dog	help	name sitter	say	thirsty	water (2 meanings)	
bottle	doll	hit	night night	see	thank you	what	
boy	door	hot	nose	shh	throw	who	
broom	drink	hug	no	shoe	toast	wipe	

Table 2a. The words produced

baabaa	baby	banana	bottle	bubbles	bye	cat	daddy	dog
eye	kitty	mommy	nose	uhoh	vroom	what	woofwoof	

Table 3. The number of times initial phonemes were understood and produced for the child in Table 2 who understood 121 words and produced only 17.

Phonemes	Example Word	understood	said
b	bin	18	6
tS	chin	2	0
d	date	7	2
D	this	0	0
f	fax	3	0
g	gap	4	0
h	help	7	0
j	yacht	2	0
dZ	gin	3	0
k	king	8	2
l	leg	3	0
m	mail	5	1
n	nose	3	1
p	pin	9	0
9r	ring	1	0
s	sing	10	0
S	shop	3	0
t	tip	8	0
T	thing	2	0
v	very	1	1
w	will	8	2
z	zip	1	0

only 7 different phonemes. This child also understood words that begin with the consonant clusters /dr/, /sp/, /str/, and /sw/. The only consonant cluster that the child produced was /vr/ in the iconic vroom sound. It is clear that the child was able to recognize words with various consonants even though none of the words produced contained these consonants. For example, this child understood the words puppy, push, hug, help, and hot, but did not say words with /p/ or /h/ in initial position or words with /p/, /sh/, /g/, or /t/ in final position. This discrepancy makes it unlikely the advantage of comprehension can be explained by motor theory or common representation theory. Given that the child did not articulate many of the segments of words that she understood, it is unlikely that motor processes intervened in the perception of these words. Similarly, as evidence against the common representation theory, we claim that it is not possible to create a minimal representation of each of the 121 words that the child understood that would allow for comprehension but not production. To distinguish among all of the words in the child's vocabulary she would require essentially a nearly complete representation for each word.

It might be argued that the child exploits situational context in word understanding. This possibility is certainly predicted by models that postulate multiple influences such as context in speech perception. Bayesian



models such as the FLMP have accounted for the simultaneous influence of bottom-up and top-down influences in language understanding (Massaro, 2012). Context might somehow facilitate the recognition of a word but not necessarily the pronunciation of that same word. The differences in the counts of understood and said words should not be due to context, however, because the MacArthur Bates Inventory requires the parent to indicate words that are produced or understood independently of context. In addition, it is not obvious how context effects would benefit perception more than production. Given that there are so many words that the child understands but is not able to produce in a way that can be understood, then the child's putative adumbrated representation of these words (McMurray et al., 2012) cannot be rich enough to mediate their perception. Similarly, these results are equally damaging to the central assumption of motor theory (Galantucci et al., 2006) that motor processes necessarily mediate speech perception. Thus the analysis of individual children's receptive and expressive language appears to be most consistent with a pattern recognition account that assumes speech perception occurs independently of speech production processes.

3.2 Group Analyses.

For each unique word in our database, we pooled the results across the 1089 children. The total number of children that were scored as saying the word and understanding the word was treated as dependent measures, as was the difference between understood and said words. Our goal was to determine what variables might influence saying or understanding a word and whether there might be dissociations between perceiving and saying a word. On the other hand, correlations between perception and production might appear to be a more direct method to assess the relationship between understood and said words. If understanding a word were intricately tied to its production (see Figure 1), then we would expect a high correlation between the words that are understood and said. As expected, the Spearman correlation between said and understood words was .757, $p < .001$.

There is necessarily a positive correlation between perception and production in this database because all words that are checked by the parents as said are also scored as understood. A positive correlation between perception and production might also reflect the influence of a third variable such as frequency of occurrence in the parents' speech. If such a variable influences speech perception we might also expect it to influence speech production even if perception and production are based on separate processes (Figure 3). For example, it is possible that consonants that are easier to perceive are also easier to produce. Thus we expect a strong positive correlation between perception and production for several possible reasons even if their processes are unrelated. Thus, we cannot use the correlation between said and understood words as deciding among the three theoretical alternatives. Our individual subject analyses have already demonstrated that a child understands many words that she is not able to produce and that the child's putative adumbrated representation of these words cannot be rich enough to mediate their perception. Another tact to distinguish among the theories

would be to find a variable that differentially influences said and understood words. One possible variable is difficulty of articulation.

3.2.1 Correlations with Difficulty of Articulation.

Another possible test of the theories is to determine how a child's comprehension and production of words relates to the difficulty of articulation of segments of the words. If speech production processes are central to word acquisition, then we would expect words easiest to produce would tend to be both perceived and produced before words more difficult to produce. Further, if the same representation of a word was used for both comprehension and production or if production mediated comprehension, then we would expect that difficulty of articulation would inversely correlate with both comprehension and production. On the other hand, if perception and production processes are mostly independent of one another then it is possible that difficulty of articulation would influence production much more than it influences perception.

Although there are currently no established measures of difficulty of articulation (or difficulty of perception), we created a metric based on several relevant studies (Kirk, 2008; Kirk & Demuth, 2005; McAllister Byun, 2012; Rvachew, Chiang, & Evan, 2007; Smit 1993; Smit et al., 1990; Stoel-Gammon, 1987; Stoel-Gammon & Buder, 1999). Difficulty of Articulation of the consonant segments was defined as a 1-7 value on a scale of easy to difficult. These difficulty values for each of the phonemes are shown in Table 4. The metric for defining the difficulty of articulation of a word was simply the sum of the difficulty of articulation of each of its consonants. In this case, larger values correspond to a more difficult articulation. Although other factors such as coarticulation might also influence difficulty, there is even

Table 4. Two difficulty of articulation measures for each consonant used in computing the difficulty of articulation of each word (given by the sum of the difficulty measures of all of the consonants in the word).

Phonemes	Example Word	Difficulty of Articulation	Shriberg Acquisition Measure
b	bin	1	1
tS	chin	7	2
d	date	1	1
D	this	3	3
f	fax	2	2
g	gap	2	2
h	help	1	1
j	yacht	2	1
dZ	gin	6	2
k	king	2	2
l	leg	4	3
m	mail	1	1
n	nose	2	1
N	long	7	2
p	pin	1	1
9r	ring	7	3
s	sing	7	3
S	shop	6	3



t	tip	2	2
T	thing	3	3
v	very	5	2
w	will	3	1
z	zip	7	3
Z	vision	7	3

less guidance on creating a metric for this factor. Similarly, there is very little basis to create difficulty values for the vowels. Accordingly, this measure of difficulty of articulation for a word is based on only the consonants in the words. The Difficulty of Articulation values for each of the words are available on request.

3.2.2 Correlations with Parental Input Frequency.

We evaluated the influence of difficulty of articulation simultaneously with parental input frequency. Our measure of parent input frequency was determined from a database that can be searched at ParentFreq (2014). We initially evaluated both linear and log frequency but as expected found that log frequency gave the best predictions. Therefore, for ease of presentation, we use only log frequencies in our analyses. We analyzed a possible influence of difficulty of articulation in the data set of the 1089 children described previously in our individual subject analyses. After combining words that were repeated twice on the check list, there were 386 words in this data set. The dependent variables were the total number of children that understood each word, said each word, and an adjusted understood measure derived from subtracting the number of children who said a word from the number of children who understood the word. Each word was also assigned a word difficulty measure and linear and log parental input frequency.

We carried out Spearman and Pearson correlations among difficulty of articulation, parental input frequency, number of children that understood each word, said each word, and an adjusted understood measure derived from subtracting the number of children who said a word from the number of children who understood the word. Pearson and Spearman correlations are given in the upper and lower quadrants of Table 5, respectively. These two types of correlations gave very similar results. As already mentioned, words said significantly correlated with words understood and with the adjusted measure of words understood. Log parental input frequency did not correlate with the number of children who said or understood the words. As can also be seen in the table, higher parental input frequencies were negatively correlated with difficulty of articulation. Most importantly, for distinguishing among the three theories, difficulty of articulation was negatively correlated with words said but not with the adjusted measure of words understood. This is true even though words said remains positively correlated with our adjusted measure of words understood. This result provides some evidence that different processes might be involved in production than in perception so that difficulty of articulation is more of an influence in production than in perception.

To better assess the contribution of difficulty of articulation independently of log parental input frequency, we carried out partial correlations shown in Table 6. Difficulty of articulation remains correlated with words said when parental input frequency was partialled out. In contrast, difficulty of articulation did not correlate with the adjusted measure of words understood when parental input frequency was partialled out.

Table 5. Correlations among difficulty of articulation, log parental input frequency, words said, words understood, and the adjusted words understood for the complete data set of 386 words. Pearson correlations are given in the upper right quadrant and Spearman correlations are given in the lower left quadrant. (Note * = $p < .05$; ** = $p < .001$).

	Difficulty of Articulation	Log Parental Input Frequency	Said Words	Understood Words	Adjusted Understood Words
Difficulty of Articulation		-.293**	-.245**	-.142*	-.045
Log Parental Input Frequency	-.283**		-.032	.018	.044
Said Words	-.208**	-.039		.748**	.412**
Understood Words	-.123*	.047	.757**		.913**
Adjusted Understood Words	-.056	.083	.576**	.952**	

Table 6. Partial correlations of difficulty of articulation with words said, words understood, and the adjusted words understood when log parental input frequency is partialled out. (Note * = $p < .05$; ** = $p < .001$).

	Difficulty of Articulation (Log Parental Input Frequency Partialled Out)
Said	-.271**
Understood	-.142*
Adjusted Understood	.031

To assess the validity of our Difficulty of Articulation measure, we compared it to a measure given by Shriberg (1993) who categorized 24 speech



segments into early, middle, and late acquisition classes with 8 segments per class.

Early 8: /m, b, j, n, w, d, p, h/

Middle 8: /t, ŋ, k, g, f, v, tʃ, dʒ /

Late 8: /ʃ, ʒ, s, z, θ, ð, l, r/

These measures have found some validity in research and applications (Goldstein and Fabiano, 2010; Shriberg et al., 1997). Across the 386 words, the Shriberg measure given in Table 4 correlated .91 with our Difficulty of Articulation measure and gave identical statistical outcomes.

Table 7. Partial correlations of Log Parental Input Frequency with words said, words understood, and the adjusted words understood when Difficulty of Articulation is partialled out. (Note * = $p < .05$; ** = $p < .001$).

	Log Parental Input Frequency (Difficulty of Articulation Partialed Out)
Said	-.113*
Understood	-.025
Adjusted Understood	.033

Table 7 gives the partial correlations between log parental input frequency and our three dependent measures. With the exception of the log parental input frequency and said words, log parental input frequency had very little influence on words understood or our difference measure when difficulty of

Table 8. Correlations among difficulty of articulation, log parental input frequency, words said, words understood, and the adjusted words understood for the set of 202 nouns. Pearson correlations are given in the upper right quadrant and Spearman correlations are given in the lower left quadrant. (Note * = $p < .05$; ** = $p < .001$).

	Difficulty of Articulation	Log Parental Input Frequency	Said Words	Understood Words	Adjusted Understood Words
Difficulty of Articulation		-.255**	-.293**	-.161*	-.044
Log Parental Input Frequency	-.221*		.313**	.263**	.178*
Said Words	-.280**	.311**		.810**	.525**

Understood Words	-.136*	.251**	.783**	.924**
Adjusted Understood Words	-.053	.207*	.640**	.971**

articulation is partialled out. This result should not be too surprising since Goodman et al. (2008) found no overall effect of log parental input frequency in a similar data set. However, they found some significant effects of log parental input frequency when they carried out the analyses within some of the data set's different lexical categories. We therefore carried out correlation analyses on subsets of our data sample. First, we created a noun category that had 202 nouns from our data sample. Table 8 gives the Pearson correlations among difficulty of articulation, parental input frequency, number of children that understood each word, said each word, and the adjusted understood measure for the subset of 202 nouns. The results shown in Table 8 replicate the differential influence on said and understood words. Difficulty of articulation correlated significantly with said words ($-.293$, $p < .001$) but not with the adjusted measure of words understood. Replicating Goodman et al.'s findings, log parental input frequency correlated significantly with our three dependent measures.

Table 9 gives the partial correlations between Difficulty of Articulation and our three dependent measures for the 202 nouns. The influence of Difficulty of Articulation holds up when the analysis is limited to just the noun category. The partial correlation of difficulty of articulation with said words was significant even when the contribution of parental input frequency was partialled out. Similarly, difficulty of articulation did not correlate with understood words or the adjusted measure of understood words. This significant correlation of difficulty of articulation with said words but not with understood words casts doubt on both the motor theory that posits similar processes in perceiving and producing words and the theory assuming equivalent representations for perception and production.

Table 9. Partial correlations of difficulty of articulation with words said, words understood, and the adjusted words understood when log parental input frequency are partialled out for 202 nouns. (Note * = $p < .05$; ** = $p < .001$).

	Difficulty of Articulation (Log Parental Input Frequency Partialled Out)
Said	-.237**
Understood	-.104
Adjusted Understood	.000



Table 10 gives the partial correlations between parental input frequency and our three dependent measures for the class of 202 nouns. When the analysis is restricted to the class of nouns, log parental input frequency has the expected effect that children are more likely to have vocabulary corresponding to what they hear.

Table 10. Partial correlations of Log Parental Input Frequency with words said, words understood, and the adjusted words understood when Difficulty of Articulation is partialled out for the 202 nouns category. (Note * = $p < .05$; ** = $p < .001$).

	Log Parental Input Frequency (Difficulty of Articulation Partialed Out)
Said	.257**
Understood	.232**
Adjusted Understood	.172*

We repeated the partial correlations on the 51 action words, the 32 descriptive words, and the 42 prepositions, quantifiers, question, and words about words. The same results were found as for the subset of 202 nouns.

3.2.3 Correlations with Imagery and Concreteness

Ma et al. (2009) found that the acquisition of both nouns and verbs in Chinese and English was correlated with the imageability ratings of adults. McDonough et al. (2011) replicated these results in English and found that both parental input frequency and form class correlated with vocabulary acquisition, as measured by the CDI. The imagery ratings were taken from Masterson and Druks (1998). Given the possible differential influence of imagery and concreteness on comprehension and production, we found 133 nouns that had concreteness and imagery ratings in the MRC psycholinguistic database and replicated our analysis with these items. We found very little influence of concreteness and imagery. As expected, concreteness and imagery correlated with one another, $r = .569$, $p < .001$. However, the only significant correlation was that higher concreteness was positively correlated with words said, $r = .194$, $p < .05$.

There is a very small range of concreteness and imageability within this class of nouns, which probably accounts for the failure to find robust effects of concreteness and imagery within the class of nouns. To expand the range of concreteness and imageability ratings in our analysis, we found 269 words in our complete database of 386 words that had concreteness and imagery ratings in the MRC psycholinguistic database and replicated our analysis with these items. These correlations are shown in Table 11. As can be seen in the table, Imagery was positively correlated with Difficulty of Articulation

and inversely correlated with Log Parental Input frequency. Imagery was positively correlated with said but not with understood words.

To assess whether imagery, log parental input frequency, and Difficulty of Articulation independently contributed to vocabulary development, partial correlations were carried out. Table 12 gives the partial correlations and shows that all three variables made significant independent contributions to vocabulary development.

Table 11. Correlations among difficulty of articulation, parental input frequency, words said, words understood, and the adjusted words understood for the subset of 269 words with imagery concrete ratings. Pearson correlations are given in the upper right quadrant and Spearman correlations are given in the lower left quadrant. (Note * = $p < .05$; ** = $p < .001$).

	Difficulty of Articulation	Log Parental Input Frequency	Imagery	Said Words	Understood Words	Adjusted Understood Words
Difficulty of Articulation		-.293**	.206**	-.245**	-.142*	Words
Log Parental Input Frequency	-.283**		-.642*	-.032	.018	-.045
Imagery	.188*	-.443**		.193**	.130*	.044
Said Words	-.208**	-.039	.305**		.748**	.081
Understood Words	-.123*	.047	.092	.757**		.412**
Adjusted Understood Words	-.056	.083	.036	.576**	.952**	.913**

Table 12. Partial correlations of Log Parental Input Frequency with words said, words understood, and the adjusted words understood when Difficulty of Articulation is partialled out for the 202 nouns category. (Note * = $p < .05$; ** = $p < .001$).

	Imagery (Difficulty of Articulation and Log Parental Input Frequency Partialled Out)	Difficulty of Articulation (Log Parental Input Frequency and Imagery Partialled Out)	Log Parental Input Frequency (Difficulty of Articulation and Imagery Partialled Out)
Said	.367**	-.253**	.232**
Understood	.269**	-.143*	.182*



Adjusted	.183*	-.072*	.133*
Understood			

To assess whether imagery, log parental input frequency, and Difficulty of Articulation independently contributed to vocabulary development, partial correlations were carried out. Table 12 gives the partial correlations and shows that all three variables made significant independent contributions to vocabulary development.

A multiple regression analysis was carried out and revealed that all three variables accounted for about 18%, 9% and 4% of the variance of said, understood, and adjusted understood words, respectively.

3.2.4 Correlations with Word Length, Segment Probability, and Neighbors

We also investigated the influence of several other possible influences: namely word length, segment probability, and number of neighbors. The number of phonemes in each word was used as a measure of word length. The positional segment frequency was computed for each sound in the target word by iterating over every entry in the corpus that is long enough to have any sound in the corresponding position (counted from the left edge of the word without respect to syllable structure) and by checking for matches against each sound in the target word. A position-sensitive sum and a biphone sum were computed for each words based on a child and adult corpus, respectively (Storkel, 2013, 2014). This measure is defined as the sum of the log frequencies of all of the words in the corpus that contain the given segment in a specific word position, divided by the sum of the log frequencies of all of the words in the corpus that contain any segment in that word position (see, Storkel, 2004b). Basically, these latter two measures give the uniqueness of a segment occurrence in a word relative to its occurrence in all of the words in the corpus. Finally, the number of neighbors was defined in the traditional manner in which neighbors were counted as any words that differed from the target word by just one phoneme (Storkel, , 2013, 2014; Storkel & Hoover, 2010).

To assess these variables, we carried out an analysis on 198 nouns that had these measures. The only significant correlation was that length was inversely correlated with said words.

We also computed these measures for 346 words from our full database from the measures given by Vaden (2009). These measures are identically defined as just described except that the segment probabilities are computed on the basis of frequency of occurrence in the Kucera-Francis database. These measures gave very little predictive value for the acquisition of either receptive or expressive language.

4. Conclusions and Discussion

The correlation analyses with all words overall and the lexical classes nouns and action words show a somewhat more robust correlation of difficulty of articulation with words said than with words understood. This result is

particularly impressive because the range of understood words was larger than that of said words so that *ceteris paribus* a larger correlation with understood words would be expected. This result points to different processes involved in speech perception and production and their acquisition. Future work should be extended to new and more elaborate databases.

4.1 Previous Literature

We might expect that input frequency might positively correlate with vocabulary acquisition. Goodman, Dale, and Li (2008) examined the influence of the frequency of occurrence on vocabulary acquisition as revealed by perception/comprehension and production. They used two valuable databases: The MacArthur-Bates Communicative Development Inventory which provides norming data on age of acquisition for 562 individual words collected from the parents of children aged 0;8 to 2;6, and the CHILDES database which provides estimates of frequency with which parents use these words with their children (age: 0;7–7;5; mean age: 36 months). Age of usage was defined as the age when 50% of the children in the MacArthur-Bates database used the word in comprehension or production, respectively. They computed Pearson correlations between the age of usage of words (for comprehension and for production) and the frequency of occurrence of those words in parents' child directed speech.

There were no overall effects of frequency so Goodman et al. (2008) carried out the analyses within different lexical categories: common nouns, people words, verbs, adjectives, closed class, and a default category. This partitioning is not commonly carried out in the adult literature because most experiments studying frequency effects use just a single class of items such as nouns. The expected correlations of frequency were now observed within some lexical categories but they were much larger for production than for comprehension. Parental input frequency was significantly correlated with all six categories of words for production but with only the category nouns for perception.

Goodman et al. (2012) explained these differences between comprehension and production in terms of the expected smaller number of experienced occurrences required to learn to comprehend a word versus the number required to produce that word. If learning is probabilistic in both cases, a variable with the smaller number of occurrences will necessarily give a smaller correlation than a variable with a larger number of occurrences. If production lags comprehension, it follows that a larger number of occurrences is needed for production than for comprehension and, therefore, more extreme correlations should be expected for production than comprehension.

An analogous explanation might explain the differences of the influence of difficulty of articulation on receptive and expressive language. If additional occurrences are required to learn to produce a word than to understand a word, then the influence of difficulty of articulation would be magnified for production relative to perception. If this explanation is valid, it adds further support for the idea that comprehending a word occurs much earlier in language acquisition than producing that word, as observed in the



MacArthur-Bates database and by Roy's (2011) observations. In addition, it adds to the broad range of results that indicate that the same representation is unlikely to be used in both receptive and expressive language, as assumed by the motor theory and the McMurray et al. (2012) formal common-representation model.

One possibility, not considered by Goodman et al. (2008), however, was that there were fewer items in the correlation analysis for the comprehension condition than for the production condition. The comprehension measures were not available for the age group 1;6 to 2;6 (the scoring for these items asked whether the word was both understood and produced). Our current study overcomes this potential confounding because we have equivalent data sets for perception and production.

Some children are extremely late talkers but seem to show normal comprehension of speech. This result would not be expected if indeed speech production processes necessarily mediated speech comprehension. Sowell (2001) gives a number of examples of late talking children. These children begin talking late in their development after ages two or three but seem to have no delay in cognitive development or even the comprehension of language. Leslie, a young girl had an IQ of 139 but she did not talk until age 2. Leslie was like many other children who talk late but have no trouble understanding what other people are saying. While Leslie repeatedly had difficulty in producing words to express her meanings and relied heavily on multi-purpose words, her passive vocabulary was in the 99th percentile, as evidenced by her score on the Peabody picture vocabulary test. Sowell concludes, "Yet another fact consistent with this hypothesis is that many bright children who have not yet begun talking have no difficulty understanding what other people are saying to them and may even follow complex instructions better than most other children their age" (Sowell, 2001. p. 95). Temple Grandin did not talk until she was four.

There is some recent evidence, however, that late talkers might be somewhat delayed in comprehension also. Weismer, Venker, Evans, and Moyleb (2013) carried out a fast mapping vocabulary learning task with late-talking (LT) toddlers and toddlers with normal language (NL) on both novel object labels and familiar words. The LT group revealed poorer learning than the NL group on both the comprehension and production of novel words and the production of familiar words. The authors suggested "that late talkers' limitations in expressive language do not just stem from lexical retrieval problems, but appear to be related to more fragile phonological, lexical and semantic representations that are reflected in subtle comprehension difficulties that, in turn, result in more substantive deficits in vocabulary production." (p. 10).

Goodman et al (2012) uncovered a potentially damaging result for motor theory. They found that closed-class words are not produced at age 2;6 even though they occur frequently in the child's parental input speech. (The closed-class words consisted of pronouns, words about time, articles, quantifiers, prepositions, and question words.) There is evidence that children by this same age do comprehend many of these words, which is additional evidence against motor theory. Toddlers (13-15 month-olds) have

only a few words in their productive vocabulary but can understand the relations among words in a spoken sentence (Hirsh-Pasek & Golinkoff, 1996). One alternative explanation that we do not accept is that children might be using telegraphic speech and intentionally omitting these words in their production even though they are perfectly capable of doing so.

The dissociation between expressive and receptive language is not limited to an advantage of receptive language. There are cases in which a correct use of expressive language can precede receptive language. As reviewed by Hendriks (2014), preschool English-speaking children correctly produce pronouns such as *me* or *him* and a reflexive such as *myself* or *himself*. On the other hand, these same children and even children a few years older misinterpret these pronouns in sentence contexts. They will often interpret an object pronoun as a reflexive. So, for example, they understand “Ernie washed him” as “Ernie washed himself.” Hendriks (2014) describes several other aspects of grammar in which expressive language precedes receptive language.

4.2 Theories Revisited

As recently observed by Hickok (2009), motor theories have a long history in behavioral science. Over 100 years ago, Walter B. Pillsbury (1911) said, “The [motor] theory is so simple and so easy to present that every one is glad to believe it. The only question that anyone cares to raise is how much of it will the known facts permit one to accept” (p. 84). It is sobering how little has changed during the intervening century in that few, in any, motor theorists would be swayed by Pillsbury’s caveat.

What is most disappointing about motor theory is that its advocates do not offer specific testable hypotheses about how it accomplishes speech perception. Most of the explanations have taken the form of analysis by synthesis models. Given a speech event, the perceiver benefits from selecting and synthesizing potential alternatives during their processing of the speech input and somehow determining which of these possibilities is the best match. To accomplish this selection, however, the perceiver still requires access to the speech signal so we have regressed back to the original challenge of understanding how the speech signal is processed.

Brain traumas have helped scientists uncover specific areas of the brain responsible for comprehension and production of sign language (Hickok, 2009). Traditionally, brain researchers and practitioners have distinguished between Broca’s area (discovered by Paul Broca in 1861) responsible for speech production and Wernicke’s area (discovered by Carl Wernicke in 1874) responsible for speech comprehension. Both of these areas are in the left hemisphere and are close neighbors of the auditory cortex. Damage to either of these areas appears to compromise only one of these two abilities and leave the other ability mostly intact. (Damage to the right hemisphere leaves speech perception and comprehension more or less intact but has grave consequences for visual-spatial behavior.) For example, persons with Broca’s aphasia can perceive and understand speech but have difficulty producing it.

Dodd (1975) provides interesting although somewhat indirect evidence that perception and production are dependent on separate processes. She



recorded three-year-old children when they named a large set of pictures. They then listened to this recording of them saying the names as well as adults saying the names. They were asked to identify what was being named for each word they heard. They could accurately recognize only 48% of the words they uttered in contrast to their 94% accuracy in recognizing the adults' words for these same objects. If production mediated perception, as assumed by motor theory, then the children should have recognized their own words at least as well as the adult words.

Successful perception requires a reasonable invariant relationship between the things being recognized and their corresponding categories. To date, there appears to be no invariant correspondence between a phoneme and its acoustic occurrence in the speech signal. Motor theory was advocated in part because it appeared to solve the lack of invariance in the speech signal. However, this finding alone does not justify a motor theory and there is now some evidence that there may be more invariance in the acoustic consequences of articulation than in the articulation itself. Vocal tract imaging and tracking techniques indicate that American speakers produce /r/ with many different tongue shapes, and yet all of these are perceived as /r/ (Nieto-Castanon, Guenther, Perkell, & Curtin, 2005). What appears to be critical for a /r/ to be perceived as such is that its acoustic stimulus must have a very low third formant (F3), although we know that many other characteristics of the spectrum are influential (such as the direction of the formant transitions). There is other evidence that talkers use what they hear to guide their speech production. If the auditory feedback given a talker is modified, then the talker will actually modify their articulation based on what they heard rather than their articulatory movements (Guenther, Hampson, & Johnson, 1998).

Mefford and Green (2010) assessed the degree to which articulatory and acoustic information specified a phonetic category and the variability within the category. They had typical talkers utter a sentence in typical speech or slower or louder than normal. These different utterance conditions naturally create articulatory and acoustic variability. For our purposes concerning sensory versus motor sources of information in speech perception, a question is whether the acoustic or the articulatory source is more informative. A source is defined as more informative to the extent it distinguishes among different phonemes and has low variability within a phoneme class. Based on their measurements of tongue movements and first and second formant transitions, the authors concluded that the acoustic input was more informative than the articulatory information.

Vihman (2002) describes how motor processes might work in speech acquisition in one of the few explicit accounts of the putative role of motor processes in speech perception. Infants practice canonical babbling and produce consonant-vowel (CV) sequences at 6-8 months of age. This practice in production sensitizes them to similar speech that is received from their caregivers. Having uttered CV sequences allows these CVs to be more easily recognized from their caregiver's speech because their familiarity from babbling allows them to pop out of the acoustic stream. Vihman and her colleagues have demonstrated an effect of infants' familiarity with their own

production patterns on attention to isolated-word lists (DePaolis, Vihman, & Nakai, 2013; Majorano, Vihman, & DePaolis, 2013) and to words embedded in passages (DePaolis, Vihman, & Keren-Portnoy, 2011). They interpret these results that their vocal practice “bring into particular focus elements of what is perceived (Thelen & Smith, 1994; Vihman, 1991, 1996)”. In another study, they showed that accuracy in speech production was correlated with the quality of speech perception. They tracked the recordings of 59 infants weekly from 9 months, to identify the age at which they had two well-practiced consonants in their speech. The authors then recorded looking times to words likely to be known or unknown to the infants. They found that those infants with two well-practiced consonants in their speech revealed a sharper differentiation between the known and unknown words. The results were interpreted as an interaction between productive and receptive knowledge in development.

However, the impressive series of results by Vihman and her colleagues do not necessarily mean that production processes are necessarily involved in speech perception. Although the babbling patterns would become familiar with practice and this increase in familiarity might facilitate perception, it does not mean that the motor processes involved in babbling were actually functional during the infant’s speech perception. We can expect that both perception and production would be modified as the child experiences more language, even though they operate independently of one another (e.g., Westerman & Miranda, 2004). It is also possible that perception of a segment or word is mastered first and that improvement in production comes later. Production might allow an infant to better learn a speech category without necessarily requiring that production processes be involved in the infant’s perception. Furthermore, infants at 6-8 months quickly learn arbitrary statistical properties of segments occurring in continuous speech, which cannot be easily explained by canonical babbling (Saffran, 2003).

If perception precedes production in a child’s language development, then the possibility of perception-based phonological development might be worth studying (Werker et al, 2012). As described earlier, one of the basic assumptions of the FLMP is that there are multiple sources of information in the speech signal and these cues are not equally informative, they are learned at different stages, and individuals differ greatly in terms of when and the degree to which they are learned. This framework is similar to Werker and Curtin’s (2005) PRIMIR model (a developmental framework for Processing Rich Information from Multi-dimensional Interactive Representations). Curtin, Byers-Heinlein, & Werker (2010) have extended their model to give an account of bilingual language development. As observed by Soderstrom et al. (2009), “infants pursue multiple analyses of many input properties, letting a thousand flowers bloom.”

We have made the theoretical and empirical arguments for an independence of perception and production processes. The significant correlation of difficulty of articulation with said words but not with understood words casts doubt on the motor theory that posits similar processes in perceiving and producing words and a theory assuming equivalent representations for perception and production. Our new results run counter to the motor theory of speech perception and to the assumption that there is a common



representation used in both perceiving and producing speech. This same evidence argues in favor of separate processes for speech comprehension and speech production.

Given that empirical and theoretical results have weakened the claims of motor theory, there have been several softer versions of the theory.

4.2.1 Modified Motor Theory: Ambiguity Backup.

In this version of motor theory, the recognition system recruits the motor representations only when there is an ambiguity in the input stream such that speech perception does not succeed based on typical sensory processing (e.g., Moulin-Frier & Arbib, 2013). In some cases, there may be a degraded speech signal so that no single alternative can be chosen from a number of similar candidate words. It should be kept in mind, however, that for typical cases when the sensory representations are sufficient for pattern recognition the motor representations are never consulted. When ambiguity persists after prototypical pattern recognition processes have been engaged, however, it might seem reasonable to consult an auxiliary source of information such as a motor representation.

More recently, Moulin-Frier & Arbib (2013) acknowledge that there is good evidence that much of speech perception can occur without any intervention of motor processes but speculate that, in some cases, motor processes can make a positive contribution. One example would be speech perception of challenging input such as the speech of a non-native speaker of the language. In this case, production processes might be instantiated that would provide a set of alternative candidates for understanding, as in traditional analysis-by-synthesis models.

Although the ambiguity backup model is attractive primarily because it allows for motor processes but does not make them mandatory, it will not be easy to test. Utilizing motor processes when the direct route to perception fails might necessarily imply that perception would be slower when the motor route is engaged than when it is not. Thus, the overall distribution of reaction times (RTs) to a successful perceptual identification should be a composite of two underlying distributions rather than just one. If indeed perception followed this dual route model, it is surprising that no one has observed a bivariate distribution of RTs in spoken word recognition.

Like the original motor theory, it is still troublesome to understand how a motor representation might help. The fallback explanation is usually the instantiation of an analysis by synthesis process. The motor system generates likely alternatives that somehow allow a comparison process that might aid in perception of the word. We believe our tests that falsify the standard motor theory are equally relevant to the Ambiguity Backup version.

4.2.2 Modified Motor Theory: Limited Influence.

We know that there are multiple influences in speech perception and it is possible that motor processes might have a small influence. Hickok, et al. (2011) are swayed by recent findings that suggest a limited and modulatory role for the motor system in speech perception. This evidence comes from neurological measurements that show activation of the motor system during

speech perception with no explicit motor task. These results, of course, could simply reflect associated activation with no functional role in speech perception.

Other evidence indicates some success with interfering with speech perception by stimulating premotor cortex, and facilitating or interfering with speech perception by stimulating motor or lip or tongue areas. It is important to note, however, that motor theory necessarily predicts that the sensory information must be more accurately resolved independently of any contribution of additional sources of information that prediction provides. It's processing cannot be simply biased in the direction of one alternative or another based on an additional source of information. None of these studies that show some role for the motor system has partialled out these two types of influence. Hickok et al. (2011) admit that all of these findings might simply reflect a non-perceptual bias rather than perceptual sensitivity (Massaro & Cowan, 1993).

Roy (2011) documents about 70 instances in which his son attempted to pronounce "water" before he was able to pronounce it correctly. Many of these instances illustrate that he was able to perceive and understand the spoken word but simply unable to produce it. In addition, Roy's child successive attempts at producing the word water did not appear to be a continuous embellishment of a single impoverished representation but rather included significantly different forms. This scientific observation revives the anecdotal one in which a father is mimicking his son's mispronunciation of the word "rabbit". His son says, "No dad, not wabbit, but wabbit." The son clearly could perceive the difference between "rabbit" and "wabbit" even though he wasn't able to accurately produce the difference.

4.2.3 Common Representation for Word Perception and Production.

This thesis of equivalent representations (Figure 2) might explain the comprehension advantage by assuming that an infant's word's representation is sketchy. This adumbrated representation is sometimes sufficient for comprehension but not for speech production. The difference between comprehension and production is simply the number of viable competitors in the task. The results from individual children, on the other hand, show that a child's good perception performance of a variety of words but little or no production contradicts the common representation explanation (see Tables 2 and 3 and corresponding explanation).

4.2.4 Pattern Recognition: Different Processing for Perception and Production

Although the three theories we have considered are still very much in contention, we favor the view that understanding speech is a prototypical pattern recognition situation. Multiple sources of information are used to make sense of the input conveying some meaningful event. This framework has been successful in describing many different speech perception experiments with both adults and children (Massaro, 1975; 1998; Movellan & McClelland, 2001). In addition to this behavioral model, the dual-route model with ventral and dorsal processing streams provides a neural level of



description (Hickok, 2008). We look forward to new results and theories to illuminate further how language is learned and used.

Acknowledgement

The authors would like to thank Larry Fenson, Virginia Marchman, and the advisory board of the MacArthur-Bates Communicative Development Inventories (CDI) for making the data from individual children available for the current analyses. We would also like to thank Philip Dale for providing the parental input frequencies and Tara McAllister Byun, Judith Goodman, Cecilia Kirk, Ping Li, Brian MacWhinney, and Susan Rvachew for helpful correspondence on this research. The first author dedicates this article to the memory of Elizabeth Bates who was a primary force in revolutionizing how we think about language acquisition.

References

- American CDI (2014).
<http://www.cdi-clex.org/vocabulary/singlewordlist/search/corpora/1>
- Bates, E. (1993). Comprehension And Production In Early Language Development: Comments On Savage-Rumbaugh et al. *Monographs of the Society for Research in Child Development*. Serial No. 233, Vol. 58, Nos. 3-4, 1993, pp. 222-242
- Bates, E., Bretherton, I. & Snyder, L. (1988). *From first words to grammar: individual differences and dissociable mechanisms*. New York, C.U.P.
- Bates, E. & Goodman, J. C. (1999). On the emergence of grammar from the lexicon. In B. MacWhinney (Ed.), *The emergence of language*. Mahwah, NJ: Lawrence Erlbaum.
- Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36, 181-253. doi:10.1017/S0140525X12000477
- Curtin, S., Krista Byers-Heinlein, K., & Werker, J. F. (2010). Bilingual beginnings as a lens for theory development: PRIMIR in focus. *Journal of Phonetics*. doi:10.1016/j.wocn.2010.12.002
- Dale, P. S., Bates, E., Reznick, J. S. & Morisset, C. (1989). The validity of a parent report instrument of child language at twenty months. *Journal of Child Language*, 16, 239-51
- Dale, P. S., and Fenson, L. (1996). Lexical development norms for young children. *Behavioral Research Methods, Instruments, & Computers*, 28, 125-127.
<http://www.cdi-clex.org/vocabulary/singlewordlist/search/corpora/1>
- DePaolis, R., Vihman, M. M. & Keren-Portnoy, T. (2011). Do production patterns influence the processing of speech in prelinguistic infants? *Infant Behavior and Development*, 34, 590-601.
- DePaolis, R., Vihman, M. M. & Nakai, S. (2013). The influence of babbling patterns on the processing of speech. *Infant Behavior and Development*, 36, 642-649 .
- Dodd, B. (1975). Children's understanding of their won phonological forms. *Quarterly Journal of Experimental Psychology*, 27, 165-172.

- Fenson, L., Fenson, L., Dale, P., Reznick, J., Bates, E., Thal, D., & Pethick, S. (1994). Variability in early communicative development. *Monographs of the Society for Research in Child Development*, Serial No. 242, Vol. 59, No. 5.
- Fenson, L., Bates, E., Dale, P., Goodman, J., Reznick, J. S., & Thal, D. (2000). Measuring variability in early child language: Don't shoot the messenger. Comment on Feldman et al. *Child Development*, 71(2), 323-328.
- Fenson, L., Dale, P. S., Reznick, J.S., Thal, D., Bates, E., Hartung, J. P., Pethick, S., & Reilly, J. S. (1993). *The MacArthur Communicative Development Inventories: User's Guide and Technical Manual*. Baltimore : Paul H. Brookes Publishing Co.
- Fenson, L., Marchman, V. A., Thal, D. J., Dale, P. S., Reznick, J. S. & Bates, E. (2007). *MacArthur-Bates Communicative Development Inventories : user's guide and technical manual*, 2nd ed. Baltimore : Paul H. Brookes.
- Friston K. (2010) The free-energy principle: a unified brain theory? *Nature Reviews Neurosciences* 11(2):127-138.
- Galantucci, B., Fowler , C. A., & Turvey, M. T. (2006). The motor theory of speech perception reviewed. *Psychonomic Bulletin & Review*,13, 361-377.
- Glenberg A. (2008). *Toward the Integration of Bodily States, Grounding: Social, Cognitive and Neuroscientific Approaches*. ch. 2, pp. 43-70, Cambridge University Press.
- Glenberg, A. M., & Gallese, V. (2012). Action-based Language: A theory of language acquisition, comprehension, and production. *Cortex*, 48, 905-922. doi:10.1016/j.cortex.2011.04.010.
- Goldstein, Brian and Fabiano, Leah (2007). Assessment and intervention for bilingual children with phonological disorders, *ASHA Leader*. 12(2), 6-27, 26-27, 31. [http://www.asha.org/Publications/leader/2007/070213/f070213a/].
- Goodman, J. C., Dale, P. S., & Li, P. (2008). Does frequency count? Parental input and the acquisition of vocabulary. *Journal of Child Language*, 35, 515 – 531.
- Guenther, F. H. (1995). Speech sound acquisition, coarticulation, and rate effects in a neural network model of speech production. *Psychological Review*, 102, 594-621.
- Guenther, F. H., Hampson, M., & Johnson, D. (1998). A theoretical investigation of reference frames for the planning of speech movements. *Psychological Review*, 105, 611-633.
- Hahn, E. R., & Gershkoff-Stowe, L. (2010). Children and adults learn actions for objects more readily than labels. *Language Learning and Development*, 6, 1-26.
- Hendriks, P. (2014). *Asymmetries between production and comprehension*. New York: Springer.
- Hickok, G. (2009). Eight problems for the mirror neuron theory of action understanding in monkeys and humans. *Journal of Cognitive Neuroscience*, 21, 1229-1243.
- Hickok, G. (2012). Computational neuroanatomy of speech production. *Nature Reviews Neuroscience*, 13, 135-145.



- Hickok, G., Houde, J. & Rong, F. (2011). Sensorimotor integration in speech processing: computational basis and neural organization. *Neuron* 69, 407–422.
- Hickok, G. & Poeppel, D. (2007). The cortical organization of speech processing. *Nature Reviews Neuroscience*, 8, 393-402.
- Hickok, G., Holt, L. L. & Lotto, A. J. (2009). Response to Wilson: What does motor cortex contribute to speech perception? *Trends in Cognitive Sciences*, 13, 330-331.
- Hirsh-Pasek, K., & Golinkoff, R. M. (1996). *The origins of grammar: Evidence from early language comprehension*. Cambridge, MA: MIT Press.
- Huttenlocher, J. (1974). The origins of language comprehension. In R. Solso (Ed.), *Theories in cognitive psychology: The Loyola Symposium*. Oxford, England: Erlbaum.
- Kirk, C. (2008). Substitution errors in the production of word-initial and word-final consonant clusters. *Journal of Speech, Language, and Hearing Research*, 51, 1-14.
- Kirk, C., & Demuth, K. (2005). Asymmetries in the acquisition of word-initial and word-final consonant clusters. *Journal of Child Language*, 32(4). 709-734.
- Locke, J. L. (1980). The prediction of child speech errors: Implications for a theory of acquisition. In G. H. Yeni-Komshian, J. F. Kavanagh, & C. F. Ferguson (Eds.), *Child phonology: Vol. 1. Production* (pp. 193-209). New York: Academic Press.
- Locke, J. L. (1983). *Phonological acquisition and change*. New York: Academic Press.
- Ma, W., Golinkoff, R.M., Hirsh-Pasek, K., McDonough, C., & Tardif, T. (2009). Imageability predicts the age of acquisition of verbs in Chinese. *Journal of Child Language*, 36(2), 405-423. doi:10.1017/S0305000908009008
- MacArthur-Bates Communicative Development Inventories (CDI). (2012). <http://www.cdi-clex.org/vocabulary/firstwords/top10/language/2>
- MacWhinney, B. (2000). *The CHILDES project (3rd Editioned.)*. Mahwah, NJ: Lawrence Erlbaum.
- Majorano, M., Vihman, M. M. & DePaolis, R. A. (2013). The relationship between infants' production experience and their processing of speech. Majorano, M., Vihman, M. M. & DePaolis, R. A. (2013). *Language Learning and Development*, DOI: 10.1080/15475441.2013.829740.
- Marchman, V. A., Dale, P. S., Reznick, J. S., Thal, D., & Bates, E. (2007). *MacArthur-Bates Communicative Development Inventories: User's Guide and Technical Manual*. 2nd Ed.. Baltimore, MD: Brookes Publishing Company.
- Massaro, D.W. (1975). *Understanding Language: An Information Processing Analysis of Speech Perception, Reading and Psycholinguistics*. New York: Academic Press.
- Massaro, D. W. (1987). *Speech Perception by Ear and Eye: A Paradigm for Psychological Inquiry*. Hillsdale, NJ: Erlbaum Associates.

- Massaro, D. W. (1994). Psychological aspects of speech perception: Implications for research and theory. In M. Gemsbacher (Ed.), *Handbook of Psycholinguistics* (pp. 219–263). New York: Academic Press.
- Massaro, D. W. (1998). *Perceiving Talking Faces: From Speech Perception to a Behavioral Principle*. Cambridge, MA: MIT Press.
- Massaro, D. W. (2012). Speech Perception and Reading: Two Parallel Modes of Understanding Language and Implications for Acquiring Literacy Naturally. *American Journal of Psychology*, 125, No. 3, 307-320.
- Massaro, D.W. & Chen, T.H. (2008). The Motor Theory of Speech Perception Revisited. *Psychonomic Bulletin & Review* (pp. 453-457), Vol. 15(2).
- Massaro, D.W., & Cowan, N. (1993). Information Processing Models: Microscopes of the Mind. *Annual Review of Psychology*, 44, 383-425.
- McAllister Byun, T. (2012). Bidirectional perception-production relations in phonological development: evidence from positional neutralization. *Clinical Linguistics and Phonetics*, 26(5), 397-413.
- McDonough, C., Song, L., Hirsh-Pasek, K., Golinkoff, R. M., & Lannon, R. (2011). An image is worth a thousand words: Why nouns tend to dominate verbs in word learning. *Developmental Science*, 14, 181-189.
- McLeod, S., & Bleile, K. (2003). Neurological and developmental foundations of speech acquisition. American Speech-Language-Hearing Association Convention, Chicago, November 2003, Invited Seminar Presentation.
- CDI (2014). <http://www.cdi-clex.org/>
- McMurray, B., Horst, J., and Samuelson, L. (2012) Word learning emerges from the interaction of online referent selection and slow associative learning. *Psychological Review*, 119, 831-877.
- Mefferd, A. S., & Jordan R.; & Green, J. R. (2010) Articulatory-to-Acoustic Relations in Response to Speaking Rate and Loudness Manipulations. *Journal of Speech, Language, 1206 and Hearing Research*, 53, 1206–1219.
- Mills, D. L., Coffey-Corina, S., & Neville, H. J. (1993). Language acquisition and cerebral specialization in 20-month-old infants, *Journal of Cognitive Neuroscience* 5, 317–34.
- Mills, D. L., Coffey-Corina, S., & Neville, H. J. (1997). Language comprehension and cerebral specialization from 13 to 20 months. *Developmental Neuropsychology* 13, 397–445.
- Moulin-Frier, C., & Arbib, M. A. (2013). Recognizing speech in a novel accent: the motor theory of speech perception reframed. *Biological Cybernetics*, 107, Issue 4, Page 421-447 DOI:10.1007/s00422-013-0557-3
- Nieto-Castanon, A., Guenther, F.H., Perkell, J.S., and Curtin, H. (2005). A modeling investigation of articulatory variability and acoustic stability during American English /r/ production, *J. Acoust Soc Am.* 117, 3196-3212.
- ParentFreq (Retrieved 2014).
<http://childes.psy.cmu.edu/derived/parentfreq.cdc>
- Quine, W. V. O. (1960). *Word and object: An inquiry into the linguistic mechanisms of objective reference*. Cambridge, MA: MIT Press.
- Quine W. V. O. (1990/1992) *Pursuit of Truth* (Harvard Univ Press, Cambridge, MA).
- Pillsbury, W. B. (1911). *The Essentials of Psychology*. New York, Macmillan.



- Roy, D. (2011). http://www.ted.com/talks/deb_roy_the_birth_of_a_word.html
- Rvachew, S., Chiang, P., & Evans, N. (2007). Characteristics of speech errors produced by children with and without delayed phonological awareness skills. *Language, Speech, and Hearing Services in Schools*, 38, 60-71.
- Savage-Rumbaugh, S., & Lewin, R. (1994). *Kanzi: The Ape At The Brink Of The Human Mind*. New York: John Wiley & Sons, Inc.
- Savage-Rumbaugh, S., Stuart, G. Shanker, S. G., & Taylor, T. J. (1998). *Apes, Language and the Human Mind*. New York: Oxford University Press.
- Shriberg, L. (1993). Four new speech and voice-prosody measures for genetics research and other studies in developmental phonological disorders. *Journal of Speech, Language, and Hearing Research*, 36, 105-140.
- Shriberg, L. D., Austin, D., Lewis, B., McSweeney, J. L., & Wilson, D. L. (1997). The percentage of consonants correct (PCC) metric: Extensions and reliability data. *Journal of Speech, Language, and Hearing Research*, 40, 708-722.
- Smit, A. B. 1993. Phonologic error distributions in the Iowa-Nebraska Articulation Norms Project: Consonant singletons. *Journal of Speech and Hearing Research* 36, 533-547.
- Smit, A. B., Hand, L., Freilinger, J. J., Bernthal, J. E. & Bird, A. (1990). The Iowa Articulation Norms Project and its Nebraska replication. *Journal of Speech and Hearing Disorders* 55, 779-98.
- Soderstrom, M., Conwell, E., Feldman, N., & Morgan, J. (2009) The learner as statistician: three principles of computational success in language acquisition. *Developmental Science*, 12, 409-411.
- Sowell, T. (2001). *The Einstein Syndrome: Bright Children Who Talk Late*. New York: Basic Books.
- Stoel-Gammon, C. (1987). Phonological skills of 2-year-olds. *Language, Speech, and Hearing Services in Schools* 18, 323-329.
- Stoel-Gammon, C. & Buder, E. (1999). Vowel length, post-vocalic voicing and VOT in the speech of two-year olds. *Proceedings of the XIIIth International Conference of Phonetic Sciences* 3, 2485-2488.
- Storkel, H. L. (2004). Methods for minimizing the confounding effects of word length in the analysis of phonotactic probability and neighborhood density. *Journal of Speech, Language, and Hearing Research*, 47, 1454-1468.
- Storkel, H. L. (2013). An online calculator to compute phonotactic probability and neighborhood density on the basis of child corpora of spoken American English. *Behavior Research*, 45:1159-1167. DOI 10.3758/s13428-012-0309-7
- Storkel, H. L. (2014). http://www.bncdnet.ku.edu/cgi-bin/DEEC/post_ccc.vi
- Storkel, H. L. & Hoover, J. R. (2010). A corpus of consonant-vowel-consonant real words and nonwords: Comparison of phonotactic probability, neighborhood density, and consonant age of acquisition. *Behavior Research Methods*, 42 (2), 497-506.
- Tomasello, M. & Mervis, C. B. (1994). The instrument is great, but measuring comprehension is still a problem. *Monographs of the Society for Research in Child Development*, Serial no 242, vol. 59, no 5.

- Tremblay S., Shiller D. M., Ostry D. J. (2003). Somatosensory basis of speech production. *Nature*, 423, 866-869.
- Trout, J. D. (2001). The Biological Basis of Speech: What to Infer from Talking to the Animals. *Psychological Review*, 108, (3), 523-549.
- Vaden, K.I. (2009). Phonological processes in speech perception. Doctoral dissertation, University of California at Irvine, Irvine, CA. (Proquest document id: 1781083831, <http://proquest.umi.com>; ISBN: 9781109155495).
- Vaden, K.I. (2012). Retrieved January 4, 2015, from <http://www.iphod.com/>.
- Vihman, M. M. (2002). The role of mirror neurons in the ontogeny of speech. In M. Stamenov & V. Gallese (Eds.), *Mirror neurons and the evolution of brain and language* (pp. 305-314). Amsterdam: Benjamins.
- Weismer, S. E., Venker, C. E., Evans, J. L., & Moyle, M. J. (2013). Fast Mapping in Late-Talking Toddlers. *Applied Psycholinguistics*, 34(1), 69-89. doi:10.1017/S0142716411000610
- Werker, J. F., & Curtin, S. (2005). PRIMIR: A developmental model of speech processing. *Language Learning and Development*, 1(2), 197-234. doi:10.1207/s15473341lld0102_4
- Werker, J. F., Yeung, H. H., & Yoshida, K. A. (2012). How Do Infants Become Experts at Native-Speech Perception? *Current Directions in Psychological Science*, 21, 221-226. <http://cdp.sagepub.com/content/21/4/221>.

