

## Citation for published version

Esteve-Gibert, N. [Nuria] & Muñoz Lahoz, C. [Carme]. (2021). Preschoolers benefit from a clear sound-referent mapping to acquire nonnative phonology. *Applied Psycholinguistics*, 42(1), 77-100. doi: 10.1017/S0142716420000600

### DOI

<https://doi.org/10.1017/S0142716420000600>

### HANDLE

<http://hdl.handle.net/10609/150778>

### Document Version

This is the Accepted Manuscript version.

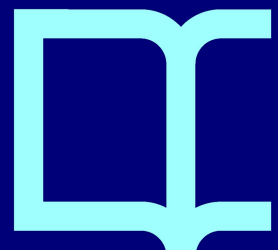
The version published on the UOC's O2 Repository may differ from the final published version.

### Copyright and Reuse

This manuscript version is made available under the terms of the Creative Commons Attribution Non Commercial No Derivatives license (CC-BY-NC-ND) <http://creativecommons.org/licenses/by-nc-nd/4.0/>, which allows others to download it and share it with others as long as they credit you, but they can't change it in any way or use them commercially.

### Enquiries

If you believe this document infringes copyright, please contact the UOC's O2 Repository administrators: [repositori@uoc.edu](mailto:repositori@uoc.edu)



## Pre-schoolers benefit from a clear sound-referent mapping to acquire non-native phonology

Journal:	<i>Applied Psycholinguistics</i>
Manuscript ID	APS-Mar-19-0046.R4
mstype:	Original Article
Specialty Area:	Child Second Language Acquisition, Phonetics and Phonology, Prosody, Speech Perception
Abstract:	<p>Previous studies have shown that visual information is a crucial input in early language learning. In the present study we examine what type of visual input helps pre-schoolers in acquiring non-native phonological contrasts. Catalan/Spanish-speaking children (4-5 years, N = 47) participated in a task to assess phonological discrimination abilities before and after a training. Three training conditions were presented: one with clear oral/visual speech information, one with an ostensive object-sound mapping, and one with a rich social interaction. Children's looking patterns were tracked to examine their focus of interest while being trained. Results revealed that pre-schoolers' discrimination abilities increase in all trained conditions, but the condition where the speaker created an ostensive object-sound mapping led to higher long-term gains (especially for younger children). Eye-tracking results further showed that children looked to the object of reference while being exposed to the novel phonological input, which may explain the higher learning gains in this condition. Our results indicate that preschoolers' learning of non-native phonological contrasts is particularly boosted when the speech input is accompanied by an object of reference that is signalled ostensively and contingently, compared to when the visual space only contains clear oral/visual speech information or social interactivity cues.</p>

SCHOLARONE™  
Manuscripts

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Pre-schoolers benefit from a clear sound-referent mapping to acquire non-native phonology**

Author 1 & Author 2

Affiliation 1, Affiliation 2

For Review Only

**Abstract**

Previous studies have shown that visual information is a crucial input in early language learning. In the present study we examine what type of visual input helps pre-schoolers in acquiring non-native phonological contrasts. Catalan/Spanish-speaking children (4-5 years,  $N = 47$ ) participated in a task to assess their phonological discrimination abilities before and after a training. Three training conditions were presented: one with clear oral/visual speech information, one with an ostensive object-sound mapping, and one with a rich social interaction. Children's looking patterns were tracked to examine their focus of interest while being trained. Results revealed that pre-schoolers' discrimination abilities increase in all trained conditions, but the condition where the speaker created an ostensive object-sound mapping led to higher long-term gains (especially for younger children). Eye-tracking results further showed that children looked to the object of reference while being exposed to the novel phonological input, which may explain the higher learning gains in this condition. Our results indicate that pre-schoolers' learning of non-native phonological contrasts is particularly boosted when the speech input is accompanied by an object of reference that is signalled ostensively and contingently in the visual space, compared to when the visual space only contains clear oral/visual speech information or social interactivity cues.

**Keywords:** contingency; visual speech; L2 acquisition; phonology; children

## Text

### INTRODUCTION

Acquiring a language is a complex process that requires a combination of the learner's linguistic, cognitive, and social abilities, together with exposure to quantitatively and qualitatively sufficient input. The present study will investigate which qualitative features of the visual input contribute the most to young children's learning of non-native phonological categories. Research on the development of non-native phonological categories in young children is still scarce (see discussions in Walley, 2008, and Erdener & Burnham, 2018). Most evidence comes from infant studies and shows that functionally-relevant phonological cues enhance learning. Young infants acquire novel phonological categories easier if they can consistently associate the acoustic distributional properties of the input to distinctive referential categories (eg. Best, 1993; Fennell & Waxman, 2010; Thiessen, 2011; Yeung & Nazzi, 2014; Yeung & Werker, 2009). The visual context accompanying a learning situation provides such associations, as infants acquire non-native phonemes when these are paired unambiguously with a visually-presented object of reference, compared to when no object of reference is present or when the pairing is inconsistent (Fennell & Waxman, 2010; Yeung & Werker, 2009). The functional dimension of phonological acquisition seems to continue in early childhood (Metsala & Walley, 1998) and adults (Feldman, Myers, White, Griffiths, Morgan, 2013). It is possibly grounded on the more general mechanism of 'acquired distinctiveness', which predicts that two phonetic stimuli are more easily distinguishable when they are consistently presented in distinctive contexts (Hall, 1991).

In natural conversations learners are exposed to additional visual cues that complement the statistical and functional dimension of the phonological acquisition process. A new word like 'flower' is not only perceived acoustically, but the learner most probably observes the speaker moving his/her lips while referring to a specific object in the space and trying to capture the learner's attention towards that object. The presence of a social partner that creates socially-engaging interactive situations (Bannard & Tomasello, 2012; Hakuno, Omori, Yamamoto, & Minagawa, 2017; Kuhl, Tsao & Yu, 2003; Kuhl, 2007; Linebarger & Vaala, 2010; Nussenbaum & Amso, 2015; Roseberry, Hirsh-Pasek, & Golinkoff, 2014), the exposure to ostensive signs of the object-label mapping (Csibra & Gergely, 2009; Hanna & Brennan, 2007; Moore, Angelopoulos, & Bennett, 1999; Triesch, Teuscher, Deák, & Carlson, 2006; Wu & Gros-Louis, 2014), and the learner's sensitivity to visual speech information (Birulés, Bosch, Brieke, Pons, & Lewkowicz, 2019; Erdener, 2007; Erdener & Burnham, 2013, 2018; Lalonde & Holt, 2015; Ter Schure, Junge, & Boersma, 2016; Weikum et al., 2007), are all factors that are found to positively impact the language acquisition process. Interestingly, speakers might not be able to spontaneously and simultaneously provide all these additional visual cues in the learning situation. For instance, an adult might create a socially-engaging situation in which a joint attentional frame is created (i.e. the adult uses eye contact or body movements to alternate his/her focus of interest between the child and the object of interest; Tomasello, 1995). But because of the alternation of the focus of interest, the adult might spend less time facing directly the child and therefore the child's exposure to visual speech input will be reduced. It is thus important to investigate the relative importance of the additional visual cues of a learning situation, to see if any of these may particularly boost the non-native phonological acquisition process.

We know that being part of socially-engaging interactive situations helps young learners to acquire non-native phonemes and novel lexical items (Bannard & Tomasello, 2012; Hakuno et al., 2017; Kuhl et al., 2003, Kuhl 2007; Linebarger & Vaala, 2010; Nussenbaum & Amso, 2015; Roseberry et al., 2014). Nine-month-old infants acquire novel (non-native) phonological contrasts when trained in live exposure situations (Kuhl, 2007), or when exposed to (media or real life) learning situations that are socially interactive, resembling real-life experiences

1  
2  
3 (Linebarger & Vaala, 2010). Bannard and Tomasello (2012) trained 2-year-old toddlers in two  
4 situations: one in which the interlocutor named the referent while alternating the gaze between  
5 the child and the object, and one in which the referent was available in the visual display, but  
6 the interlocutor did not look at it while naming it. While children showed implicit knowledge of  
7 the word-referent mapping in the two situations equally (by looking at the right referent in both  
8 situations), only if trained in the socially-engaging condition had children overly pointed at the  
9 right referent when asked to do so.

11 The exposure to ostensive signals that reinforce the object-label mapping also impacts the early  
12 acquisition of non-native phonemes. In first language acquisition, for instance, caregivers  
13 naturally speak to their infants using more redundancy and with exaggerated prosody (eg.  
14 Fernald, 1993; Fernald & Mazzie, 1991; Saint-Georges et al., 2013; Soderstrom, 2007). These  
15 social interactive features help infants' acquisition of phonological, syntactic, and lexical  
16 categories because they highlight these units in speech, boost the label-referent association,  
17 and promote infants' engagement in the communicative interaction (Golinkoff, Can,  
18 Soderstrom, & Hirsh-Pasek, 2015; Spinelli, Fasolo, & Mesman, 2017). At a visual level, adults  
19 naturally provide ostensive cues like eye gaze to indicate the relevant focus of attention, to  
20 disambiguate and reinforce the object-label associations, and to help young learners  
21 comprehend the meaning of what is being said (Hanna & Brennan, 2007; Wu,  
22 Tummeltshammer, Gliga, & Kirkham, 2014). Infants can follow their interlocutor's gaze very  
23 early in development (Brooks & Meltzoff, 2005), and eye-tracking results have shown that they  
24 use this ability to attend to the relevant object of reference (Senju & Csibra, 2008). Other eye-  
25 tracking studies have also shown that infants learn new words better when speech input is  
26 accompanied by ostensive signs indicating its functional value (Yoon, Johnson, Csibra, 2008; Wu,  
27 Gopnik, Richardson, & Kirkham, 2011; Wu & Kirkham, 2010), and the timing in which the adult  
28 establishes the object-label association is highly relevant. Previous studies in vocabulary learning  
29 find consistent evidence that in the ideal learning situation the adult provides the new linguistic  
30 input while the infant is already attending to the relevant referent, as opposed to trying to  
31 redirect the infant's focus of interest by providing a new linguistic input. This phenomenon has  
32 been called 'social contingency' or 'parental responsiveness' (e.g. Bannard & Tomasello, 2012;  
33 Hakuno et al., 2017; McGillion, Pine, Herbert, & Matthews, 2017; Nussenbaum & Amso, 2015;  
34 Roseberry et al., 2014, see Mermelshstine, 2017, for a review).

35 Phonetic information is not only perceived through the auditory modality. Listeners perceive  
36 speech sounds also through the visual channel, as the inspection of mouth movements while  
37 speaking provides redundant information to the acoustic signal (Gogate, Walker-Andrews, &  
38 Bahrack, 2001). Previous findings show that listeners discriminate and identify phonemes with  
39 more accuracy when visual (eg. Alm, Behne, Wang, Eg, 2009; Schwartz, Berthommier, &  
40 Savariaux, 2004) and haptic (Gick & Derrick, 2009) information is presented next to the acoustic  
41 information. Infants are sensitive to the visual aspects of speech from very early on (see Esteve-  
42 Gibert & Guellai, 2018, for a review), and they use the information from lip and head movements  
43 to acquire novel phonological contrasts. Learners perceive and identify novel consonants and  
44 vowels better after being exposed to audio-visual input, compared to audio-only input, evidence  
45 coming from young infants (eg.; Ter Schure et al., 2016; Weikum et al., 2007), school-aged  
46 children (eg. Erdener, 2007; Erdener & Burnham, 2013), and adults (eg. Aliaga-Garcia & Mora,  
47 2009; Cebrian & Carlet, 2012; Hardison, 2003; Hazan, Sennema, Faulkner, Ortega-Llebaria, Iba,  
48 & Chung, 2006; Ortega-Llebaria, Faulkner & Hazan, 2001). Interestingly, the very few studies  
49 exploring pre-schoolers seem to indicate that the benefit of visual input on top of the auditory  
50 cues may be less clear in this age range. Erdener (2007) tested 48 3- and 4-year-old English-  
51 speaking children in a non-native (Thai) phoneme discrimination task and found that children  
52 discriminated non-native phonemes better when presented in an audio-visual condition, and  
53 that their performance in the audio-visual condition was predicted by their ability in the  
54  
55  
56  
57  
58  
59  
60

auditory-only perception task (but not by their ability in the visual-only –lipreading– perception task).

In the current study we investigate if any of these additional visual cues (the presence of a socially interactive interlocutor, the interlocutor's production of ostensive signals of the object-label, and the exposure to visual speech information) has a stronger impact in young children's acquisition of non-native phonemes. Young children is an interesting population to investigate because preschool is for many children the time when they start being exposed to a new language, either because the language at school is different from the one spoken in their home environment, or because they start formal instruction in a second language (L2). And yet, as reviewed, most of previous research on the role visual input in phonological acquisition has been conducted with infants.

A word-learning task was chosen for the training because the presence of word referents enhances the learners' creation of non-native representations (e.g. Yeung, Chen, & Werker, 2014; Yeung & Werker, 2009). Young children's phonetic discrimination abilities are compared after being trained in one of these three distinct word-learning conditions: a *socially-engaging condition* in which the adult alternates his/her attention between a referent and his/her interlocutor, an *ostensive-cueing condition* in which the adult establishes a clear link between the referent and the linguistic input, and a *visual-speech condition* in which the adult faces the child without being socially engaging nor providing ostensive cues of the speech-meaning mapping. During the training children were exposed to minimal pairs of nonce words, which included English phonological contrasts that do not exist in their L1s and which were consistently paired with an object of reference. The third training session was conducted with an eye-tracking system to check for the children's eye-gaze patterns during the learning situation.

We predicted that if children only need an easy access to the visual speech input to acquire non-native phonological contrasts, their gains after being trained with the 'visual-speech' trials will be higher than with any other trained condition. Instead, if social engagement is more crucial, children's gains after being trained with the 'socially-engaging' condition will be higher than with the other conditions. Yet, if the ideal situation for the children's acquisition of non-native phonemes is one in which the L2 speaker provides the critical linguistic input while unambiguously referring to the meaning of the speech material, the 'ostensive-signalling' condition will lead to higher gains. In terms of the children's eye gaze preferences, if children use lipreading to learn novel phonological contrasts, we expect them to look more at the mouth in conditions that elicit the highest acquisition gains. Instead, if a higher social engagement is more relevant, we expect more gaze shifts from mouth to object in conditions that lead to higher gains. Finally, if children benefit from an ostensive signal of the speech-meaning association, we expect more gazes at the object in conditions that lead to higher gains.

## METHODS

### Participants

A total of 47 Catalan/Spanish-speaking children participated in the study (21 4-year-olds and 26 5-year-olds, 23 boys and 24 girls). Children were recruited from a school at a 1-hour radius of [removed for review]. The sole language of instruction at the school was Catalan. The children's parents signed a consent form and filled in a language background questionnaire. Parents reported that their child was Catalan dominant (N=22) or bilingual (N=25, 17 Catalan/Spanish and 8 Catalan/Spanish/other<sup>1</sup>). All children were included in the final sample because neither

---

<sup>1</sup> Other languages included French, Romanian, Galician, Portuguese, and Arabic. Although the phonemic inventory of French, Romanian, and Portuguese include /v/ (one of the novel phonemes to be learned,

Catalan nor Spanish include the English phonological contrasts to be learned, and because additional analyses revealed no effect of the other languages spoken on our results<sup>2</sup>.

### Stimuli

Three non-native phonological contrasts were studied: a consonant contrast (/b/ vs. /v/), a vowel contrast (/i:/ vs. /ɪ/), and a lexical stress contrast (trochaic word vs. iambic word). The training and test materials contained these critical phonemes in the form of minimal pairs. In the consonant contrast, the following 4 minimal pairs of non-words were used: *baggy-vaggy* (/bægi/-/vægi/), *boddy-voddy* (/ba:di/-/va:di/), *billy-villy* (/bili/-/vili/), *benny-venny* (/beni/-/veni/). In the vowel contrast, the following 4 minimal pairs of non-words were used: *teaggy-tiggy* (/tigi/-/tigi/), *deaddy-diddy* (/didi/-/didi/), *leanny-linny* (/lini/-/lini/), *seabby-sibby* (/sibi/-/sibi/). In the lexical stress contrast, cognate words were used that have a trochaic (Strong-Weak, SW) pattern in English but an iambic (Weak-Strong, WS) pattern in Catalan or Spanish. Catalan and Spanish are languages with lexical stress, and so both SW and WS patterns are already present in the young children's L1 vocabulary. The crucial difference between Catalan/Spanish and English is that many cognate words have a WS pattern in Catalan/Spanish but an SW pattern in English (e.g. 'acTOR' in Spanish but 'ACTor' in English, where capital letters indicate stress position), and learners need to learn to relocate the stress syllable in order to produce the contrastive metrical pattern. Thus, the following 4 minimal pairs were used for the stress contrast: *dolphin-dolphin*, *crocodile-crocodile*, *penguin-penguin*, *elephant-elephant*.

These three contrasts were chosen because adult learners discriminate and identify them better if the acoustic signal is accompanied by visual information (see Figure 1 for a display of the video frames corresponding to the points of maximal visual differentiation for each contrast). The two phonemes in the /b/-/v/ contrast have a different place of articulation (bilabial and labiodental, respectively) and thus their lip configuration varies. Previous studies on Catalan and Spanish learners of English suggest that visual information (lipreading) influences the learners' identification of these target consonants (Cebrian & Carlet, 2012; Hazan et al., 2006; but see Pons, Lewkowicz, Soto-Faraco, & Sebastian-Galles, 2009, for contradictory evidence in young infants and adults). The two phonemes /i:/-/ɪ/ also vary in terms of lip configuration (lips are wider spread in /i:/ than in /ɪ/, and there is a longer opening of the mouth in /i:/ than in /ɪ/) and Catalan/Spanish adult learners of English are found to rely on these visual differences for phoneme discrimination (Aliaga-Garcia, 2017; Flege, 1989; Ortega-Llebaria et al., 2001). For the lexical stress contrast, the visual cues are not related to the configuration of lips but to the movement of the head. When we speak we produce body gestures that are timely aligned with landmarks in speech. Head nods, one of these body gestures, are found to co-occur with prominent (i.e. stressed and/or pitch-accented) syllables in speech (e.g. Esteve-Gibert, Borràs-Comes, Asor, Swerts, & Prieto, 2017; Hadar, Steiner, Grant, & Rose, 1983; Ishi, Ishiguro, & Hagita, 2014), and listeners rely on these timely-aligned co-speech body movements to detect prominent syllables in speech (Krahmer & Swerts, 2007).

---

see Materials), parents reported that children were exposed to these other languages less than 50% of their daily life.

<sup>2</sup> Learning gains in children with the complex multilingual background did not differ significantly with respect to the other children (pre-test vs. post-test:  $\chi^2(1)=1.2205$ ,  $p=.2693$ ; pre-test vs. delayed post-test:  $\chi^2(1)=0.6101$ ,  $p=.4347$ ), and no interaction was found between linguistic background and the other predictors in our study (linguistic\_background\*contrast:  $\chi^2(2)=4.4867$ ,  $p=.1061$ ; linguistic\_background\*age:  $\chi^2(1)=2.4621$ ,  $p=.1166$ ; linguistic\_background\*condition:  $\chi^2(2)=0.8902$ ,  $p=.6408$ ).



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Insert Figure 1 about here

### *Stimuli for the training phrase*

One training trial consisted of 6 repetitions of a critical word that was inserted in the context of a meaningful phrase, as these contexts are found to help word-object mappings (Fennell & Waxman, 2010; Namy & Waxman, 2000). Each training trial contained two parts (one per item of the minimal pair) and had the following shape: «Look, it's a *voddy*! Look, a *voddy*! *Voddy* is nice! Nice *voddy*! Hey *voddy*! *Voddy*! [2-sec pause] [Part 2:] Look! It's a *boddy*! Look, a *boddy*! *Boddy* is nice! Nice *boddy*! Hey *boddy*! *Boddy*!». To avoid any **learning** bias, half of the trials had the /b/-/i/-WS words in the 1<sup>st</sup> part of the trial and the /v/-/i/-SW words in the second part, and the other half of the trials followed the opposite pattern.

The visual display in each training trial consisted of the native speaker appearing in the middle of the screen, plus the two objects referred to being displayed on the left bottom corner (first half of the trial) and right bottom corner of the screen (second half of the trial) (see Figure 2). In the consonant and vowel training trials, the object images were taken from the Novel Object and Unusual Name (NOUN) Database (Horst & Haust, 2016). In the lexical stress training trials, critical words were cognates and therefore children already had a cognitive representation for them in their L1. We consequently used images depicting the real meaning of the word for the real word of the minimal pair (e.g. a drawing of a dolphin for the critical word *dolphin*), and a drawing resembling the real meaning for the counterpart words in the minimal pair (e.g. a drawing of a dolphin-like animal for the critical word *dolphin*; **see Figure 2 as an example**). Adobe Premiere Pro was used to insert the object of reference into the visual display.

**Insert Figure 2 about here**

Each training trial was presented in 3 different visual conditions: a 'socially-engaging' condition, an 'ostensive-cueing' condition, and a 'visual-speech' condition. In the 'socially-engaging' condition, the speaker (native Northern American English) uttered the sentences containing the target phonemes and then alternated her gaze between the child and the object of reference (Figure 3, left panel). The speaker alternated the gaze 6 times during each training trial (one gaze alternation after each repetition of the critical word). In the 'visual-speech' condition, the speaker only directed her gaze to the child but not to the object of reference during the production of the sentences (Figure 3, middle panel). In the 'ostensive-cueing' condition, the speaker only directed her gaze to the object of reference but not to the child during the production of the sentences (Figure 3, right panel). In total, the speaker produced 72 training trials (3 visual conditions x 3 phonological contrasts x 8 critical words per contrast).

Each training condition was designed to enhance one of the three visual cues that were our focus of interest, despite some degree of cue-overlap between them. In all three training conditions the speaker's mouth was visible, but in the visual-speech condition young children were exposed longer (during the entire trial) and more clearly (front view) to visual speech cues. In all three conditions there was a triadic social interaction (child-speaker-object), but only in the socially-engaging situation the speaker acted in a referential and engaging way (the visual-speech situation was non-referential and in the ostensive-cueing the speaker did not establish eye contact with the child). Finally, in both the socially-engaging and the ostensive-cueing condition the speaker looked at the object of reference, but only in the ostensive-cueing

1  
2  
3 condition the speaker uttered the target sound while staring at it (in the socially-engaging  
4 situation the speaker first named the object and then directed her gaze towards it).  
5

6 The speaker, a trained prosodist, was asked to use the same intonation and prosodic features  
7 across visual conditions to avoid children being attentive to specific word-object relations due  
8 to acoustic salience. We nonetheless acoustically analysed all instances of critical words to check  
9 for any inconsistency in pitch range across visual conditions. Whenever a significant change was  
10 observed, we manipulated the pitch range in Praat (PSOLA Manipulation) to accommodate it to  
11 the mean. Table 1 shows the pitch range values of the critical words in the final stimuli across  
12 conditions. A linear regression analysis was applied to the data (with Pitch Range of the CW as  
13 the dependent variable and with Visual Condition as fixed factor) and results showed that the  
14 pitch range values of the final critical words did not vary across conditions in any occurrence  
15 within the training trial (see results in Table 1).  
16  
17  
18  
19

20 Insert Table 1 about here  
21

22  
23  
24 Insert Figure 3 about here  
25  
26  
27  
28

### 29 *Pre-, post-, and delayed post-training discrimination tests*

30  
31 A same-different AX task was used to test children's discrimination abilities. To create the stimuli  
32 for the pre- and post-training discrimination tests, the last repetition of the critical word in each  
33 training trial of the 'visual-speech' video-recordings was extracted. Stimuli were presented in an  
34 audio-visual format to children, but the visual display only showed the native speaker in the  
35 centre of the screen, not the objects referred to.  
36

37 In the AX task each child was presented with 24 test trials in a randomized order. There were 12  
38 AA trials (4 trials x 3 phonological contrasts) and 12 AB trials (4 trials x 3 phonological contrasts).  
39 In AA trials each phoneme/phonological pattern was represented twice (e.g. for the consonant  
40 contrast, there were 2 AA trials with /v/-/v/ and 2 AA trials with /b/-/b/). In AB trials the order  
41 of the phoneme/phonological pattern combination was also repeated twice (e.g. for the vowel  
42 contrast, there were 2 AB trials with a /i:-/i/ order and 2 AB trials with a /i/-/i:/ order) (see  
43 Appendix 1 for a list of test trials in the pre- and post-tests).  
44  
45  
46

### 47 **Procedure**

48 Children were assessed on their discrimination abilities before and after the word-label phonetic  
49 training. The pre-test occurred one week before the training. There were 3 training sessions,  
50 spaced at about 7-9 working day intervals. The post-test took place one week after the last  
51 training session and was followed by a delayed post-test three weeks later to check for long-  
52 term gains. Children were tested in a silent room at the school setting and they always wore  
53 Beyerdynamic DT-770 closed (noise-cancelling) headphones during the tasks.  
54  
55

### 56 *Pre-, post-, and delayed post-training discrimination tests*

57 A total of 8 practise trials (4 AA and 4 AB) preceded the beginning of the 24 test trials, to  
58 familiarize the child with the task and to check their ability to distinguish the concept 'same' and  
59 'different'. Practise trials contained phonological contrasts that were non-critical and that  
60

1  
2  
3 belonged to the children's L1 (e.g. /s/-/k/ or /k/-/t/). If the child did not respond correctly to any  
4 practise trial, the Experimenter gave feedback to the child by emphasizing the different  
5 phoneme in AB trials, or emphasizing the phonological coincidence in AA trials. Only when the  
6 Experimenter was sure that the child eventually understood why the two target words in the  
7 minimal pair were the same or different, the Experimenter moved to the next trial.  
8

9 A Power-Point presentation was used to advance through the trials. The Experimenter, who was  
10 the first author of the study, sat next to the child and handled the Power-Point presentation  
11 using a wireless mouse. To keep children interested in the task, test trials were alternated with  
12 attention-getting slides. In these slides animals were hidden behind coloured squares and  
13 children had to point at the square to be uncovered, one at a time. The test session lasted about  
14 15-20 minutes for 5-year-old children and 20-25 minutes for 4-year old children (younger  
15 children spent more time looking at the attention-getting slides).  
16  
17

18 After each trial children were asked to say whether they heard an exact repetition of the same  
19 word (target response in AA trials) or, instead, two words that varied slightly (target response in  
20 AB trials). The children's responses were manually coded in a score sheet by the Experimenter.  
21 To rule out the possibility that they would respond incorrectly because they would not be able  
22 to choose the appropriate verbal label for 'same' and 'different', their response was behavioural  
23 rather than verbal: using a small set of Lego pieces, they were asked to give 2 same pieces to  
24 the Experimenter if they heard 2 same words, or 2 different pieces if they heard 2 different  
25 words. They scored '0' if they responded 'same' in an AB trial or 'different' in an AA trial, and '1'  
26 if they responded 'different' in an AB trial or 'same' in an AA trial. If the child responded 'I don't  
27 know', the Experimenter played the stimulus again. If the child did not respond in the second  
28 repetition, the Experimenter asked the child to make a guess about the response. If the child  
29 responded but in an uncertain way (according to the Experimenter, as perceived by his/her facial  
30 expression or the prosodic features of the voice), the Experimenter asked 'Are you sure?' to the  
31 child. If the answer was 'yes', the Experimenter coded the response in the score sheet; if the  
32 answer was 'no', the Experimenter played the trial again to the child and asked him/her to  
33 respond again. In case of a divergence, the response that was coded was always the latest. This  
34 procedure had to be applied in 1.4% of the trials.  
35  
36  
37  
38

### 39 *Training sessions*

40  
41 Each of the 3 training sessions consisted of 12 training trials divided into 3 blocks, each block  
42 including 4 trials that trained one specific phonological contrast in a specific visual condition.  
43 Contrast and visual condition were counterbalanced and randomized across participants in a  
44 Latin Square design and resulted in three different lists so that each block changed according to  
45 the list. Participants were randomly assigned to the lists (see summary in Table 2).  
46  
47  
48

49 Insert Table 2 about here  
50  
51  
52

53 A Power-Point presentation was used to advance through the trials. The Experimenter sat next  
54 to the child and handled the Power-Point presentation using a wireless mouse. Children were  
55 asked to attend to the screen during the training sessions and did not have to perform any  
56 activity. Attention-getting slides were alternated between blocks to make sure children kept  
57 being attentive and maintained their interest in the trials. Attention-getting slides consisted of  
58 a display of several objects (e.g. fruits) that sometimes were repeated, and children had to point  
59 at repeated ones. The Experimenter only interacted with the child during the attention-getting  
60

1  
2  
3 slides to animate the task and make sure the child kept on being attentive. When the attention-  
4 getting task finished, the Experimenter told the child that the presentation would move to the  
5 next trial, and did not interact with the child until the next between-block pause. One training  
6 session lasted about 20-25 for 5-year-old children and 25-30 minutes for 4-year-old children  
7 (younger children were again slower in the attention-getting slides).  
8

9 The last training session was conducted with an eye-tracking system (Tobii X 120) to check for  
10 the children's gaze patterns during the learning situation. Because children showed traces of  
11 fatigue during the last trials of the second training session, this third training session was  
12 designed as significantly shorter: it included only three training trials, one trial per block (instead  
13 of the 4 trials per block in the previous session). Children's gaze patterns were recorded during  
14 the third training session because we estimated it was the session where children would display  
15 less anxiety (as they were already familiarized with the Experimenter, the experimental setting,  
16 and the task), and so their behaviour would be more natural and the gaze patterns would be  
17 better indicators of the children's phonological learning.  
18  
19

## 20 21 22 23 **RESULTS**

### 24 *Accuracy scores*

25  
26 Children's accuracy in discriminating between phonological contrasts was calculated using  $d'$   
27 scores, an unbiased signal detection theory measure that corrects for any potential bias in  
28 learners' responses (MacMillan & Creelman, 2005). For that, the hit rate (proportion of  
29 'different' responses in AB trials) and false alarm rate (proportion of 'different' responses in AA  
30 trials) were calculated for each participant and contrast, and then z-transformed. The final score  
31 represents a subtraction of the z-transformed false alarm rate from the z-transformed hit rate.  
32  
33

34 Figure 4 shows the  $d'$  scores across the three different time points (pre-test, post-test and  
35 delayed post-test) and as a function of type of visual condition. It reveals that children were  
36 more accurate in the post-test compared to the pre-test, independently of the visual condition  
37 in which they were trained. It also shows that children's accuracy upheld at the delayed post-  
38 test only when being trained with the 'ostensive-cueing' condition.  
39  
40  
41  
42  
43

44 Insert Figure 4 about here  
45  
46  
47  
48  
49

50 To investigate whether children's gains across testing sessions significantly varied as a function  
51 of the trained visual condition, the phonological contrast, and the children's age, two linear  
52 mixed effects models were fit using the *lmer* function of the lme4 package (Bates, Maechler, &  
53 Bolker, & Walker, 2015) in R (R Core Team, 2014). One model included gains between pre-test  
54 and post-test as the dependent variable, while the other model included gains between pre-test  
55 and delayed post-test as the dependent variable. In both models the fixed factors were visual  
56 condition (3 levels: 'socially-engaging', 'visual-speech', 'ostensive-cueing'), age (2 levels: 4 year  
57  
58  
59  
60

olds, 5 year olds), and phonological contrast (3 levels: consonant, vowel, lexical stress), and had a by-participant random slope for the effect of visual condition<sup>3</sup>.

The first model revealed that gains between pre-test and post-test were not affected significantly by visual condition ( $\chi^2(2) = .24, p = .88$ ), age ( $\chi^2(1) = .01, p = .92$ ), phonological contrast to be learned ( $\chi^2(2) = 1.20, p = .55$ ), nor by any interaction between these factors ( $\chi^2(2) = 4.42, p = .11$  for visual condition x age;  $\chi^2(4) = 5.68, p = .22$  for visual condition x phonological contrast;  $\chi^2(2) = 2.19, p = .33$  for phonological contrast x age). The second model showed that gains between pre-test and delayed post-test were significantly affected by visual condition ( $\chi^2(2) = 7.17, p < .05$ ), by an interaction between visual condition and age ( $\chi^2(2) = 8.47, p < .05$ ), and by an interaction between visual condition and phonological contrast ( $\chi^2(4) = 9.49, p = .05$ ). An inspection of the estimated coefficients of the second model showed, first, that the main effect of visual condition was due to the fact that children learned significantly more in the 'ostensive-cueing' condition compared to the 'socially-engaging' ( $\beta = -1.38, SE = .57, t = -2.42, p < .05$ ) and 'visual-speech' ( $\beta = -1.21, SE = 0.57, t = -2.13, p < .05$ ) conditions, the last two not differing between each other ( $\beta = .16, SE = .57, t = .29, p = .77$ ). Second, and as illustrated by Figure 5, we found that visual condition and age interacted in that 4-year-olds learned significantly more than 5-year-olds in the 'ostensive-cueing' condition ( $\beta = -3.01, SE = .93, t = -3.23, p < .01$ ), all other comparisons being non-significant (age 4 vs. age 5 in 'visual-speech' condition:  $\beta = -.26, SE = .93, t = -.28, p = .77$ ; age 4 vs. age 5 in 'socially-engaging' condition:  $\beta = -.17, SE = .93, t = -.18, p = .85$ ). Third, the interaction between visual condition and phonological contrast was due to the 'ostensive-cueing' condition leading to higher gains in the consonant contrast ('ostensive-cueing' vs. 'socially-engaging':  $\beta = -2.45, SE = 1.13, t = -2.16, p < .05$ ; 'ostensive-cueing' vs. 'visual-speech':  $\beta = -1.81, SE = 1.13, t = -1.59, p = .11$ ; 'socially-engaging' vs. 'visual-speech':  $\beta = .64, SE = 1.16, t = .56, p = .58$ ), and in the vowel contrast ('ostensive-cueing' vs. 'socially-engaging':  $\beta = -2.48, SE = 1.13, t = -2.19, p < .05$ ; 'ostensive-cueing' vs. 'visual-speech':  $\beta = -0.80, SE = 1.15, t = -0.69, p = .49$ ; 'visual-speech' vs. 'socially-engaging':  $\beta = -1.68, SE = 1.13, t = -1.48, p = .14$ ). Instead, for learning the lexical stress contrast, the comparison across conditions was not significantly different ('socially-engaging' vs. 'visual-speech':  $\beta = -1.92, SE = 1.13, t = -1.69, p = .09$ ; 'socially-engaging' vs. 'ostensive-cueing':  $\beta = -.89, SE = 1.15, t = -.77, p = .44$ ; 'visual-speech vs. 'ostensive-cueing':  $\beta = 1.03, SE = 1.13, t = .91, p = .36$ ). Neither age nor phonological contrasts, nor the interaction between these two came out as significant ( $\chi^2(2) = 2.67, p = .11$ ;  $\chi^2(2) = 3.71, p = .16$ ;  $\chi^2(2) = .84, p = .66$ , respectively).

Insert Figure 5 about here

### *Children's gaze patterns*

Three Areas of Interest (Aoi) were defined in the training materials: the speaker's mouth, the speaker's eyes, and the object of reference (either on the bottom left or on the bottom right corner of the screen). Because the position of the speaker's mouth and eyes slightly varied across video frames, especially in the 'socially-engaging' condition, these Aois were set in a dynamic way in order to account for the distinct positions of the target Aoi in the visual space. Assuming that it takes about 200 ms for the eyes to program a saccade in reaction to a linguistic stimulus (e.g. Altmann & Kamide, 2004; Matin, Shao, & Boff, 1993; Salverda, Kleinschmidt, & Tanenhaus, 2014), we extracted children's gaze patterns at the onset of the critical words and

<sup>3</sup> Item was not included as a random factor because item variation was removed when calculating  $d'$  scores, since the proportion of false alarms and hits was calculated for each condition (all items together) by participant.

1  
2  
3 until the critical word ended, adding a leeway of 200 ms from the offset of the target word. For  
4 the analyses we only considered children's gaze patterns during critical words, and hence looks  
5 during the sentence context of these words were excluded.  
6

7 Figure 6 shows the amount of time spent on each AoI for each trained visual condition. Overall  
8 children looked more at the mouth than at the other AoI in the three trained conditions. Within  
9 each AoI, some differences can be observed across trained conditions: children seemed to look  
10 longer at the eyes in the socially-engaging condition, and this same condition also seemed to  
11 elicit more looks at the mouth than any other condition. In contrast, children seemed to spend  
12 more time looking at the object of reference in the ostensive-cueing condition (a bit less in the  
13 visual-speech condition and even less in the socially-engaging condition).  
14  
15  
16  
17

18 Insert Figure 6 about here  
19  
20  
21  
22

23 Three linear mixed models were applied to the data to explore children's looking patterns across  
24 visual conditions, phonological contrasts, and age, using the *lmer* function of the *lme4* package  
25 (Bates, Maechler, & Bolker, & Walker, 2015) in R (R Core Team, 2014). The first model explored  
26 the odds ratio of time looking at the mouth versus at the other two AoI; the second model  
27 explored the odds ratio of time looking at the object of reference versus at the other AoI; the  
28 third model investigated odds ratio of looking shifts between mouth and object of reference  
29 versus gazes to mouth, object or eyes (without shifting between regions). In all models fixed  
30 factors were visual condition (3 levels: 'socially-engaging', 'visual-speech', 'ostensive-cueing'),  
31 age (2 levels: 4 year olds, 5 year olds), and phonological contrast (3 levels: consonant, vowel,  
32 lexical stress). Participant and item were set as random factors. The structure of the models was  
33 determined by our research predictions: if the relevant visual information is the processing of  
34 mouth and lip speech movements, fixations on the mouth would be higher in conditions with  
35 higher phonological gains (first model); instead, if the relevant visual information is an ostensive  
36 cueing of the relation between phonetic input and object of reference, fixations on the object  
37 of reference would be higher in conditions with higher phonological gains (second model); or,  
38 in contrast, if the relevant visual information is the social interaction and triadic joint attentional  
39 frame, gaze shifts between mouth and object would be higher in conditions with higher  
40 phonological gains (third model).  
41  
42  
43

44 The results of the first model (time looking at mouth vs. at other AoI) showed a main effect of  
45 phonological contrast ( $\chi^2(2) = 17.38, p < .01$ ) and a marginal effect of visual condition ( $\chi^2(2) =$   
46  $5.1991, p = .07$ ), but no main effect of age ( $\chi^2(1) = .24, p = .62$ ) nor any interaction between the  
47 three factors (age x visual condition:  $\chi^2(2) = 1.05, p = .59$ ; age x phonological contrast:  $\chi^2(2) = .67,$   
48  $p = .72$ ; visual condition x phonological contrast:  $\chi^2(4) = 3.38, p = .50$ ). The estimated coefficients  
49 reveal that children spent more time looking at the mouth when presented with the vowel  
50 contrast and when being trained with the 'socially-engaging' condition (see estimated  
51 coefficients of this and the other two models in Table 3).  
52  
53

54 The second model (time looking at object of reference vs. at other AoI) revealed a main effect  
55 of visual condition ( $\chi^2(2) = 8.15, p < .05$ ) and a main interaction between visual condition and  
56 phonological contrast ( $\chi^2(4) = 15.33, p < .01$ ), all other main effects and interactions being non-  
57 significant (main effect of age:  $\chi^2(1) = 1.15, p = .28$ ; main effect of phonological contrast:  $\chi^2(2) =$   
58  $3.51, p = .17$ ; age x visual condition:  $\chi^2(2) = .88, p < .25$ ; age x phonological contrast:  $\chi^2(2) = .74, p$   
59  $= .69$ ). The estimated coefficients indicate that children looked more at the object of reference  
60 in the 'ostensive-cueing' condition than in the other conditions, and that in the 'ostensive-

1  
2  
3 cueing' condition, the consonant contrast triggers more looks at the object of reference than  
4 the vowel or lexical stress contrasts (see Table 3 and Figure 7).  
5  
6  
7

8 Insert Figure 7 about here  
9

10  
11 Insert Table 3 about here  
12  
13

14 The third model (looking shifts between mouth and object of reference vs. 'static' gazes to  
15 mouth, object, or eyes) showed no main effects of visual condition ( $\chi^2(2) = 3.33, p = .19$ ),  
16 phonological contrast ( $\chi^2(2) = 3.70, p = .16$ ) or age ( $\chi^2(1) = 1.50, p = .22$ ), nor any interaction  
17 between any of these factors (visual condition x age:  $\chi^2(2) = 2.14, p = .34$ ; visual condition x  
18 phonological contrast:  $\chi^2(4) = 1.12, p = .89$ ; phonological contrast x age:  $\chi^2(2) = 4.56, p = .10$ ).  
19 These results indicate that children shifted their gaze between the speaker's mouth and object  
20 of reference at a similar proportion in all training visual conditions and independently of the  
21 phonological contrasts to be learned and of their age (see all coefficients in Table 3).  
22  
23

## 24 25 DISCUSSION

26  
27 The present study aimed at investigating which enriching visual cues is particularly helpful for  
28 training pre-schoolers' perception of non-native phonemes: that in which learners can easily  
29 access visual speech information, that in which novel phonemes are presented in a socially-  
30 engaging situation, or that in which there is a clear link between a reference entity and the  
31 phonological input. We designed a training study in which Catalan/Spanish-speaking pre-  
32 schoolers were presented with three novel English phonological contrasts in a context of an  
33 object-labeling task: the /b/-/v/ distinction (Catalan/Spanish learners of English assimilate /v/ to  
34 /b/), the /i:/-/i/ contrast (Catalan/Spanish learners of English assimilate /i:/ to the native /i/ and  
35 /i/ either to /i/ or to /e/), and the SW-WS contrast (although Catalan, Spanish and English are  
36 languages with lexical stress, many cognate words have the opposite stress pattern in both  
37 languages: WS in Catalan and Spanish but SW in English). The pre-schoolers' accuracy in  
38 perceiving these contrasts was assessed before and after training, to evaluate which trained  
39 enriching visual cue leads to higher phonological gains, and the children's gaze preferences  
40 during the training phase were recorded using an eye-tracker, to investigate the relation  
41 between the children's focus of attention and their learning gains.  
42  
43

44 The analysis revealed that children's learning of non-native phonological contrasts is boosted in  
45 the 'ostensive-cueing' condition, as this visual condition contributed to higher long-term gains  
46 than the other visual conditions. When children were assessed immediately after the training  
47 (post-test), similar gains were observed across all trained visual conditions, children's age, or  
48 nature of the novel phonological contrast (consonant, vowel, stress position). However, when  
49 children's gains were evaluated some weeks after the training (delayed post-test), a main effect  
50 of visual condition arose and the 'ostensive-cueing' condition emerged as the ideal frame for  
51 acquiring L2 phonological contrasts. Our results on pre-schoolers align with previous results  
52 reporting crucial effects of referential cues for young infants' phonological categorization in  
53 object-labeling tasks (Fennell & Waxman, 2010; Yeung, Chen, & Werker, 2014; Yeung & Werker,  
54 2009). In all learning conditions the object of reference was visually available when the novel  
55 phonological information was uttered, but only in the 'ostensive-cueing' stimuli the L2 speaker  
56 was ostensibly looking at the object of reference during the entire trial, and therefore also when  
57 naming the target object and producing the non-native phonological input to be learned. We  
58  
59  
60

1  
2  
3 argue that this visual situation helped establishing an unequivocal link between phonological  
4 input and its meaning in the real world, and hence enhanced phonological learning.  
5

6 Previous studies on young infants had also reported that referential and lexical factors influence  
7 positively the acquisition of novel phonological categories at early stages in development. At the  
8 level of the word form, it has been found that infants form novel phonological categories when  
9 the fine-grained detailed acoustical information is paired to distinct lexical contexts (Feldman et  
10 al., 2013; Thiessen, 2011). At the level of word meaning, infants acquire novel phonemes if these  
11 have a functional value, i.e. if there is a consistent pairing between novel phonemes and object  
12 of reference that reinforce phonetic distinctiveness (Fennell & Waxman, 2010; Yeung & Werker,  
13 2009). Interestingly, the importance of meaningful interactions for language learning does not  
14 seem to decline with age. During preschool years, shared attention cues and adult  
15 responsiveness are a better predictor of language skills than simple exposure (Romeo et al.,  
16 2018; Rowe & Best, 2020), and our study suggests that pre-schoolers' phonological learning is  
17 boosted when the presence of an object of reference is accompanied by ostensive unequivocal  
18 signs of the sound-object relation.  
19  
20

21 Next to the essential status of referential clarity, previous work had found that learners use  
22 social interactive joint-attentional frames (Bannard & Tomasello, 2012; Kuhl et al., 2003; Kuhl,  
23 2007; Roseberry & Kuhl, 2013) and visual speech information (i.e. the observation of the  
24 interlocutor's lip movement when speaking) to discriminate between non-native phonemes (eg.  
25 Hardison, 2003; Hazan et al., 2006; Ortega-Llebaria et al., 2001; Ter Schure et al., 2016; Weikum  
26 et al 2007). In our study we also found that children looked at the speaker's mouth more than  
27 her eyes or the object of reference, when the critical phonological input unfolded. Because this  
28 happened in the three trained visual conditions, we can presuppose that pre-schoolers attend  
29 to the speaker's mouth by default in any learning situation and independently of the visual  
30 context. This means that providing clear visual-speech information does not imply that children  
31 will look more at the mouth and that they will improve their non-native phoneme perception.  
32 Similarly, the socially-engaging situation we designed did not significantly contribute to make  
33 children shifting their gaze between the L2 speaker and the object of reference, nor boosted  
34 their non-native phoneme acquisition. Instead, providing ostensive cues to the object-label  
35 mapping did enhance children's long-term learning of non-native phonemes.  
36  
37  
38

39 Why is it that, compared to ostensive-cueing signs, visual speech and socially-engaging cues  
40 have a lesser contribution to long-term phonological gains? Several explanations may be  
41 suggested. One relates to the learners' age. Existing theoretical accounts propose that the value  
42 of any feature of the linguistic input for learning varies depending on the learner's age and the  
43 learner's linguistic ability (Rowe & Snow, 2020). Longitudinal studies provide the ideal evidence  
44 to see which features of the visual input matter more for non-native phonological acquisition at  
45 each developmental stage, but they are scarce. In his seminal study, McGurk and MacDonald  
46 (1976) tested pre-schoolers, school-aged children and adults, and found that adults were more  
47 influenced by visual speech information than school-aged or preschool children when perceiving  
48 speech sounds. Similarly, Erdener and Burnham (2013) compared speech perception by 5-, 6-,  
49 7-, 8-year-old children and adults across two conditions (matching vs. mismatching audio-visual  
50 McGurk stimuli), and found that adults were more influenced than children by visual speech  
51 information. Cross-sectional studies with young infants find that visual speech input has a  
52 positive effect on the listener's perception of non-native sounds (Teinonen, Aslin, Alku, 2008;  
53 Ter Schure et al., 2016; but see divergent results in Pons et al., 2009), and cross-sectional studies  
54 with adults obtain similar findings (Cebrian & Carlet, 2012; Ortega-Llebaria et al., 2001).  
55 Tentatively, it could be suggested that the importance of visual speech information follows a u-  
56 shaped pattern: stronger effect in infancy and adulthood, and a temporary regression in  
57 childhood. Indeed, some accounts propose that distinct cues matter differently in the course of  
58 development (Hollich, Hirsh-Pasek, & Golinkoff, 2000), and non-linear developmental  
59  
60



trajectories have been observed in various dimensions of language acquisition (Gershkoff-Stowe & Thelen, 2004; Marcus, 2004). Future longitudinal studies that include infants, children, an adult populations are needed to investigate this potential age-related effect in more detail.

The non-linearity explanation seems to be less appropriate to interpret the lack of boosting effect of the socially-engaging situation. Previous cross-sectional studies find that social interactions are crucial in young infants' language learning (Bannard & Tomasello, 2012; Kuhl et al., 2003; Kuhl, 2007; Roseberry & Kuhl, 2013), and that the importance of social cues in language learning does not decrease with age (Rowe & Snow, 2020). It could be, however, that the specific features of the social interaction are determinant. In our 'socially-engaging' training stimuli the L2 speaker uttered the phrase containing the critical phonological contrast and then turned her gaze towards the object of reference. From a novice learner point of view, this might have blurred the relation between linguistic input and object, as children might have not known if any portion of the carrier sentence referred to the object. Previous studies have observed that learning takes place when the young learner is focusing their attention on the relevant referent and at the same exact time they perceive the linguistic input that refers to this element (Yu & Smith, 2012). In the successful ostensive-cueing condition the adult speaker provided the phonological input while directing her gaze to the object of reference, and this might have reduced the referential ambiguity of the non-native phonological input. The children's fixations on the object of reference are in line with these time-related effects: young children fixated on the object of reference significantly more in the ostensive-cueing condition, increasing the children's chances to process the non-native phoneme while attending to its object of reference.

Providing learners with the relevant linguistic information in synchrony with the learners' focus of interest is typical of contingent contexts. In a contingent situation the adult follows the infant's focus of interest and therefore provides linguistic input that refers to meanings that are highly relevant for the child. Our eye-tracking results show that the most contingent visual context in our study was the ostensive-cueing situation, as the speaker provided the linguistic input that referred to the actual child's focus of interest. Previous research reported positive correlations between social contingency and infants' linguistic development (e.g. Bannard & Tomasello, 2012; Hakuno et al., 2017; McGillion et al., 2017; Nussenbaum & Amso, 2015; Roseberry et al., 2014, see Mermelshtine, 2017, for a review). Our study adds to this body of findings by supporting the positive effects of contingent behaviours in young children's non-native phonological acquisition, but also call for a more precise picture: we found that the younger the learners, the more they may need contingent interactions to acquire sounds that do not belong to their native phonological system.

This study investigated segmental and suprasegmental contrasts. At the segmental level, children learned to discriminate non-native consonant and vowel contrasts, and at the suprasegmental level children learned to relocate the stressed syllable in cognate words. Our eye-tracking results show that the vowel contrast triggered more looks at the mouth than any other contrast, while the consonant contrast triggered more looks at the object of reference than any other contrast. The acoustic salience of each contrast, and its relative distance with the native category, might explain this pattern of results, as predicted by the Perceptual Assimilation Model (Best & Tyler, 2007). As a trade-off effect, learners might have looked more to the visual speech information (the mouth) when there was less acoustic salience, and they might have looked less at the mouth (and more at the object) when acoustic salience was higher. The tense /i:/ - lax /ɪ/ vowel contrast is particularly challenging for Catalan/Spanish learners of English, as neither Catalan nor Spanish have any tense-lax contrast in their vocalic system. It has been found that the acoustic signal is not highly reliable for learners, as they are found to associate the lax /ɪ/ vowel to the /i/ and /e/ categories (Cebrian, 2006), so the children in the present study might have focused more on the mouth to compensate for the lack of acoustic reliability. Instead, Catalan and Spanish consonantal inventory does not include the voiced labiodental fricative /v/,

1  
2  
3 but both languages have bilabial-labiodental contrast (as in /b/-/f/). The consonant contrast  
4 might have been more acoustically salient, and so they focused less on visual speech input and  
5 more on the object of reference.  
6  
7

8 One of the limitations of the current study is the lack of speaker variability during the distinct  
9 sessions. The same L2 speaker produced the stimuli during the training sessions and was also  
10 presented in the pre- and post- discrimination tests. Our current results do not show if children  
11 generalize the learned patterns to new speakers. Likewise, we did not include non-trained  
12 stimuli in the post-test discrimination tasks to look for generalization. While the inclusion of the  
13 delayed post-test task was an attempt to investigate long-term more stable phonological gains,  
14 we agree that future work should address these concerns.  
15  
16  
17

18 For many children in the world the preschool period is the time when they first get in contact  
19 with an additional language, either because the language used at home differs from that in the  
20 school setting, or because it is when the school system starts the formal instruction of a second  
21 language. Among other things, these children need to be sensitive to the fact that the new  
22 language has a distinct phonological system, and therefore need to construct new phonological  
23 categories that might not exist in their native language. Investigating how they manage to do it  
24 was one of the motivations for the present study. Our study showed that young pre-schoolers  
25 learn non-native phonemes better if they are presented with ostensive signs of the relevant  
26 object of reference, as in contingent situations. The presence of clear and contingent mapping  
27 between linguistic input and referential function seems to outrank the availability of clear visual  
28 speech input or the amount of social interactivity, at least for the contrasts we studied, at  
29 preschool age, and when gains are measured in the long-term. We do not claim that visual  
30 speech cues or social engagement could or should be erased from a learning situation, only that  
31 an ideal learning situation for pre-schoolers is one in which children have the opportunity to  
32 clearly match what the interlocutor says with what the interlocutor means.  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **Acknowledgements**  
4

5 Thanks to teacher, parents, and children at the [removed for review] school in [removed for  
6 review]. We thank [removed for review] for his advice in analyzing the results. This research was  
7 funded by the [removed for review] grant and [removed for review].  
8  
9

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Review Only

1  
2  
3 **References**  
4

- 5 Aliaga-Garcia, C. (2017). *The effect of auditory and articulatory phonetic training on the*  
6 *perception and production of L2 vowels by Catalan-Spanish learners of English*. PhD  
7 Dissertation. Universitat de Barcelona.  
8
- 9 Aliaga-Garcia, C. & Mora, J. C. (2009). Assessing the effects of phonetic training on L2 sound  
10 perception and production. In: A M. A. Watkins, A. S. Rauber, & B. O. Baptista (Eds.),  
11 *Recent Research in Second Language Phonetics/Phonology: Perception and Production*.  
12 Newcastle upon Tyne: Cambridge Scholars Publishing (pp. 2-31).  
13
- 14 Alm, M., Behne, D. M., Wang, Y., & Eg, R. (2009). Audio-visual identification of place of  
15 articulation and voicing in white and babble noise. *The Journal of the Acoustical Society of*  
16 *America*, 126(1), 377–387.  
17
- 18 Altmann, G. T. M., & Kamide, Y. (2004). Now You See It, Now You Don't: Mediating the Mapping  
19 between Language and the Visual World. In J. Henderson & F. Ferreira (Eds.), *The interface*  
20 *between language, vision, and action: Eye movements and the visual world*. New York;  
21 Hove, UK: Psychology Press (pp. 347–386).  
22
- 23 Bannard, C., & Tomasello, M. (2012). Can We Dissociate Contingency Learning from Social  
24 Learning in Word Acquisition by 24-Month-Olds? *PLoS ONE*, 7(11), e49881.  
25
- 26 Bates, D., Maechler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models  
27 Using lme4. *Journal of Statistical Software*, 67(1), 1-48.  
28
- 29 Best, C. T. (1993). Emergence of language-specific constraints in perception of non-native  
30 speech: A window on early phonological development. In: De Boysson-Bardies, B., de  
31 Schonen, S., Jusczyk, P., MacNeilage, P., & Morton, J. (Eds.). *Developmental*  
32 *neurocognition: Speech and face processing in the first year of life* (pp. 289–304).  
33 Dordrecht, The Netherlands: Academic Publishers B.V.  
34
- 35 Best, C. T., & Tyler, M. D. (2007). Nonnative and second-language speech perception:  
36 Commonalities and complementarities. In M. J. Munro & O.-S. Bohn (Eds.), *Second*  
37 *language speech learning: The role of language experience in speech perception and*  
38 *production* (pp. 13–34). Amsterdam, The Netherlands: Benjamins.  
39
- 40 Birulés, J., Bosch, L., Brieke, R., Pons, F., & Lewkowicz, D. J. (2019). Inside bilingualism: Language  
41 background modulates selective attention to a talker's mouth. *Developmental Science*,  
42 22(3), 1–12.  
43
- 44 Brooks, R., & Meltzoff, A. (2005). The development of gaze following and its relations to  
45 language. *Developmental Science*, 8, 535–543.  
46
- 47 Cebrian, J. (2006) Experience and the use of duration in the categorization of L2 vowels. *Journal*  
48 *of Phonetics* 34(3), 372-387.  
49
- 50 Cebrian, J., & Carlet, A. (2012). Audiovisual perception of native and non-native sounds by native  
51 and non-native speakers 1. In S. Martín Alegre, M. Moyer, E. Pladevall, & S. Tubau (Eds.),  
52 *At a Time of Crisis: English and American Studies in Spain* (pp. 300–307).  
53
- 54 Csibra, G., & Gergely, G. (2009). Natural pedagogy. *Trends in Cognitive Sciences*, 13(4), 148–153.  
55
- 56 Erdener, D., & Burnham, D. (2013). The relationship between auditory-visual speech perception  
57 and language-specific speech perception at the onset of reading instruction in English-  
58 speaking children. *Journal of Experimental Child Psychology*, 116(2), 120–138.  
59  
60

1  
2  
3 Erdener, D., & Burnham, D. (2018). Auditory–visual speech perception in three- and four- year-  
4 olds and its relationship to perceptual attunement and receptive vocabulary. *Journal of*  
5 *Child Language*, 45(2), 273–289.

6  
7 Erdener, V. D. (2007). *Development of Auditory-Visual speech perception in young children*.

8  
9 Esteve-Gibert, N., Borràs-Comes, J., Asor, E., Swerts, M., & Prieto, P. (2017). The timing of head  
10 movements: The role of prosodic heads and edges. *The Journal of the Acoustical Society*  
11 *of America*, 141(6), 4727–4739.

12  
13 Esteve-Gibert, N., & Guellaï, B. (2018). Prosody in the auditory and visual domains: A  
14 developmental perspective. *Frontiers in Psychology*, 9(March), 1–10.

15  
16 Feldman, N. H., Myers, E. B., White, K. S., & Morgan, J. L. (2013). Word-level information  
17 influences phonetic learning in adults and infants. *Cognition*, 127(3), 427–438.

18  
19 Fennell, C. T., & Waxman, S. R. (2010). What paradox? Referential cues allow for infant use of  
20 phonetic detail in word learning. *Child Development*, 81(5), 1376–1383.

21  
22 Fernald, A. (1993). Approval and disapproval: infant responsiveness to vocal affect in familiar  
23 and unfamiliar languages. *Child Development*, 64, 657–674.

24  
25 Fernald, A., & Mazzie, C. (1991). Prosody and focus in speech to infants and adults.  
26 *Developmental Psychology*, 27(2), 209–221.

27  
28 Flege, J. (1989). Differences in inventory size affect the location but not the precision of tongue  
29 positioning in vowel production. *Language and Speech*, 32(2), 123–147.

30  
31 Gershkoff-Stowe, L., & Thelen, E. (2004). U-shaped changes in behavior: A dynamic systems  
32 perspective. *Journal of Cognition and Development*, 5, 11–36.

33  
34 Gick, B., & Derrick, D. (2009). Aero-tactile integration in speech perception. *Nature*, 462(7272),  
35 502–504.

36  
37 Gogate, L. J., Walker-Andrews, A. S., & Bahrick, L. E. (2001). The intersensory origins of word  
38 comprehension: An ecological-dynamic systems view. *Developmental Science*, 4(1), 1–18.

39  
40 Golinkoff, R. M., Can, D. D., Soderstrom, M., & Hirsh-Pasek, K. (2015). (Baby)Talk to Me: The  
41 Social Context of Infant-Directed Speech and Its Effects on Early Language Acquisition.  
42 *Current Directions in Psychological Science*, 24, 339–344.

43  
44 Hadar, U., Steiner, T. J., Grant, E. C., & Rose, F. C. (1983). Head movement correlates of juncture  
45 and stress at sentence level. *Language and Speech*, 26(2), 117–129.

46  
47 Hakuno, Y., Omori, T., Yamamoto, J., & Minagawa, Y. (2017). Social interaction facilitates word  
48 learning in preverbal infants: Word–object mapping and word segmentation. *Infant*  
49 *Behavior and Development*, 48(Part B), 65–77.

50  
51 Hall, G. F. (1991). *Perceptual and associative learning*. Oxford, UK: Clarendon Press.

52  
53 Hanna, J.E. & Brennan, S.E. (2007). Speakers' eye gaze disambiguates referring expressions early  
54 during face-to- face conversation. *Journal of Memory and Language*, 57(4), 596–615.

55  
56 Hardison, D. M. (2003). Acquisition of second-language speech: Effects of visual cues, context,  
57 and talker variability. *Applied Psycholinguistics*, 24(4), 495–522.

58  
59 Hazan, V., & Barrett, R. (2000). The development of phonemic categorization in children aged 6-  
60 12 years. *Journal of Phonetics*, 28, 377–396.

- 1  
2  
3 Hazan, V., Sennema, A., Faulkner, A., Ortega-Llebaria, M., Iba, M., & Chung, H. (2006). The use  
4 of visual cues in the perception of non-native consonant contrasts. *The Journal of the*  
5 *Acoustical Society of America*, 119(3), 1740–1751.  
6  
7 Hollich, G. J., Hirsh-Pasek, K., & Golinkoff, R. M. (2000). Breaking the Language Barrier: An  
8 Emergentist Coalition Model of Word Learning (Monographs of the Society for Research  
9 in Child Development). *Monographs of the Society for Research in Child Development*,  
10 65(3), 1–135.  
11  
12 Horst, J. S., & Hout, M. C. (2016). Novel Object & Unusual Name (NOUN) Database: A collection  
13 of novel images for use in experimental research. *Behavior Research Methods*, 48(4),  
14 1393–1409.  
15  
16 Ishi, C. T., Ishiguro, H., & Hagita, N. (2014). Analysis of relationship between head motion events  
17 and speech in dialogue conversations. *Speech Communication*, 57, 233–243.  
18  
19 Krahmer, E., & Swerts, M. (2007). The effects of visual beats on prosodic prominence: Acoustic  
20 analyses, auditory perception and visual perception. *Journal of Memory and Language*,  
21 57(3), 396–414.  
22  
23 Kuhl, P. K. (2007). Is speech learning “gated” by the social brain? *Developmental Science*, 10,  
24 110–120.  
25  
26 Kuhl, P. K., Tsao, F.-M., & Liu, H.-M. (2003). Foreign-language experience in infancy: effects of  
27 short-term exposure and social interaction on phonetic learning. *Proceedings of the*  
28 *National Academy of Sciences of the United States of America*, 100, 9096–9101.  
29  
30 Lalonde, K., & Holt, R. F. (2015). Preschoolers Benefit From Visually Salient Speech Cues. *Journal*  
31 *of Speech, Language, and Hearing Research*, 58(1), 135–150.  
32  
33 Linebarger, D. L., & Vaala, S. E. (2010). Screen media and language development in infants and  
34 toddlers: An ecological perspective. *Developmental Review*, 30(2), 176–202.  
35  
36 MacMillan, N.A. & Creelman, C.D. (2005) *Detection theory: A user's guide* (2nd ed.). Mahwah,  
37 NJ.: Lawrence Erlbaum Associates.  
38  
39 Marcus, G. F. (2004). What's in a U? The Shapes of Cognitive Development. *Journal of Cognition*  
40 *and Development*, 5(1), 119–122.  
41  
42 Matin, E., Shao, K. C., & Boff, K. R. (1993). Saccadic overhead: Information-processing time with  
43 and without saccades. *Perception and Psychophysics*, 53, 372–380.  
44  
45 Maye, J., Weiss, D.J., & Aslin, R.N. (2008). Statistical phonetic learning in Infants: facilitation and  
46 feature generalization. *Developmental Science*, 11(1), 122–134.  
47  
48 Maye, J., Werker, J. F., & Gerken, L. (2002). Infant sensitivity to distributional information can  
49 affect phonetic discrimination. *Cognition*, 82, B101–B111  
50  
51 McGillion, M., Herbert, J. S., Pine, J., Vihman, M., Keren-portnoy, T., & Matthews, D. (2017).  
52 What Paves the Way to Conventional Language? The Predictive Value of Babble, Pointing,  
53 and Socioeconomic Status. *Child Development*, 88(1), 156–166.  
54  
55 McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature*, 264, 746–748.  
56  
57 Mermelshstine, R. (2017). Parent–child learning interactions: A review of the literature on  
58 scaffolding. *British Journal of Educational Psychology*, 87(2), 241–254.  
59  
60 Metsala, J. L., & Walley, A. C. (1998). Spoken vocabulary growth and the segmental restructuring  
of lexical representations: Precursors to phonemic awareness and early reading ability. In

- 1  
2  
3 J. L. Metsala & L. C. Ehri (Eds.), *Word recognition in beginning literacy* (pp. 89 – 120).  
4 Mahwah, NJ: LEA.  
5
- 6 Moore, C., Angelopoulos, M., & Bennett, P. (1999). Word learning in the context of referential  
7 and salience cues. *Developmental Psychology*, 35, 60–68.  
8
- 9 Namy, L.L. & Waxman, S.R. (2000). Naming and exclaiming: Infants' sensitivity to naming  
10 contexts. *Journal of Cognition and Development*, 1, 405–428.  
11
- 12 Nittrouer, S. (1996). Discriminability and perceptual weighting of some acoustic cues to speech  
13 perception by 3-year-olds. *Journal of Speech and Hearing Research*, 39, 278–297.  
14
- 15 Nittrouer, S. (2005). Age-related differences in weighting and masking of two cues to word-final  
16 stop voicing in noise. *Journal of the Acoustical Society of America*, 118, 1072–1088.  
17
- 18 Nussenbaum, K., & Amso, D. (2015). An Attentional Goldilocks Effect: An Optimal Amount of  
19 Social Interactivity Promotes Word Learning From Video. *Journal of Cognition and  
20 Development*, 17(1), 30–40.  
21
- 22 Ohde, R. N., & German, S. R. (2011). Formant onsets and formant transitions as developmental  
23 cues to vowel perception. *Journal of the Acoustical Society of America*, 130, 1628–1642.  
24
- 25 Ortega-Llebaria, M., Faulkner, A., & Hazan, V. (2001). Auditory-Visual L2 Speech Perception:  
26 Effects of Visual Cues and Acoustic-Phonetic Context for Spanish Learners of English.  
27 *Speech, Hearing and Language: Work in Progress*, 13, 3951.  
28
- 29 Pons, F., Lewkowicz, D. J., Soto-Faraco, S., & Sebastián-Gallés, N. (2009). Narrowing of  
30 intersensory speech perception in infancy. *Proceedings of the National Academy of  
31 Sciences of the United States of America*, 106(26), 10598–10602.  
32
- 33 Roseberry, S., Hirsh-Pasek, K., & Golinkoff, R. M. (2014). Skype Me! Socially Contingent  
34 Interactions Help Toddlers Learn Language. *Child Development*, 85(3), 956–970.  
35
- 36 Saint-Georges, C., Chetouani, M., Cassel, R., Apicella, F., Mahdhaoui, A., Muratori, F., Laznik, M.  
37 C., & Cohen, D. (2013). Motherese in Interaction: At the Cross-Road of Emotion and  
38 Cognition? (A Systematic Review). *PLoS ONE*, 8(10), 1–18.  
39
- 40 Salverda, A. P., Kleinschmidt, D., & Tanenhaus, M. K. (2014). Immediate effects of anticipatory  
41 coarticulation in spoken-word recognition. *Journal of Memory and Language*, 71(1), 145–  
42 163.  
43
- 44 Schwartz, J. L., Berthommier, F., & Savariaux, C. (2004). Seeing to hear better: Evidence for early  
45 audio-visual interactions in speech identification. *Cognition*, 93(2), 69–78.  
46
- 47 Senju, A., & Csibra, G. (2008). Gaze following in human infants depends on communicative  
48 signals. *Current Biology*, 18(9), 668–671.  
49
- 50 Soderstrom, M. (2007). Beyond babytalk: Re-evaluating the nature and content of speech input  
51 to preverbal infants. *Developmental Review*, 27(4), 501–532.  
52
- 53 Spinelli, M., Fasolo, M., & Mesman, J. (2017). Does prosody make the difference? A meta-  
54 analysis on relations between prosodic aspects of infant-directed speech and infant  
55 outcomes. *Developmental Review*, 44, 1–18.  
56
- 57 Stager, C.L. & Werker, J.F. (1997). Infants listen for more phonetic detail in speech perception  
58 than in word-learning tasks. *Nature*, 388(6640), 381–382.  
59
- 60 Teinonen, T., Aslin, R. N., Alku, P., & Csibra, G. (2008). Visual speech contributes to phonetic  
learning in 6-month-old infants. *Cognition*, 108(3), 850–855.

- 1  
2  
3 Ter Schure, S., Junge, C., & Boersma, P. (2016). Discriminating non-native vowels on the basis of  
4 multimodal, auditory or visual information: Effects on infants' looking patterns and  
5 discrimination. *Frontiers in Psychology*, 7(APR), 1–11.  
6  
7 Thiessen, E. D. (2011). When Variability Matters More Than Meaning: The Effect of Lexical Forms  
8 on Use of Phonemic Contrasts. *Developmental Psychology*, 47(5), 1448–1458.  
9  
10 Triesch, J., Teuscher, c., Deak, G. O., & Carlson, E. (2006). Gaze following: Why (not) learn it?  
11 *Developmental Science*, 9, 125-147.  
12  
13 Tomasello, M. (1995). Joint attention as social cognition. In C. M. & P. Dunham (Ed.), *Joint*  
14 *attention: Its origins and role in development*. Lawrence Erlbaum Associates, Inc. (pp.  
15 103–130).  
16  
17 Walley, A. C. (2008). Speech Perception in Childhood. In D. B. Pisoni & R. E. Remez (Eds.), *The*  
18 *Handbook of Speech Perception*. Blackwell Publishing (pp. 1–37).  
19  
20 Weikum, W. M., Vouloumanos, A., Navarra, J., Soto-Faraco, S., Sebastián-Gallés, N., & Werker,  
21 J. F. (2007). Visual Language Discrimination in Infancy. *Science*, 316, 1159.  
22  
23 Wu, R., Gopnik, A., Richardson, D. C., & Kirkham, N. Z. (2011). Infants learn about objects from  
24 statistics and people. *Developmental Psychology*, 47, 1220–1229.  
25  
26 Wu, Z., & Gros-Louis, J. (2014). Infants' prelinguistic communicative acts and maternal  
27 responses: Relations to linguistic development. *First Language*, 34(1), 72–90.  
28  
29 Wu, R., & Kirkham, N. Z. (2010). No two cues are alike: Depth of learning during infancy is  
30 dependent on what orients attention. *Journal of Experimental Child Psychology*, 107, 118–  
31 136.  
32  
33 Wu, R., Tummeltshammer, K. S., Gliga, T., & Kirkham, N. Z. (2014). Ostensive signals support  
34 learning from novel attention cues during infancy. *Frontiers in Psychology*, 5, 251.  
35  
36 Yeung, H. H., Chen, L. M., & Werker, J. F. (2014). Referential Labeling Can Facilitate Phonetic  
37 Learning in Infancy. *Child Development*, 85(3), 1036–1049.  
38  
39 Yeung, H. H., & Nazzi, T. (2014). Object labeling influences infant phonetic learning and  
40 generalization. *Cognition*, 132(2), 151-163.  
41  
42 Yeung, H. H., & Werker, J. F. (2009). Learning words' sounds before learning how words sound:  
43 9-Month-olds use distinct objects as cues to categorize speech information. *Cognition*,  
44 113, 234-243.  
45  
46 Yoon, J. M. D., Johnson, M. H., & Csibra, G. (2008). Communication-induced memory biases in  
47 preverbal infants. *Proceedings of the National Academy of Sciences, USA*, 105, 13690–  
48 13695.  
49  
50 Yu, C., & Smith, L. B. (2012). Embodied attention and word learning by toddlers. *Cognition*,  
51 125(2), 244–262.  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 **Appendix 1.** List of test trials  
4

5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Contrast	Items	IPA transcription
/b/-/v/	billy-villy baggy-vaggy benny-venny boddy-voddy	/'bɪli/-/'vɪli/ /'bægi/-/'vægi/ /'bɛni/-/'vɛni/ /'bɒdi/-/'vɒdi/
/i:-/i/	leanny-linny deaddy-diddy teaggy-tiggy seabby-sibby	/'li:ni/-/'lɪni/ /'di:di/-/'dɪdi/ /'ti:gi/-/'tɪgi/ /'si:bi/-/'sɪbi/
SW-WS	<u>crocodile-crocodile</u> <u>penguin-penguin</u> <u>dolphin-dolphin</u> <u>elephant-elephant</u>	/'krɒkədail/-/'krɒkə'daɪl/ /'peŋgwɪn/-/'peŋ'gwɪn/ /'dɒlfɪn/-/'dɒl'fɪn/ /'eləfənt/-/'elə'fənt/

## Tables

*Table 1.* Mean pitch range values in Hertz (SD in parentheses) of the critical word across visual conditions and for all carrier sentences within a training trial

	<i>'socially-engaging'</i>	<i>'visual-speech'</i>	<i>'ostensive-cueing'</i>	<i>Anova results</i>
<i>Look, it's a CW!</i>	246.9 (34)	249.4 (30)	225.2 (32)	F(2,42)=2.78, p = n.s., $\eta^2$ = .117
<i>Look, a CW!</i>	253.7 (33)	228.5 (34)	230.1 (33)	F(2,42)=3.3132, p = n.s., $\eta^2$ =.130
<i>CW is nice!</i>	134.7 (31)	146.7 (24)	128.8 (28)	F(2,42)=2.210, p = n.s., $\eta^2$ =.095
<i>Nice CW!</i>	38.8 (12)	45.1 (14)	45.1 (13)	F(2,42)=1.143, p = n.s., $\eta^2$ =.052
<i>Hey CW!</i>	12.2 (8)	11.5 (7)	14.8 (7)	F(2,42)=.837, p = n.s., $\eta^2$ =.038
<i>CW!</i>	192.5 (51)	192.5 (47)	210.5 (50)	F(2,42)=.714, p = n.s., $\eta^2$ =.033

For Review Only

*Table 2.* Summary of the various trained conditions across blocks and lists.

	List 1	List 2	List 3
Block 1	/b-v/ + 'socially-engaging'	SW-WS + 'visual-speech'	/i:-l/ + 'ostensive-cueing'
Block 2	/i:-l/ + 'visual-speech'	/b-v/ + 'ostensive-cueing'	SW-WS + 'socially-engaging'
Block 3	SW-WS + 'ostensive-cueing'	/i:-l/ + 'socially-engaging'	/b-v/ + 'visual-speech'

For Review Only

Table 3. Coefficient effects of all main effects and 2-way interactions of the three models exploring the children's looking patterns across visual conditions, phonological contrasts, and age. Model 1 explores amount of time looking at the mouth vs. at other Aol; Model 2 explores amount of time looking at the object of reference vs. other Aol; Model 3 explores looking shifts between mouth and object of reference vs. 'static' gazes to mouth, object, or eyes. \*\*\*  $p < .001$ , \*\*  $p < .01$ , \*  $p < .05$

	Model 1			Model 2			Model 3		
	$\beta$	SE	t	$\beta$	SE	t	$\beta$	SE	t
<b>Main effect of Phonological Contrast</b>									
Consonant (Intercept)	2.67	0.43	6.19***	0.45	0.12	3.77***	0.09	0.02	5.70***
Vowel	1.32	0.54	2.43*	-0.03	0.12	-0.22	0.00	0.02	0.14
Stress	-0.46	0.50	-0.90	0.03	0.12	-0.27	-0.03	0.02	-1.53
<b>Main effect of Visual Condition</b>									
Socially-engaging (Intercept)	3.56	0.46	7.71***	0.17	0.04	4.46***	0.06	0.02	3.65***
Ostensive-cueing	-1.29	0.55	-2.35*	0.13	0.05	2.77*	0.03	0.02	1.41
Visual speech	-0.48	0.55	-0.87	0.07	0.05	1.51	0.04	0.02	1.61
<b>Main effect of Age</b>									
Age 4 (Intercept)	3.30	0.51	6.51***	0.45	0.09	4.70***	0.07	0.01	4.97***
Age 5	-0.30	0.69	-0.43	-0.19	0.13	-1.41	0.02	0.02	1.23
<b>Interaction between Phonological Contrast and Visual Condition</b>									
Consonant:socially-engaging (Intercept)	3.84	0.76	5.07***	0.08	0.15	0.52	0.06	0.03	2.22*
Consonant:ostensive-cueing	-1.89	1.11	-1.71	0.38	0.23	1.67	0.04	0.04	1.01
Consonant:visual-speech	-1.79	1.11	1.61	0.47	0.23	2.05*	0.07	0.04	1.68
Vowel:socially-engaging	-0.01	1.11	-0.01	0.13	0.23	0.56	0.02	0.04	0.51
Vowel:ostensive-cueing	1.12	1.72	0.64	-0.41	0.35	-1.16	-0.05	0.06	-0.89
Vowel:visual-speech	2.90	1.69	1.72	-0.12	0.34	-0.36	-0.01	0.06	-0.10

Stress:socially-engaging	-0.34	1.11	-0.31	0.15	0.23	0.66	-0.02	0.04	-0.49
Stress:ostensive-cueing	0.71	1.69	0.42	0.05	0.34	0.15	-0.03	0.06	-0.60
Stress:visual-speech	0.81	1.72	0.46	-0.48	0.35	-1.37	-0.02	0.06	-0.34

---

**Interaction between Age and Phonological Contrast**

Age 4:consonant (Intercept)	3.04	0.68	4.44***	0.36	0.14	2.54*	0.06	0.02	2.48*
Age 4:vowel	0.92	0.79	1.16	0.12	0.18	0.67	0.02	0.03	0.63
Age 4:stress	-0.14	0.79	-0.18	0.16	0.18	0.91	0.01	0.03	0.47
Age 5:consonant	-0.69	0.94	-0.74	-0.01	0.19	-0.06	0.07	0.03	2.06*
Age 5:vowel	0.75	1.09	0.68	-0.28	0.24	-1.14	-0.03	0.04	-0.73
Age 5:stress	0.44	1.09	0.40	-0.24	0.24	-1.00	-0.10	0.04	-2.12*

---

**Interaction between Age and Visual Condition**

Age 4:socially-engaging (Intercept)	3.73	0.69	5.44***	0.12	0.23	0.52	0.07	0.02	2.74**
Age 4:ostensive-cueing	-0.96	0.80	-1.19	0.36	0.31	1.16	0.00	0.03	0.04
Age 4:visual-speech	-0.35	0.80	-0.43	0.33	0.31	1.16	0.01	0.03	0.43
Age 5:socially-eng.	-0.01	0.94	-0.00	-0.06	0.29	-0.21	-0.00	0.03	-0.28
Age 5:ostensive-cueing	-0.63	1.10	-0.57	-0.02	0.43	-0.04	0.07	0.05	1.47
Age 5:visual-speech	-0.25	1.10	-0.22	-0.24	0.24	-1.02	0.03	0.05	0.73

---

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Figures

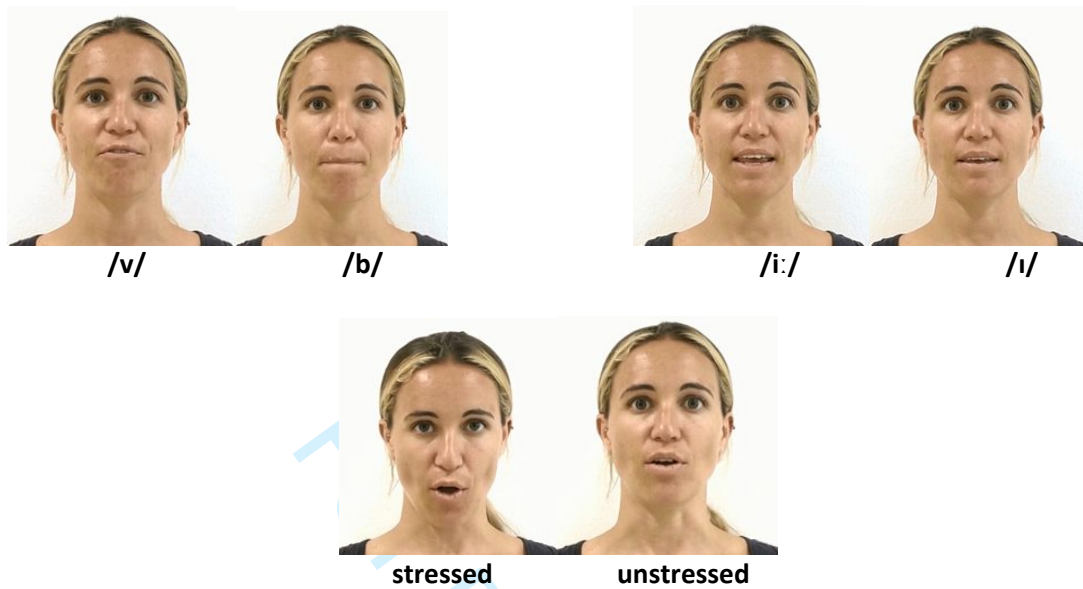


Figure 1. Video frame of the point of maximal visual differentiation for each contrast. On the top left, lip configuration at the onset of consonants /v/ and /b/ in the minimal pair *venny-benny*. On the top right, lip configuration in the middle of the vowels /i:/ and /ɪ/ in the minimal pair *seabby-sibby*. On the bottom, head posture during the production of the 'dol-' syllable in a stressed and unstressed context in the minimal pair *dolphin-dolphin*.



Figure 2. Example of images depicting the meaning of target words of a minimal pair in the lexical stress training trials in which cognates were used. On the left, a drawing depicting the meaning of the real word of a minimal pair (e.g. a dolphin). On the right, a drawing depicting the meaning of the counterpart (non-real) word in the minimal pair (e.g. a dolphin).

For Review Only



Figure 3. Visual display of the three visual conditions during the training trials. Left panel, example of the 'socially-engaging' condition; middle panel, example of the 'visual-speech' condition; right panel, example of the 'ostensive-cueing' condition. The dashed arrows in the left panel indicate a dynamic movement by which the speaker alternated her gaze between the object and the observer and did not appear in the real display.

For Review Only



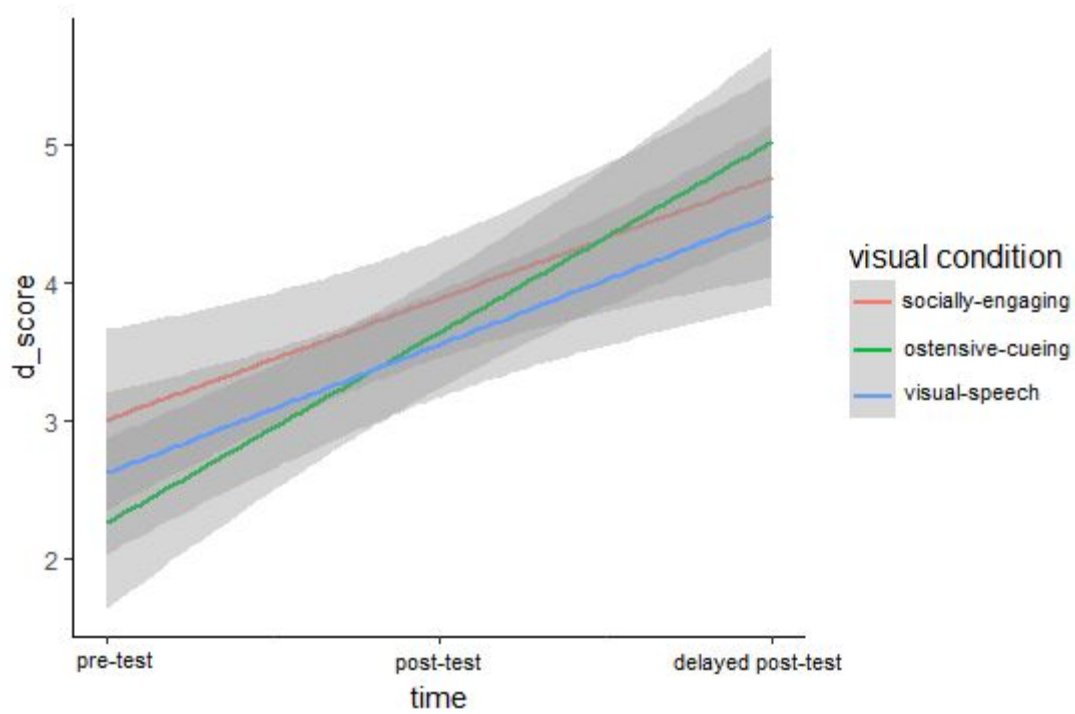


Figure 4. Children's accuracy in distinguishing the phonological contrasts (as measured by  $d'$  scores), across the three testing sessions and as a function of the social interaction in which they were trained.

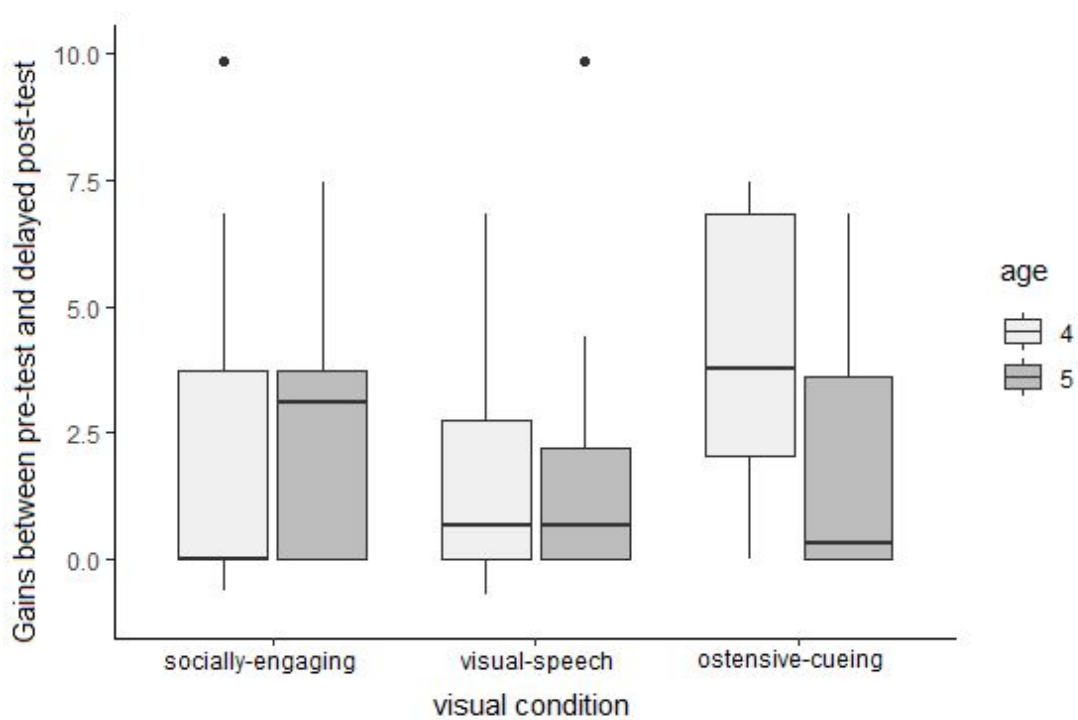


Figure 5. Boxplots depicting gains between pre-test and delayed post-test (measured by  $d'$  scores), as a function of the three distinct visual conditions and of the children's age.

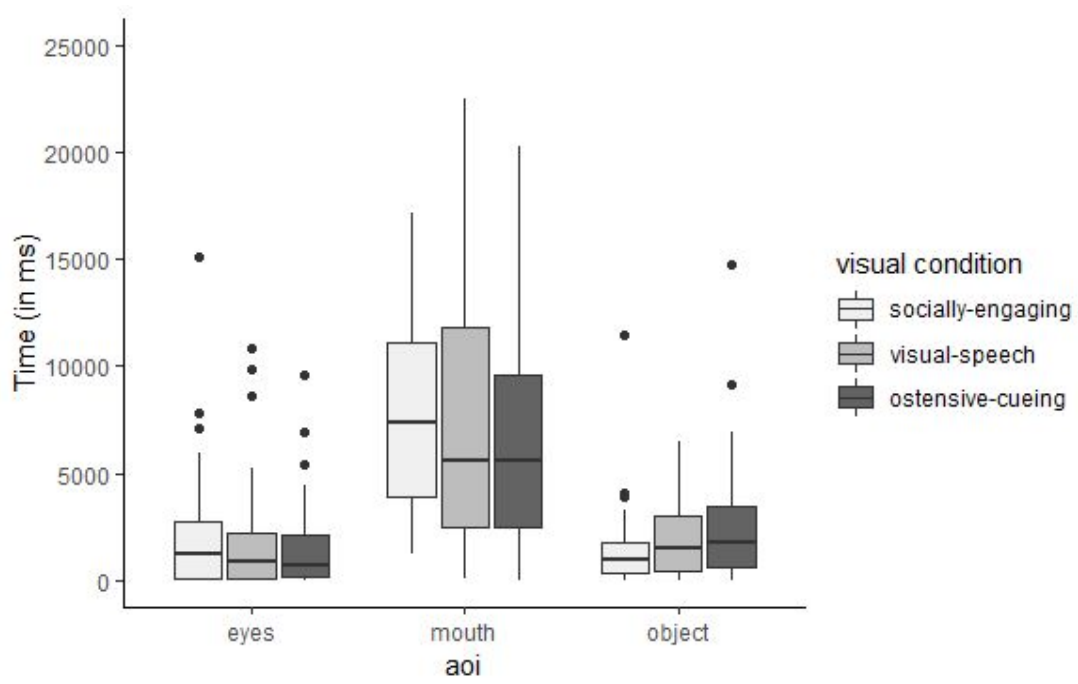


Figure 6. Boxplots displaying the amount of time (in milliseconds) children spent looking to each AOI across the three distinct trained visual conditions.

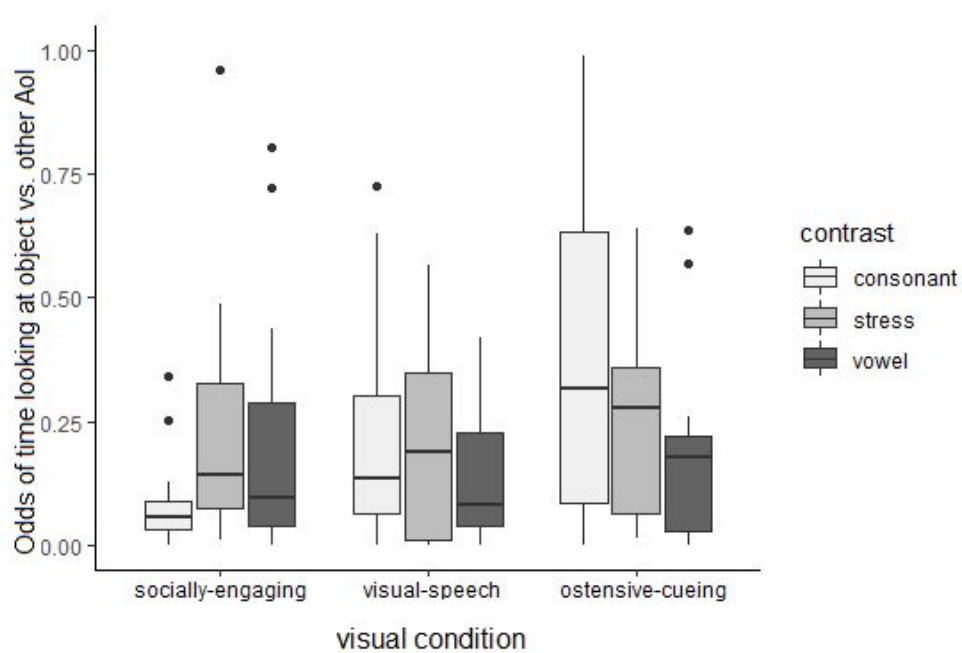


Figure 7. Box plots representing the odds ratio of time looking at the object of reference (vs. other Aol), across visual conditions and phonological contrasts.