1

2

3

4 **Reconciling the contradictory effects of production on word learning:**

5 **Production may help at first, but it hurts later**

6

7

8 Efthymia C. Kapnoula

9 Basque Center on Cognition, Brain and Language; Ikerbasque

10

11 and

12

13 Arthur G. Samuel

14 Basque Center on Cognition, Brain and Language; Stony Brook University; Ikerbasque

15

16 Running Head: PRODUCTION HELPS AT FIRST, BUT HURTS LATER

17
18
19
20
21
22
23 Corresponding Author:
24 Efthymia C. Kapnoula
25 Paseo Mikeletegi 69
26 Basque Center on Cognition, Brain and Language (BCBL)
27 20009 San Sebastián - Donostia
28 Spain
29 kapnoula@gmail.com
30

1                                   **Abstract**

2          Does saying a novel word help to recognize it later? Previous research on the effect of

3    production on this aspect of word learning is inconclusive, as both facilitatory and detrimental

4    effects of production are reported. In a set of three experiments, we sought to reconcile the

5    seemingly contrasting findings by disentangling the production from other effects. In Experiment

6    1, participants learned eight new words and their visual referents. On each trial, participants

7    heard a novel word twice: either (1) by hearing the same speaker produce it twice (*Perception-*

8    *Only condition*) or (2) by first hearing the speaker once and then producing it themselves

9    (*Production condition*). At test, participants saw two pictures while hearing a novel word and

10   were asked to choose its correct referent. Experiment 2 was identical to Experiment 1, except

11   that in the Perception-Only condition each word was spoken by two different speakers

12   (equalizing talker variability between conditions). Experiment 3 was identical to Experiment 2,

13   but at test words were spoken by a novel speaker to assess generalizability of the effect.

14   Accuracy, RT, and eye-movements to the target image were collected. Production had a

15   facilitatory effect during early stages of learning (after short training), but its effect became

16   detrimental after additional training. The results help to reconcile conflicting findings regarding

17   the role of production on word learning. This work is relevant to a wide range of research on

18   human learning in showing that the same factor may play a different role at different stages of

19   learning.

20   Keywords: word learning, spoken word recognition, production, mental lexicon, visual world
21   paradigm

1                                      **Introduction**

2        To learn a novel word, we need to integrate it into our mental lexicon. The trajectory of

3 lexical integration and the factors upon which it depends are hotly debated topics. Here we

4 examine the role of *production* (i.e., the action of producing a word out loud compared to just

5 hearing it) on word learning, using spoken word recognition as a proxy of lexical integration.

6        There is a widely held assumption that producing a novel word helps to build its lexical

7 representation, which can then support lexical functions, including its recognition. In accordance

8 with this assumption, when people learn a new word, they are often asked to repeat it. This is a

9 common practice, for example, in second language learning contexts, where instructors often ask

10 their students to repeat new words immediately after encountering them for the first time (Duff,

11 2000; Kadota, 2019). In line with this idea, current approaches to second language learning, such

12 as Communicative Language Teaching, Task-Based Language Teaching, and Communicative

13 Competence emphasize learners' immediate communication needs and encourage production

14 from the earliest moments of instruction (for a review see Lightbown & Spada, 2013).

15        The idea that production helps word learning is in line with a substantial body of research

16 showing that production enhances memory. First reported by Hopkins and Edwards (1972), the

17 key finding is that material read aloud is better remembered. This facilitatory *production effect*

18 has been documented many times. For example, Gathercole and Conway (1988) conducted a

19 series of experiments in which adults were presented with a set of words (one by one) and,

20 depending on the experiment, were asked to read a word silently, read it out loud, mouth it, read

21 it and hear it, only hear it, or write it (seeing it or not). At test, participants had to indicate

22 whether a given word was new or old. Across the experiments, performance at test was better for

23 words that were read out loud compared to all other conditions (see also Dodson & Schacter,

1  2001; Gathercole & Conway, 1988; P. MacDonald & MacLeod, 1998; MacLeod & Bodner,

2  2017).

3      Given the robustness of this phenomenon, it seems intuitively reasonable to assume that

4  similar facilitation should apply to novel words; however, research on this topic has been

5  inconclusive. There is some evidence that production facilitates word learning, but there is also

6  evidence that production can actually impair word learning. The difficulty in assessing the role

7  of production in word learning stems mainly from the fact that very few studies have tried to

8  isolate the effect of production independently of other effects, such as the testing effect

9  (Karpicke & Roediger, 2008), the finding that recall from memory benefits the retention of novel

10  information.

11      To date, studies that have looked more closely at the effect of production have reported

12  contradictory results. For example, Zamuner et al. (2016) report evidence from eye-movements

13  that production has a facilitatory effect on novel word learning. However, Zamuner et al. (2018,

14  using the same task, but with children rather than adults) reported the reverse pattern, i.e., a

15  detrimental effect of production. Similarly, Leach and Samuel (2007) also showed that

16  production during training leads to weaker effects of lexical engagement, meaning that produced

17  words were not as well integrated into the mental lexicon, compared to words that were only

18  heard during training.

19      The present study investigates the source(s) of these discrepancies to shed light on the

20  role of production on the learning of novel spoken words, as reflected by their recognition. Given

21  the existence of both positive and negative consequences of production, we begin by briefly

22  presenting a number of mechanisms via which production may facilitate or impair word learning.

23  Our goal was to identify possible confounds and isolate the direct effect(s) of production per se.

1  **Reasons why production may help word learning**

2  As mentioned above, even though there is evidence consistent with a facilitatory role of

3  production in word learning, it is difficult to assess whether production per se is the critical

4  factor behind these effects and, more importantly, to pinpoint the exact mechanism. Work on the

5  production effect on familiar words has generated support for a *distinctiveness* account,

6  according to which, the act of producing a word enhances its distinctiveness and thus strengthens

7  the corresponding mnemonic trace (Gathercole & Conway, 1988; MacLeod & Bodner, 2017;

8  Ozubko et al., 2014; Ozubko & Macleod, 2010). However, learning a novel word involves more

9  than just generating a mnemonic trace; it requires creating an entirely novel lexical

10  representation and integrating it into the mental lexicon. Thus, it is unknown whether and how

11  production may help different aspects of word learning. Below we describe a few different

12  mechanisms via which this facilitation may occur.

13  First, production may boost word learning via the creation of *articulatory*

14  *representations*. According to the model of speech production proposed by Hickok and

15  colleagues (Hickok, 2012, 2014; Hickok et al., 2011), phonological representations act as hubs

16  that are used to map sound to meaning and to speech articulation. Even though the two mappings

17  are computationally distinguishable (Nora et al., 2015), they can both be viewed as components

18  of lexical representations. Thus, learning and practicing the articulatory sequence that

19  corresponds to a novel word may serve as an additional dimension of the newly formed lexical

20  representation, which can then be used to bootstrap further integration of the word into the

21  mental lexicon. An articulatory locus for the facilitatory effect of production would also be

22  consistent with work on second-language acquisition showing a learning advantage for overt but

1    not covert repetition (Mattys & Baddeley, 2019). This mechanism could be considered a

2    facilitatory effect of production per se.

3        Second, it has been proposed that word learning is tightly linked to and supported by

4    *phonological short-term memory* (PSTM; for review, see Baddeley, Papagno, & Vallar, 1988;

5    Gathercole, 2006). The general idea is that in order to produce a word, the corresponding

6    phonological sequence is briefly maintained in PSTM, which gradually leads to longer-term

7    learning. However, the details of this mechanism vary across different models (for a review see

8    Thorn & Page, 2008). For example, according to the primacy model, proposed by Page and

9    Norris (Page & Norris, 1998, 2008, 2009), phonological word-form learning is a more

10   naturalistic version of the Hebb repetition effect (Hebb, 1961; which refers to the finding that

11   immediate serial recall of a list of familiar items, such as digits, gradually improves over

12   multiple repetitions). On the other hand, Gupta's model (Gupta, 2003, 2008; Gupta & Tisdale,

13   2009) involves a short-term sequencing mechanism that builds associations between sublexical

14   sequences and patterns of lexical-level activation. In that respect, Gupta's model makes a

15   distinction between sublexical sequences and lexical word-forms. In this case the facilitatory

16   effect would again be inherently linked to production.

17       Third, a broader mechanism that may be responsible for the facilitatory effects of

18   production is *attention*. The idea that language production requires higher levels of attention is

19   not only intuitive, but also supported by data (e.g., see Boiteau, Malone, Peters, and Almor,

20   2014). Despite this, few experiments have controlled for this factor in assessing the role of

21   production in word learning. For example, in the Gathercole and Conway (1988) experiments

22   mentioned above, words were either preceded by the critical instruction (e.g., "say aloud"), or

23   they were only presented in one experimental condition for a given participant. Thus, in both

1    cases participants knew if they were required to produce the word ahead of time. Knowing that

2    they will need to say the target word out loud may have led to increased attention – e.g., towards

3    the phonological structure of the word, its acoustic implementation, its possible semantic

4    associations, or any combination of these – which could explain the robust facilitatory effect of

5    production during testing. Thus, any facilitation caused by attention would be an indirect

6    consequence of production.

7         Fourth, often (but not always[1]) production involves recall from memory. This is the case,

8    for example, when a novel word-form is linked to a visual referent, which is subsequently (at

9    training and/or at testing) used to prompt the production of its newly learned label. Retrieval

10   practice is known to lead to better retention of information (i.e., the *testing effect*; Roediger &

11   Karpicke, 2006). In support of this hypothesis, Karpicke and Roediger (2008) trained English-

12   speaking adults in Swahili-English word pairs using training regimes that differed in whether

13   they involved repeated testing (involving recall) versus repeated studying (not involving recall).

14   In contrast to repeated studying, which had no effect on delayed recall, repeated testing had a

15   large facilitative effect. Similarly to the previous case, this kind of facilitation would be viewed

16   as a by-product of production.

17        Lastly, there is an additional way in which production can differ from passive exposure.

18   In training paradigms that involve auditory presentation of the target words, when participants

19   are asked to produce a word themselves, they also hear it in a new voice (their own). In these

20   cases, production is confounded with increased *speaker variability*. When learning new words,

21   listeners encode voice-related information (Creel & Tumlin, 2011; Houston & Jusczyk, 2000;

22   Kapnoula & Samuel, 2019), which is why variability in this dimension may affect word learning.

---

[1] When production immediately follows presentation of the target word, no retrieval from long-term memory is necessary.

1    Indeed, there is evidence in favor of a beneficial role of talker variability in word learning. For

2    example, Rost and McMurray (2009) showed that when infants learn similar words (like *buk* and

3    *puk*) spoken by multiple speakers, they later discriminate between them better than when they

4    have only heard the words spoken by one person (see also Höhle et al., 2020). Similarly,

5    Richtsmeier et al (2009) exposed 4-year-olds to novel words that were spoken either by one or

6    10 different talkers. At test, children were faster and made fewer errors in producing the words

7    that had been spoken by many talkers. These findings suggest that hearing a novel word spoken

8    by different talkers can lead to a more robust and abstract lexical representation. A proposed

9    mechanism behind this effect is that increased variability in irrelevant dimensions helps the

10   listener to identify the relevant dimensions (Rost & McMurray, 2010; Singh, 2008). In this case,

11   facilitation would be caused by the additional variability that comes with production, rather than

12   production itself.

13      Based on the points outlined above, it is clear that including a production requirement in

14   a training regime may lead to better word learning, but, depending on the details of the

15   procedure, this may be due to a number of different mechanisms. Thus, in order to fully

16   understand the role of production in word learning, we need to use experimental designs that take

17   these points into consideration.

18   **Reasons why production may hinder word learning**

19      Although counterintuitive, it is also theoretically possible that production could impede

20   word learning. Again, there are a number of different ways in which this could happen. First,

21   production may interfere with the *encoding* of a novel word-form at the earliest moments of

22   learning. By encoding we refer to learning the sound pattern of a word-form, which can be

23   viewed as the bare minimum amount of information that is necessary to recognize a word. This

1    maps onto what Leach and Samuel (2007) refer to as *lexical configuration* and it is thought to

2    correspond to early stages of word learning (i.e., when the word is first added into the mental

3    lexicon, after a handful of exposures to its spoken form). That is, immediately after hearing the

4    new word, listeners may benefit from having a moment in which no further input (or output) is

5    processed. During this period, the system can make the necessary adjustments (e.g., adjust the

6    connection weights between speech sounds and lexical levels) that correspond to the successful

7    encoding of the novel word-form. Specifically, production may hinder this process by blocking

8    access to the echoic trace of the stimulus. This may in turn impede encoding directly, by taking

9    away the input of encoding, and/or indirectly, by taking away the input of sub-vocal rehearsal,

10    which would be expected to boost encoding. In these ways, immediate production may interfere

11    with the encoding process, yielding non-optimal learning outcomes.

12          Second, and relatedly, if learners are simultaneously dealing with the need to learn a

13    word perceptually and to learn how to produce it, any mismatch between the memory

14    representations and/or cognitive processes needed for these two tasks can lead to *interference*

15    between them. Given that perception is based on auditory codes and production is based on

16    motor codes, there is inherently some mismatch. The fact that both codes need to refer to the

17    same object at some level means that they may be particularly vulnerable to a form of lateral

18    inhibition. The idea that production may interfere with perception during early stages of learning

19    is not new. Krashen's Input Hypothesis (1985), for example, makes this point in the context of

20    second language learning. According to this hypothesis, when learning a new language one

21    should not rush into producing new words during the very first stages of language learning.

22    Instead, production should follow after a "silent period" has passed. Even though this is based on

23    developmental observations of L1 acquisition, the idea is applicable to adult language learning.

1          Third, production may interfere with word learning at a later stage, when the word-form

2     is being linked to its *semantic referent*. Here, the underlying assumption is that a new word-form

3     is first encoded (in some form of proto-lexical representation) and then mapped to a referent

4     (Fernandes et al., 2009; Rodriguez-Fornells et al., 2009, but see François et al., 2017, for

5     evidence that they can also happen in parallel). Indeed, a number of studies have shown that

6     mapping to meaning is facilitated when it follows speech segmentation (Graf Estes et al., 2007;

7     Hay et al., 2011; Mirman, Magnuson, et al., 2008). The rationale and corresponding processes of

8     how production interferes with word learning would be very similar to the kind of interference

9     described above (during encoding), but in this case production would overlap with the

10    (subsequent) mapping of the novel word to its meaning.

11        Fourth, it might also be the case that production exposes learners to *poor input*, which

12    leads to poorer learning. That is, the input that is presented to the participants during training is

13    usually comprised of clean, carefully manipulated, high-quality stimuli. In contrast, the nature of

14    a learner's own production is out of the experimenter's control and can thus vary substantially in

15    terms of quality. Noisy output can act as input, for example when the participant is asked to read

16    a word, or hear and repeat it. As a result of being exposed to noisy input, learning may be

17    negatively affected.

18        Indeed, the idea that the quality of the input can affect processing of spoken language is

19    intuitive and supported by empirical findings. Relevant work has mostly looked at *clear speech*,

20    which is slower and hyperarticulated compared to plain speech (Bradlow et al., 1996; Bradlow &

21    Hayes, 2003). Speakers typically adopt this speaking style when the listener is thought to face

22    communication-related difficulties, e.g., nonnative language or hearing impairment (Smiljanić &

23    Bradlow, 2009), and indeed clear speech input appears to facilitate comprehension (e.g., see

1 Payton et al., 1994). More pertinent to our study, Riley and McGregor (2012) examined the

2 effects of speaking style on children's word learning. They found that new words heard in clear

3 speech were later produced more accurately; however no effect of speaking style was found for

4 perception (see also Baese-Berk & Samuel, 2016, for a discussion of the role of input quality in

5 perceptual learning).

6 **Previous research on the role of production in word learning is inconclusive**

7 It should not be surprising that studies on the role of production on word learning have

8 yielded contradictory findings; given the variety of mechanisms in which production may affect

9 word learning, seemingly small differences between experimental designs and procedures may

10 lead to large differences in (or even reversal of) the obtained pattern of results.

11 Leach and Samuel (2007) evaluated how a number of different factors affected novel

12 word learning. In this work, the authors focused on two lexical properties: *configuration* and

13 *engagement*. In the current context, as mentioned above, lexical configuration corresponds to

14 building a phonological representation. In contrast, lexical engagement refers to the ways in

15 which a word interacts with other representations (e.g., inhibiting other words, or boosting the

16 activation of speech sound representations). The latter property is taken as a stronger marker of

17 word learning, as it reflects deeper integration of a novel item into the mental lexicon.

18 In the Leach and Samuel (2007) study, participants were either trained with a phoneme

19 monitoring task (Exp.1) or a word-picture-association task (Exp.2), coupled with a production

20 requirement (Exp.4) or not (Exp.5). After training, lexical configuration and lexical engagement

21 were assessed separately. Lexical configuration was assessed in a three alternative recognition

22 judgment (in which participants would hear a newly learned word along with two similar lures

23 and had to choose which of the three items they had just learned) and a word-in-noise task (in

1  which participants had to recognize the critical items buried in progressively lower levels of

2  white noise). To assess lexical engagement the authors measured the ability of the new items to

3  drive phonemic restoration (Samuel, 1996; Warren, 1970; the finding that when part of a word is

4  missing or replaced by a different sound, listeners still report hearing it) and perceptual learning

5  (Norris, McQueen, & Cutler, 2003; the finding that repeated exposure to an ambiguous sound

6  embedded in real words changes the way listeners identify this sound in a later task).

7  Interestingly, the results revealed a dissociation: production boosted lexical configuration, but

8  hindered lexical engagement.

9  　　　　Hopman and Macdonald (2018) also looked at the role of production, but unlike the

10  Leach and Samuel study, their production task required participants to recall the critical

11  information (i.e., the newly learned words). Word learning was assessed via a vocabulary test,

12  which required the comprehension of individual words within a phrase context. Their results

13  revealed a facilitatory effect of production on word learning. However, it remains unclear

14  whether this effect was due to production per se, or can, for example, be attributed to repeated

15  retrieval (i.e., the testing effect). Additionally, it could be argued that their measure of word

16  learning assessed speed of phrase comprehension, rather than word learning per se.

17  　　　　Finally, Zamuner et al. (2016) looked at the role of production in word learning using

18  eye-tracking. During training each new word was presented along with its visual referent. For

19  half of the items, participants heard the new word twice (Heard-Only condition) and for the other

20  half they heard it once and were required to repeat it themselves (Produced condition), thus

21  equalizing the number of times each item was presented in each condition. At test, participants

22  heard each word and had to select its correct referent, given two options. Eye-movements during

23  testing were analyzed using growth curve analysis (GCA; Mirman, Dixon, & Magnuson, 2008),

1 which revealed a significant difference in the shape of the looking curves (quadratic term)

2 between conditions. The authors interpreted this difference as evidence for a facilitatory effect of

3 production on word learning.

4 A possible concern is that production was confounded with speaker variability; Heard-

5 Only items were heard twice by the same speaker (and the same recorded token), whereas

6 Produced items were heard by two speakers (the voice played to them and the participant

7 themselves). Thus, it is unclear whether the significant difference between experimental

8 conditions was driven by production or by input variability. In addition, training was limited to

9 two trials per item (which is much lower than the number of training repetitions typically used in

10 word learning studies, e.g., Gaskell & Dumay, 2003; Kapnoula et al., 2015; Leach & Samuel,

11 2007). This leaves open the possibility that this effect only appears at very early stages of word

12 learning, which may not reflect integration of the new items into the mental lexicon.

13 **Present study**

14 The main goal of the present study is to examine the effect of production on word

15 learning independently of other commonly confounding factors (talker variability, the testing

16 effect, attention, etc.), in order to offer an account that reconciles previous results. In addition, to

17 achieve a more comprehensive understanding of the effect, we looked at whether/how this effect

18 is modulated by the amount of training.

19 We adopted Zamuner et al.'s (2016) approach, which combines a number of strengths,

20 such as: 1) the number of presentations is equal across conditions, 2) retrieval is not required for

21 production (i.e., production is not confounded with retrieval practice), 3) participants are not

22 instructed about the mode of response ahead of time (i.e., minimizing differences in attention),

23 and 4) eye-movement data can be used as a proxy of lexical activation, allowing us to track

1  lexical activation in real time (Allopenna et al., 1998; Salverda & Tanenhaus, 2017). In addition,

2  by adopting this design we can more directly compare our results to those of the original study.

3  We conducted three experiments with the same general structure/design, each one

4  focusing on a different question. Experiment 1 aimed at replicating Zamuner et al. and

5  examining whether/how the results change as a function of the amount of training. To test this,

6  we added further training and testing after what corresponded to the end of the Zamuner et al.

7  experiment. Experiment 2 examined the effect of production on word learning while controlling

8  for speaker variability. To achieve this, we introduced speaker variability[2] in the Perception-Only

9  condition by playing each word in two different voices (thus matching the variability present in

10  the Produced condition). In Experiment 3, we used a novel voice at test to examine whether word

11  learning with/without production generalizes differently to novel speakers (e.g., whether

12  production leads to better generalization to novel speakers).

13  In all three experiments, we used the visual word paradigm (VWP) to track activation of

14  the target word in real time. In this paradigm, the underlying hypothesis is that the probability of

15  looking at an object increases as a function of the activation of the corresponding lexical item.

16  Based on this linking hypothesis, fixation proportions over time can be used as a direct index of

17  lexical activation (Allopenna et al., 1998; Salverda & Tanenhaus, 2017; Tanenhaus et al., 1995;

18  see also Magnuson, 2019 for a review of alternative hypotheses). Across experiments, the effect

19  of Production is defined as the difference between the Production and Perception-Only

20  conditions, since our question was how producing a new word affects learning compared to just

21  hearing it.

---

[2] Note that in other studies the typical number of talkers in high-variability conditions is much higher than two; however, in contrast to those studies, our goal was to control for talker variability, rather than test its effect.

1                                          **Experiment 1**

2          Experiment 1 is intended to replicate the facilitatory effect of production on word

3     learning reported by Zamuner et al. (2016). Note that Zamuner et al. (2016) used two training

4     trials per item, which is much lower than the typical amount of training used in word learning

5     studies (11-24 trials; e.g., see Gaskell & Dumay, 2003; Kapnoula et al., 2015; Kapnoula &

6     McMurray, 2016; Leach & Samuel, 2007). Thus, it is possible that the facilitatory effect of

7     production applies specifically to early stages of lexical acquisition (e.g., lexical configuration, if

8     we adopt the terminology proposed by Leach and Samuel). To test this, we asked whether this

9     effect is modulated by the amount of training.

10    **Method**

11    *Participants*

12         Forty (31 females; mean age = 25.8 years) native speakers of Spanish participated in

13    Experiment 1. Power analyses were conducted on data from three previous eye-tracking

14    experiments (reported in Kapnoula et al., 2015; Kapnoula & McMurray, 2016) that used a

15    different within-subject manipulation. Given the absence (to our knowledge) of a well-tested

16    method of sample size estimation for curve-fitting analyses (which was our primary analytical

17    approach), we conducted analyses for repeated-measures, within-subjects ANOVA (which was

18    our secondary analytical approach). These analyses indicated that a power of .95 requires a

19    sample size of 31 to 41. All analyses were conducted in G*Power (Faul et al., 2009, 2007).

20         Most participants were also fluent in Basque, which was foreseen and taken into account

21    in selecting the stimuli (see *Materials* below). All participants had normal/corrected-to-normal

22    vision and no known hearing or neurological impairments. Participants underwent informed

1   consent and were remunerated for their participation. All experimental procedures were

2   approved by the BCBL ethics committee.

3   *Design*

4        Experiment 1 had two Phases, each consisting of one training and one testing block (see

5   Table 1). Phase 1 of the training had the same number of training trials used by Zamuner et al.,

6   2016, while Phase 2 had five times that number of trials. Thus, there were 12 repetitions across

7   Phase 1 and Phase 2 training, which matches the typical amount of training used in previous

8   word learning studies (11-24 trials). In such previous studies, learners typically show asymptotic

9   performance after approximately 8-10 trials (e.g., Leach & Samuel, 2007; Samuel & Larraza,

10  2015). The testing blocks were identical across the two Phases and the test trials in each Phase

11  matched the number of test trials (32) in the Zamuner et al. study.

12
13  Table 1. *Number of training and testing trials per phase*

|          | Phase 1 | | Phase 2 | |
|----------|---------|--------|---------|--------|
| **Training** | 8 words × 2 repetitions = | 16 trials | 8 words × 10 repetitions = | 80 trials |
| **Testing**  | 8 words × 4 repetitions = | 32 trials | 8 words × 4 repetitions = | 32 trials |

14
15
16        *Training*. Participants learned eight novel words. Each word-form was used as the label

17  for an unfamiliar object. The correspondence between novel words and objects was randomized

18  across participants using a Latin Square.

19        Crucially, for each participant, half of the words were assigned to the Perception-Only

20  condition, and half to the Production condition (the assignment of each word to one condition or

21  the other was randomized across participants). In Perception-Only training trials, the picture of

22  one unfamiliar object was presented and the corresponding novel word was heard twice. In

1    Production training trials, the only difference was that the corresponding novel word was heard

2    once and then repeated by the participant.

3         In Phase 1, we used the same number of training trials used by Zamuner et al. (2016),

4    which was two trials per word (8 words × 2 repetitions = 16 training trials). In contrast, the

5    training block of Phase 2 consisted of 10 trials per word (8 words × 10 repetitions = 80 training

6    trials). Thus, cumulatively, there were 96 training trials across Phases.

7         *Testing*. Participants were tested twice, once at the end of each Phase. All testing trials

8    were identical across conditions and phases. In each testing trial, participants saw two of the

9    objects while hearing one of the novel words and had to select the picture that was the correct

10   referent of that word. Each target was repeated four times in each testing block, resulting in 64

11   testing trials (8 words × 4 repetitions × 2 testing blocks) across Phases.

12   **Materials**

13        All novel words were CVCV items: /bopa/, /ʧofa/, /deɾa/, /guθa/, /kiða/, /reka/, /tuma/,

14   and /jata/. The items were checked by a Spanish-Basque bilingual research assistant to make sure

15   that they were nonwords in both Spanish and Basque, but morphologically consistent with

16   Spanish. Spoken stimuli were recorded by a native female speaker of Spanish in a sound-

17   attenuated room, sampling at 44,100Hz. We collected multiple recordings and chose one

18   recording per item based on sound quality. Chosen recordings were cut, cleaned (background

19   noise and occasional click/pop sounds removed), and intensity-scaled. Finally, 50 ms of silence

20   was added before and after each word. The average duration of the final stimuli (including the

21   100 ms of silence) was 714 ms.

22        Visual stimuli consisted of color pictures of eight unfamiliar objects. All images were
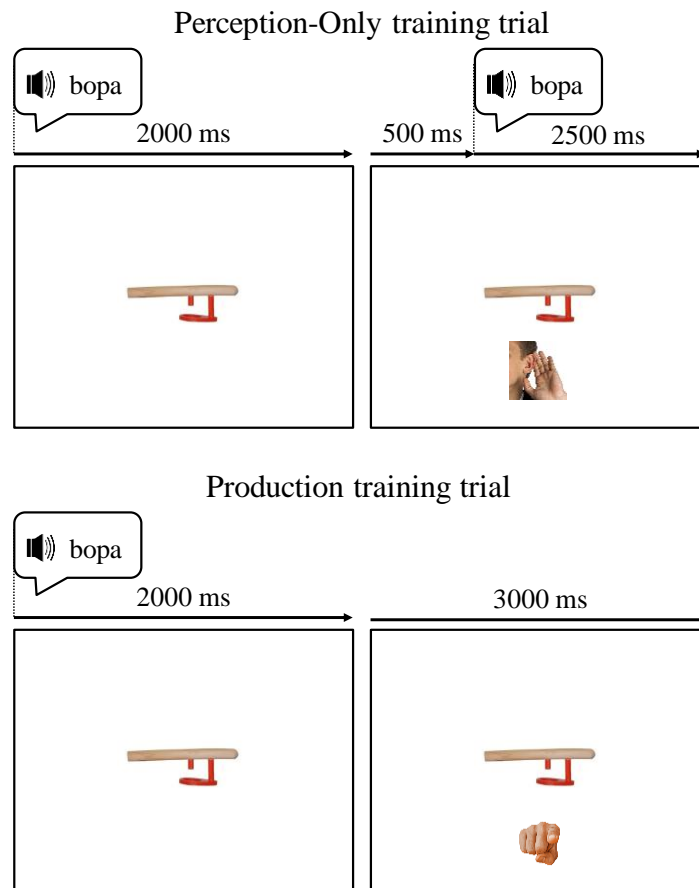
1    taken from Horst & Hout (2016). Images measured 300 × 300 pixels during presentation.

2    ***Procedure***

3         Participants were seated in front of a computer screen and were instructed that their task

4    would be to learn a set of new words and their meanings. They were then fitted with an SR

5    Research EyeLink 2K eye-tracker, a system with remote desktop mounting.

6         After calibration, participants were given instructions for the task and did a short training-

7    and-testing practice. The practice stimuli consisted of six images of fruits (piña [pineapple],

8    melon [melon], mango [mango], fresa [strawberry], uvas [grapes], and pera [pear]) and their

9    corresponding names, recorded by the same speaker as the experimental stimuli.



10

11         *Figure 1.* Visual description of training trial per condition

1    At the beginning of each training trial, one picture was presented at the center of the

2    screen of a 19" monitor operating at a resolution of 1204 × 768 pixels. Simultaneously, the

3    auditory label of the picture was played through high quality headphones. A prompt image

4    appeared below the image of the unfamiliar object 2,000 ms after the onset of the word. In

5    Perception-Only trials, the prompt image showed a hand next to an ear making the gesture for

6    "listen". When seeing this prompt, participants had to remain silent and listen to the word again.

7    The novel word was repeated by the computer 500 ms after the presentation of the prompt. In

8    production trials, the prompt image showed a finger pointing at the participant. When seeing this

9    prompt, participants had to repeat the word out loud into a microphone attached to their

10   headphones. The prompt images were presented and explained to the participants during the

11   initial instructions. For both conditions, the image of the unfamiliar object remained on the

12   screen for 3,000 ms after the presentation of the prompt (see Figure 1). Perception-Only and

13   Production trials were randomly intermixed. The first training Phase lasted approximately 2-3

14   minutes and the second one approximately 10 minutes.

15   At the beginning of each testing trial, pictures of two objects were presented in the two

16   horizontal ends of the screen, spaced 424 pixels apart. One picture was the target item for that

17   trial (i.e., the picture assigned to that word during training). The other picture was one of the

18   other seven images. Items were paired so that each target was always presented with the same

19   competitor. This was done to minimize further learning opportunities during testing[3]. The

20   position of the target was randomized across trials. Along with the presentation of the pictures, a

21   blue circle appeared at the center of the screen. After 500 ms, the circle turned red, cueing the

22   participant to click on it to start the trial. This allowed the participants to briefly look at the

---

[3] That is, if the item pairs were not fixed, participants would be able to figure out the word-picture mappings during testing as a result of cross-situational learning (Yu & Smith, 2007).

1   pictures before hearing anything, thus minimizing eye movements due to visual search (rather

2   than lexical processing). As soon as participants clicked on the red circle, it disappeared and an

3   auditory stimulus was played through the headphones. Participants then clicked on the picture

4   they believed to be the referent of the word. No feedback was provided during testing. Each

5   testing block lasted approximately 5 minutes.

6       Participants completed all four blocks within the same session, and were given a chance

7   to take a break every 16 trials (both for the training and testing Phases). There was no time limit

8   on the trials; however, participants typically responded in less than 2 sec (M = 1,176 ms, SD =

9   256 ms).

10  ***Eye-tracking Recording and Analysis***

11      Participants were calibrated using the standard 9-point display and monocular eye

12  movements were recorded at a sampling rate of 1,000 Hz (but were resampled at 250 Hz during

13  pre-processing, which is standard for this type of data). As in previous studies (Kapnoula &

14  Samuel, 2019; McMurray et al., 2002), this was automatically parsed into saccades and fixations

15  using default psychophysical parameters. Adjacent saccades and fixations were combined into a

16  single "look" that started at the onset of the saccade and ended at the offset of the fixation.

17      Eye movements were recorded from the onset of the trial (presentation of unfamiliar

18  object for training trials; red circle for testing trials) through the participant's response (mouse

19  click). This resulted in a variable trial offset time, depending on the individual response time. We

20  adopted the approach of many prior studies (Allopenna et al., 1998; McMurray et al., 2002) by

21  setting a fixed trial duration of 2,000 ms. If a trial ended before this point, we extended the last

22  eye movement; trials longer than 2,000 ms were truncated. This approach assumes that any

1    fixations made in the very late portions of a trial reflect the word the participant settled on and

2    should, thus, be interpreted as an estimate of the final state of the system.

3        In converting the coordinates of each look to the object being fixated, the boundaries of

4    the regions of interest containing the objects were extended by 100 pixels in order to account for

5    noise and/or head-drift in the eye-track record. This did not result in any overlap between the

6    objects (the dead space between pictures was 224 pixels).

7    **Results**

8        Two participants were excluded from the analyses of fixations (but were included in the

9    analyses of responses) due to eye-tracking problems[4].

10   **Analyses of responses**

11       *Training*. Participants performed the task without difficulties and responded in a prompt

12   manner. In addition, their responses were checked offline by a trained research assistant, who

13   verified that they were doing the task as requested. Specifically, spoken responses from the

14   production task were processed with CheckVocal (Protopapas, 2007) to check accuracy[5] and

15   placement of response time (RT) marks. Accuracy was at 100% across participants, and average

16   RT was 648 ms (SD = 129 ms).

17       *Testing*. Average accuracy in testing was 98.1% (SD = 4.9%), which corresponds to 1.2

18   error trials (out of 64) per participant. Only correct trials were included in the reaction time (RT)

19   analyses. Average RT was 1,210 ms (SD = 267 ms).

20       We assessed the effects of training condition (Perception-Only versus Production) and

21   length of training (Phase 1 versus Phase 2) on accuracy (logit-transformed) and RT using $2 \times 2$

---

[4] These participants seemed to use their peripheral vision instead of looking directly at the pictures. The same applies to the exclusion of participants due to eye-tracking problems in Experiments 2 and 3.
[5] An utterance was marked as correct only if all phonemes were pronounced correctly.

1     repeated measures ANOVAs. For accuracy, Condition was not significant, $F_{(1,39)}$=0.241, p=.626,

2     $\eta^2$=.006, but Phase was, $F_{(1,39)}$=9.144, p=.004, $\eta^2$=.190, with participants showing higher

3     accuracy in Phase 2 (99.1% compared to 97.1% in Phase 1). The interaction was not significant,

4     $F_{(1,39)}$=0.218, p=.643, $\eta^2$=.006. Similarly, for RT, Condition was not significant, $F_{(1,39)}$=0.097,

5     p=.757, $\eta^2$=.002, but Phase was, $F_{(1,39)}$=45.952 p<.001, $\eta^2$=.541, with participants giving faster

6     responses in Phase 2 (1,092 ms compared to 1,327 ms in Phase 1). The Phase×Condition

7     interaction was not significant, $F_{(1,39)}$=0.041, p=.841, $\eta^2$=.001. Overall, the behavioral results

8     showed that participants were faster and more accurate on the Phase 2 test.

9     **Analyses of fixations**

10     Next, we analyzed participants' eye-movements. These analyses only include testing

11     trials. For all analyses of eye-movements, we adopted the linking hypothesis originally proposed

12     by Allopenna et al. (1998), according to which fixation proportions can be used as a direct

13     measure of lexical activation.

14     *Replication of Zamuner et al. (2016).* We started by looking only at the data directly

15     comparable to the data reported by Zamuner et al. (2016). These correspond to the testing block

16     of Phase 1. Similarly to Zamuner et al. (2016), we opted for an analysis sensitive to the dynamic

17     changes of lexical activation over time. However, in contrast to the original study, we fitted our

18     data using a nonlinear curve-fitting approach (Farris-Trimble & McMurray, 2013; McMurray et

19     al., 2010; Seedorff et al., 2018), rather than a growth curve analysis (GCA; Mirman et al., 2008).

20     Curve-fitting, like GCA, is not restricted to a specific time window; this analytical approach is

21     based on taking each participant's fixation data (per condition) and finding a set of parameter

22     values that best describes the shape of the fixation curve as a whole. This means that we do not

23     test for differences in specific time-windows, but rather differences in aspects of the overall
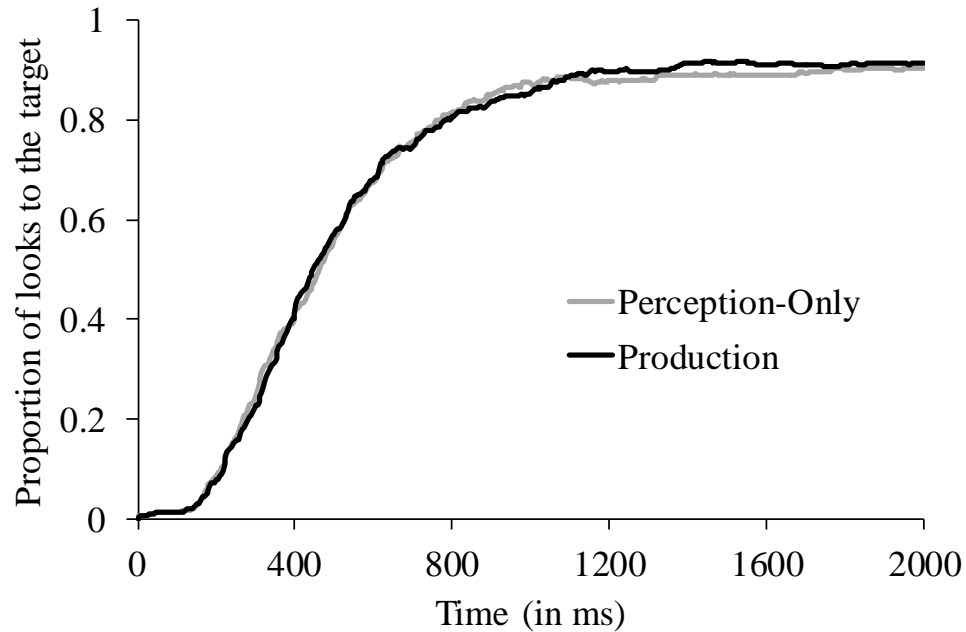
1  shape of the data. Specifically, curve-fitting parameters can be mapped onto psychologically

2  meaningful aspects of lexical activation. For example, given the linking hypothesis described

3  above (according to which, fixation proportions map to lexical activation), the steepness of the

4  ascending slope can be mapped to speed of activation. This transparency provides a more

5  straight-forward comparison between experimental conditions (Farris-Trimble & McMurray,

6  2013; McMurray et al., 2008; Scheepers et al., 2008).

7       We fit our data using a four-parameter logistic function (see Eq.1. in McMurray et al.,

8  2010). In this equation, the lower/higher asymptotes correspond to the baseline/peak of the curve

9  (i.e., minimum and maximum lexical activation) respectively, the slope reflects how quickly

10  lexical activation builds up in time, and the crossover corresponds to the point in time when

11  activation crosses from the lower half of the range to the higher half (e.g., if baseline is 0 and

12  peak is 1, then the crossover would correspond to the point in time when activation is .5).

13       All curve-fitting analyses were implemented using the bdots R package (Oleson,

14  Cavanaugh, McMurray, & Brown, 2015; Seedorff, Oleson, Brown, Cavanaugh, & McMurray,

15  2017). First, we computed the average proportion of looks to the target for each time point along

16  the entire time-course of the trial (i.e., 0 – 2,000 ms) separately for each subject and each training

17  condition (see Figure 2).

1

2   *Figure 2*. Proportion of looks to the target in time for each training condition (Perception-Only
3   versus Production) in the first testing block of Experiment 1.

4

5        Then, we used the bdots logistic.fit function to find the four-parameter logistic function

6   that provided the best fit for each curve. Using the bdots logistic.boot function (for paired data),

7   we tested the effect of training condition on each of the three parameters of interest (peak, slope,

8   and crossover; see Seedorff, Oleson, Cavanaugh, & McMurray, 2017; Seedorff et al., 2018, for a

9   presentation of the bdots package, its conceptual implementation, and a discussion of the

10   statistical approach).

1     Table 2. *Comparisons for the logistic parameters used to describe target activation in the testing*
2     *block of Phase 1 (Experiment 1)*

| Parameter | Difference between conditions (Perception-Only minus Production) | $t_{(35)}$ | SE | p |
|---|---|---|---|---|
| Peak (*p*) | -0.0275 | -7.121 | 0.004 | <0.001 |
| Slope (*s*) | -0.0002 | -1.962 | <0.001 | 0.058 |
| Crossover (*c*) | -6.8892 | -0.713 | 9.662 | 0.481 |

3     *Note 1*. All *t* tests rely on simulations; they are calculated using the bootstrapped means and are adjusted to account
4     for the additional variance around the parameter estimate (Seedorff et al., 2017, 2018).
5     *Note 2*. Differences in dfs between phases are due to different number of excluded bad fits.
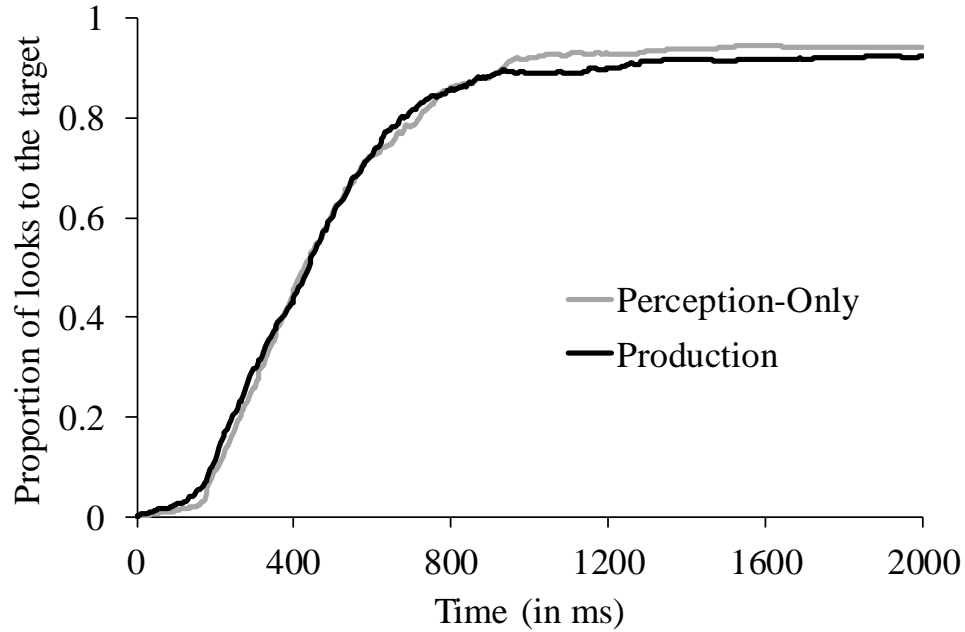6     These notes also apply to Tables 3, 4, 5, 6, and 7.
7

8         As seen in Table 2, there was a significant difference between conditions in their

9     asymptotes (reflecting maximum lexical activation). The direction of this effect indicates an

10    advantage for the Production condition, meaning ultimately higher lexical activation for words

11    that were produced during training. In addition, there was a difference in the slopes (reflecting

12    speed of activation build-up) in the same direction (i.e., Production advantage), but it was not

13    significant. These results are consistent with the findings reported by Zamuner et al. (2016).

14         *Modulation of the production effect by length of training*. Next, we asked whether this

15    positive effect of production is maintained after additional training[6]. We thus looked at the trials

16    from the testing block of Phase 2. Average proportions of looks to the target are plotted in Figure

17    3.

---

[6] In all three Experiments, we examined the effect of Condition as a function of training length by testing the difference between Conditions in each Phase and then comparing the results between Phases. A direct statistical comparison of the Condition effect between Phases is reported below in the **Summary and additional analyses across Experiments**; see full 3×2×2×2 (Experiment×Condition×Phase×TimeWindow) ANOVA.

1

2  *Figure 3.* Proportion of looks to the target in time for each training condition (Perception-Only
3  versus Production) in the second testing block of Experiment 1.

4

5          The same analytical approach was adopted as above.

6  Table 3. *Comparisons for the logistic parameters used to describe target activation in the testing*
7  *block of Phase 2 (Experiment 1)*

| Parameter | Difference between conditions (Perception-Only minus Production) | $t_{(34)}$ | SE | p |
|---|---|---|---|---|
| Peak (*p*) | 0.0144 | 3.722 | 0.004 | <0.001 |
| Slope (*s*) | -0.0001 | -1.531 | <0.001 | 0.135 |
| Crossover (*c*) | 2.9198 | 0.277 | 10.534 | 0.783 |

8

9          As seen in Table 3, there was again a significant difference between conditions in their

10  asymptotes. However, in this case the direction of the effect indicated an advantage for the

11  Perception-Only condition (i.e., ultimately higher lexical activation for words assigned to the

1    Perception-Only condition, or put another way, lower lexical activation for words that were

2    produced during training).

3    **Discussion**

4         In line with Zamuner et al. (2016), Experiment 1 showed an early advantage for novel

5    words that were repeated during training; they were more strongly activated after a few training

6    trials. However, this effect was reversed after additional training; with more training, production

7    had a detrimental effect.

8         One possible interpretation of this reversal is that the effect of production depends on the

9    learning stage. That is, production may be particularly helpful during the first stages of word

10   learning, perhaps facilitating the very first encoding of a novel word; in contrast, when it comes

11   to integrating the lexical representation into the mental lexicon, its effect may be more harmful

12   than helpful. This later detrimental effect could be due to 1) production being disruptive, 2)

13   perception-only being more helpful, or 3) both (see further discussion in the General

14   Discussion).

15        As discussed in the Introduction, one limitation of the Zamuner et al. (2016) design is

16   that it confounds production with speaker variability. That is, words assigned to the production

17   condition were also spoken by one additional voice during training. Prior work has shown a

18   facilitatory role of talker variability in word learning (Richtsmeier et al., 2009; Rost &

19   McMurray, 2009). Thus, any advantage for the Production condition could be due to the

20   additional talker in training. To address this, we ran Experiment 2, which was otherwise identical

21   to Experiment 1, but controlled for the effect of speaker variability.

1                                          **Experiment 2**

2          Experiment 2 assesses the effect of production independently of speaker variability and

3     provides an additional test of the reversal of the effect as a function of amount of training. To

4     match the number of talkers between conditions, we added an additional talker in the Perception-

5     Only condition. If the facilitatory effect of the Production condition was due to talker variability,

6     it should disappear in Experiment 2.

7     **Method**

8     *Participants*

9          Forty-one (28 females; mean age = 24.3 years) native speakers of Spanish participated in

10    Experiment 2. Experiment 2 was identical to Experiment 1 in terms of participant characteristics,

11    compensation, and ethical approval procedures.

12    *Design*

13         The same design as that of Experiment 1 was used with one critical difference: auditory

14    stimuli in the Perception-Only condition were presented in two different speaker voices. As a

15    result of this change, participants in both groups heard two different voices on each training trial

16    (either two different voices played to them, or one played to them plus their own voice).

17    *Materials*

18         All items were the same as in Experiment 1. The only difference was that all items were

19    recorded by an additional speaker of a different gender (a native male speaker of Spanish). In

20    addition, we controlled for both inter-and intra-talker variability; to match the acoustic variability

21    naturally present in the participants' own utterances (i.e., the fact that each time a speaker says a

22    word the acoustics are slightly different), we selected multiple recordings of the words, such that

1    each token was only heard once during the entire experiment (i.e., there were twelve different

2    tokens for each word, as many as the training repetitions per word). The new recordings were

3    pre-processed following the same steps as in Experiment 1. The average duration of the new

4    stimuli (including the 100 ms of silence) was 694 ms. The new items were presented during

5    training in the Perception-Only condition after the "listen" prompt (i.e., for the second

6    presentation of each word). The original stimuli from Experiment 1 were used in training for the

7    first presentation (i.e., before the visual prompt) for both conditions, and during testing.

8       Visual stimuli were the same as in Experiment 1.

9    *Procedure*

10      The procedure was identical to that of Experiment 1.

11    *Eye-tracking Recording and Analysis*

12      Eye-tracking recording and pre-processing were identical to those of Experiment 1.

13    **Results**

14    **Analyses of responses**

15      *Training*. Participants performed the task without problems and their responses were

16    checked offline by a trained research assistant, who verified that they were doing the task as

17    requested. That is, as in Experiment 1, spoken responses from the production task were

18    processed with CheckVocal (Protopapas, 2007). Accuracy was at 100% across participants,

19    while average RT was 705 ms (SD = 128 ms).

20      *Testing*. Average accuracy in testing was 96.4% (SD = 5.3%), which corresponds to 2.3

21    error trials (out of 64) per participant. Any trials with incorrect responses were excluded from RT

22    analyses. In the remaining trials, average RT was 1,320 ms (SD = 271 ms).

1      We assessed the effects of training condition and length of training on accuracy (logit-

2   transformed) and RT using the same analytical approach as in Experiment 1. For accuracy,

3   Condition was not significant, $F_{(1,40)}=0.291$, p=.592, $\eta^2=.007$, but Phase was, $F_{(1,40)}=16.956$,

4   p<.001, $\eta^2=.298$, with participants showing higher accuracy in Phase 2 (99.2% compared to

5   93.7% in Phase 1). The interaction was not significant, $F_{(1,40)}=1.321$, p=.257, $\eta^2=.032$. For RT,

6   Condition was significant, $F_{(1,40)}=6.186$, p=.017, $\eta^2=.134$, with participants giving faster

7   responses for items that had been in the Perception-Only condition during training (1,284 ms

8   compared to 1,371 ms for Production items). Phase was also significant, $F_{(1,40)}=31.220$, p<.001,

9   $\eta^2=.438$, with participants giving faster responses in Phase 2 (1,171 ms compared to 1,484 ms in

10  Phase 1). The interaction did not reach significance, $F_{(1,40)}=2.849$, p=.099, $\eta^2=.066$. Even though

11  the interaction was not significant, in the interest of comprehensiveness, we conducted

12  Bonferroni-corrected post-hoc comparisons, which revealed a detrimental effect of Production in

13  Phase 1, $F_{(1,40)}=4.830$, p=.034, $\eta^2=.108$, which was in the same direction, but not significant, in

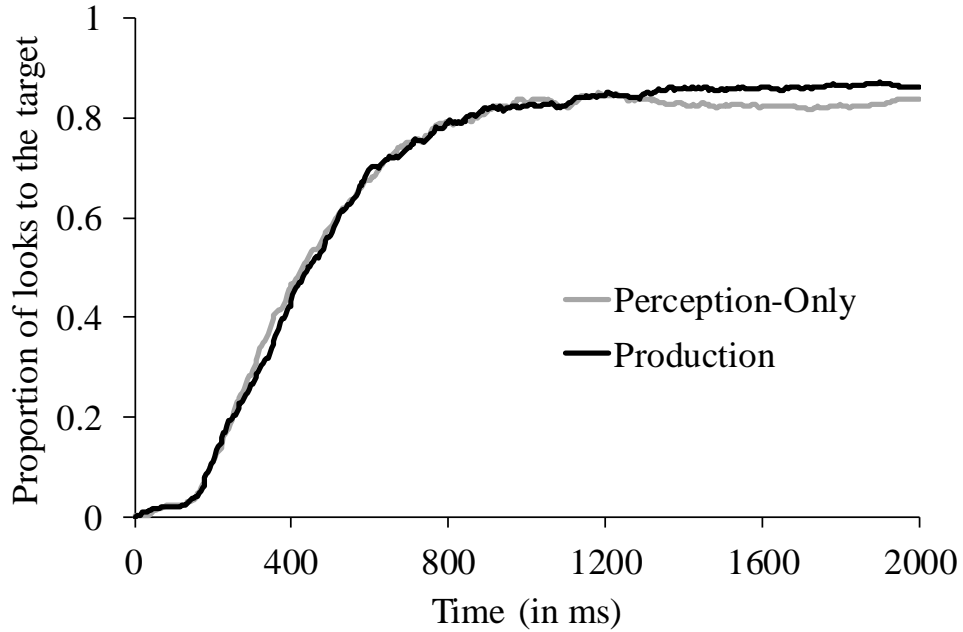14  Phase 2, $F_{(1,40)}=3.617$, p=.064, $\eta^2=.083$.

15      Overall, the behavioral results showed that participants were faster and more accurate in

16  Phase 2, but they were also consistently faster in recognizing items that had been assigned to the

17  Perception-Only training condition.

18  **Analyses of fixations**

19      Next, we analyzed participants' eye-movements during testing.

20  *Modulation of production effect by length of training.* We first looked at the effect of

21  production as a function of training length (i.e., Phase 1 versus Phase 2; average proportions of

22  looks to the target in time for each Phase are plotted in Figures 4 and 5). The same analytical

23  approach was adopted as in Experiment 1 (i.e., curve-fitting).

1



2    *Figure 4.* Proportion of looks to the target in time for each training condition (Perception-Only
3    versus Production) in the first testing block of Experiment 2.

4



5

6    *Figure 5.* Proportion of looks to the target in time for each training condition (Perception-Only
7    versus Production) in the second testing block of Experiment 2.

1

2 Table 4. *Comparisons for the logistic parameters used to describe target activation in the testing*

3 *block of Phase 1 (Experiment 2)*

| Parameter | Difference between conditions (Perception-Only minus Production) | $t_{(32)}$ | SE | p |
|---|---|---|---|---|
| Peak (*p*) | -0.0188 | -2.749 | 0.007 | 0.001 |
| Slope (*s*) | 0.0010 | 8.644 | <0.001 | <0.001 |
| Crossover (*c*) | -12.2863 | -0.542 | 22.674 | 0.592 |

4

5

6 Table 5. *Comparisons for the logistic parameters used to describe target activation in the testing*

7 *block of Phase 2 (Experiment 2)*

| Parameter | Difference between conditions (Perception-Only minus Production) | $t_{(36)}$ | SE | p |
|---|---|---|---|---|
| Peak (*p*) | 0.0090 | 2.714 | 0.003 | 0.010 |
| Slope (*s*) | -0.0001 | -1.197 | <0.001 | 0.239 |
| Crossover (*c*) | 53.188 | 5.546 | 9.256 | <0.001 |

8

9        The results were similar to Experiment 1 in two respects. First, there was an early

10 facilitatory effect of production (significantly higher activation peak for produced items after the

11 first training Phase; see Table 4). Second, after the second training Phase, the efficacy of the two

12 training conditions reversed, resulting in higher lexical activation for words assigned to the

13 Perception-Only condition (see Table 5).

14        In contrast to Experiment 1, there was a significant difference in slope in the first testing

15 Phase, indicating faster activation for Perceived-Only items, and a significant difference in

16 crossover in the second testing phase, indicating an earlier onset of activation for produced items

1    (see Discussion below). A direct comparison across all three experiments is presented in the

2    **Summary and additional analyses across Experiments** section below.

3    **Discussion**

4         The critical difference between Experiments 1 and 2 was that the latter controlled for the

5    effect of speaker variability. This allowed us to better examine the effect of production in and of

6    itself. Experiment 2 replicated the reversal of the production effect that was observed in

7    Experiment 1: an early facilitatory effect of production turned into a detrimental effect after

8    additional training (i.e., a significantly higher activation peak for Perceived-Only items in Phase

9    2).

10        In Experiment 2, we found a steeper slope for items assigned to the Perception-Only

11   condition (i.e., faster activation). Even though this effect may seem inconsistent with the early

12   facilitatory effect of production (i.e., higher activation asymptote for produced items), it could

13   reflect the way in which novel items are gradually integrated into the mental lexicon. For

14   example, one possibility is that production facilitates early encoding of novel lexical

15   representations (the "lexical configuration" stage proposed by Leach & Samuel, 2007), but

16   uninterrupted perception leads to better overall integration into the mental lexicon (the "lexical

17   engagement" stage proposed by Leach & Samuel, 2007). The automatization of lexical

18   processing is considered to be a marker of deeper lexical integration (i.e., full lexical

19   engagement; see discussion by McMurray, Kapnoula, and Gaskell, 2016). From this perspective,

20   the slope should in fact be steeper for Perception-Only items to the extent that it reflects

21   automatization of processing.

22        In Phase 2, our analysis showed a later crossover for items assigned to the Perception-

23   Only condition. At first, this seems to be a surprising result (given that in the same Phase we

1 observed overall higher activation of Perception-Only items), but it is in fact consistent with the

2 rationale laid down in the previous paragraph, according to which production may facilitate early

3 encoding (lexical consolidation), but perception may lead to better overall integration (lexical

4 engagement). In other words, the *onset* of activation (better reflected by the crossover parameter)

5 may rely on the configuration status of a novel word, whereas the *speed* of activation (better

6 reflected by the slope parameter) should depend on the degree of automatization. Thus, if

7 production facilitates lexical configuration, it would make sense to see a facilitatory effect of

8 production on crossover.

9      Taken together, the results from Experiments 1 and 2 suggest that producing a novel

10 word may have an early advantage, but its effect becomes detrimental with additional training. In

11 Experiment 3, we consider the possibility that production may play a positive role by promoting

12 the abstraction of newly acquired lexical representations, in which case we should observe better

13 generalization of learning to novel instances of learned words. In Experiments 1 and 2 the testing

14 voice was the same as at least one of the voices used in training. That means that during testing

15 all items (in both conditions) were heard in a familiar voice, one that had been heard before. Our

16 new question is whether Production (versus Perception-Only) might help learners *generalize* to

17 novel speakers. To address this question, in Experiment 3, all testing items were presented in a

18 new voice. If Production helps generalization, we should see a stronger facilitatory effect of

19 Production in Phase 1 and perhaps a weaker detrimental effect in Phase 2. In contrast, if

20 Perception-Only leads to better generalization via stronger lexical engagement, we should see the

21 opposite pattern.

1    **Experiment 3**

2    Experiment 3 tests whether production helps listeners in recognizing novel words spoken

3    by an unfamiliar talker (i.e., generalization). To assess this, test stimuli were presented in a novel

4    voice. As in Experiment 2, we matched the number of talkers between training conditions by

5    having an additional talker in the Perception-Only condition.

6    **Method**

7    *Participants*

8    Forty-one (28 females; mean age = 24.9 years) native speakers of Spanish participated in

9    Experiment 3. Experiment 3 was identical to Experiments 1 and 2 in terms of participant

10   characteristics, compensation, and ethical approval procedures.

11   *Design*

12   The same design as that of Experiment 2 was used with one critical difference: auditory

13   stimuli presented in testing were spoken by one of two new speakers (one male and one female).

14   Twenty participants were randomly assigned to one speaker and the rest to the other.

15   *Materials*

16   All training stimuli were identical to those of Experiment 2. The only difference between

17   Experiments 2 and 3 was in the testing stimuli. All items were recorded by two additional native

18   speakers of Spanish (one male and one female). These items were used to replace the testing

19   stimuli used in Experiments 1 and 2. As with the male training items of Experiment 2, we

20   selected multiple recordings of the words, such that each token was only heard once during the

21   entire experiment (i.e., there were eight different tokens for each word). This was done to

22   increase acoustic variability of the testing stimuli, as in Experiment 2. The new recordings were

1    pre-processed following the same steps as in Experiments 1 and 2. The average duration of the

2    new stimuli (including the 100 ms of silence) was 624 ms for the male and 840 ms for the female

3    speaker.

4          Visual stimuli were the same as in Experiment 1.

5    *Procedure*

6          The procedure was identical to that of Experiments 1 and 2.

7    *Eye-tracking Recording and Analysis*

8          Eye-tracking recording and pre-processing were identical to that of Experiments 1 and 2.

9    **Results**

10          Two participants were excluded from the analyses of fixations (but were included in the

11    analyses of responses) due to eye-tracking problems.

12    **Analyses of responses**

13          *Training*. Participants performed the task without problems and their responses were

14    checked offline by a trained research assistant, who verified that they were doing the task as

15    requested. Once again, spoken responses from the production task were processed with

16    CheckVocal (Protopapas, 2007). Accuracy was at 100% across participants, while average RT

17    was 689 ms (SD = 144 ms).

18          *Testing*. Average accuracy in testing was 97.3% (SD = 4.5%), which corresponds to 1.7

19    error trials (out of 64) per participant. Any trials with incorrect responses were excluded from RT

20    analyses. In the remaining trials, average RT was 1,256 ms (SD = 227 ms).

21          We assessed the effects of training condition and length of training on accuracy (logit-

22    transformed) and RT using the same analytical approach as in Experiments 1 and 2. Accuracy

1    was higher for Production (97.9%) compared to Perception-Only items (96.9%), but this

2    difference was not significant, $F_{(1,40)}=3.121$, p=.085, $\eta^2=.072$. Phase was significant,

3    $F_{(1,40)}=21.168$, p<.001, $\eta^2=.346$, with higher accuracy in Phase 2 (99.7% compared to 95.0% in

4    Phase 1). The interaction was not significant, $F_{(1,40)}=0.205$, p=.653, $\eta^2=.005$. For RT, Condition

5    was not significant, $F_{(1,40)}=0.710$, p=.404, $\eta^2=.017$, but Phase was, $F_{(1,40)}=49.961$, p<.001,

6    $\eta^2=.555$, with participants giving faster responses in Phase 2 (1,125 ms compared to 1,398 ms in

7    Phase 1). The interaction was not significant, $F_{(1,40)}=1.378$, p=.247, $\eta^2=.033$.

8         Overall, the behavioral results showed that participants were faster and more accurate in

9    Phase 2.

10   **Analyses of fixations**

11        Next, we analyzed participants' eye-movements during testing.

12        *Modulation of production effect by length of training*. The same analytical approach was

13   adopted as in Experiments 1 and 2 (i.e., curve-fitting). Average proportions of looks to the target

14   in time for each Phase are plotted in Figures 6 and 7.

1

*Figure 6.* Proportion of looks to the target in time for each training condition (Perception-Only versus Production) in the first testing block of Experiment 3.

4

*Figure 7.* Proportion of looks to the target in time for each training condition (Perception-Only versus Production) in the second testing block of Experiment 3.

7

1 Table 6. *Comparisons for the logistic parameters used to describe target activation in the testing*
2 *block of Phase 1 (Experiment 3)*

| Parameter | Difference between conditions (Perception-Only minus Production) | $t_{(34)}$ | SE | p |
|---|---|---|---|---|
| Peak (*p*) | -0.0063 | -0.978 | 0.007 | 0.335 |
| Slope (*s*) | 0.0001 | 0.747 | <0.001 | 0.460 |
| Crossover (*c*) | 0.369 | 0.044 | 10.444 | 0.972 |

3
4

5 Table 7. *Comparisons for the logistic parameters used to describe target activation in the testing*
6 *block of Phase 2 (Experiment 3)*

| Parameter | Difference between conditions (Perception-Only minus Production) | $t_{(37)}$ | SE | p |
|---|---|---|---|---|
| Peak (*p*) | 0.0115 | 3.117 | 0.004 | 0.004 |
| Slope (*s*) | -0.0001 | -1.142 | <0.001 | 0.261 |
| Crossover (*c*) | 21.6634 | 3.937 | 5.502 | <0.001 |

7

8    In contrast to Experiments 1 and 2, there was no early facilitatory effect of production

9 (see Table 6). However, in line with the previous experiments, we again observed higher lexical

10 activation for words assigned to the Perception-Only condition in the testing block of Phase 2

11 (see Table 7). In addition, similarly to Experiment 2, there was again a significant crossover

12 difference in Phase 2, indicating an earlier activation onset for produced items (see Table 7).

13 **Discussion**

14    As in Experiments 1 and 2, we again observed a detrimental effect of production after

15 additional training (i.e., significantly higher activation peak for Perception-Only items in Phase

16 2). In addition, as in Experiments 1 and 2, the activation peak was higher for Production

17 compared to Perception-Only items in Phase 1, though in this case the effect was not significant.

1    As in Experiment 2, we again found an earlier crossover for produced items in Phase 2.

2   This effect is in line with the idea (see the discussion in Experiment 2) that the onset of lexical

3   activation (as reflected by the crossover parameter) depends on the early configuration of a novel

4   word, rather than its integration into the lexicon.

5                    **Summary and additional analyses across Experiments**

6        Experiments 1 and 2 replicated the finding reported by Zamuner et al. (2016) that

7   production helps at the early stages of word learning. In Experiment 3, the results from Phase 1

8   were in the same direction, but were not significant. In contrast, and critically, we observed a

9   reliable detrimental effect of production across all three experiments after additional training

10  (i.e., higher activation of non-produced words in Phase 2). This reversal of the effect seems to

11  reflect a dissociation regarding the time course of lexical integration: *Even though early*

12  *encoding may be facilitated by production, further integration is better served by perception*. In

13  line with this interpretation, we also observed a facilitatory effect of production on the crossover

14  (likely reflecting earlier activation onset), which contrasted with a facilitatory effect of

15  perception on the slope (likely reflecting automatization of processing). That is, the facilitation

16  of production on lexical encoding may be reflected by the earlier activation onset (crossover),

17  whereas the facilitation of perception on lexical integration may be reflected in the activation

18  speed (slope).

19       We conducted additional analyses to examine these patterns in greater detail.

20  Specifically, we were interested in testing for any significant differences between early and late

21  stages of lexical processing. To do so, we split each trial into two (early/late) parts. Since this

22  was a post-hoc, exploratory analysis, we chose to avoid splitting the trials based on an arbitrary,

1   experimenter-driven criterion. Instead, each trial was split based on the offset of the auditory

2   stimulus (corrected for 200 ms oculomotor delay). This allowed us to use a flexible, stimulus-

3   driven time window and account for any variability between experiments, speakers, and stimuli[7].

4   In addition, this analysis allowed us to directly test all possible interactions between independent

5   variables.

6   **Results**

7         Here we focused on the effect of training condition (Perception-Only versus Production)

8   on the activation of novel words as a function of 1) amount of training and 2) time-point within a

9   trial. Results from all three experiments were included in the analyses (see Figure 8). We started

10  by running the full 3 (Experiment: 1/2/3) $\times$ 2 (Condition: Perception-Only/Production) $\times$ 2

11  (Phase: 1/2) $\times$ 2 (TimeWindow: early/late) repeated-measures ANOVA with average proportion

12  of fixations to the target (empirical-logit-transformed) as the DV. Detailed results of this

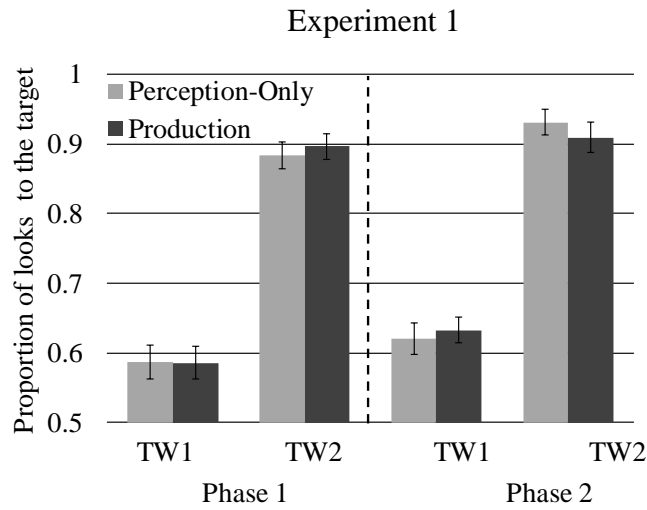13  ANOVA and its follow-ups are listed in the Appendix.

14        Phase was significant, $F_{(1,114)}=45.665$, $p<.001$, $\eta^2=.286$, indicating that participants were

15  better at activating the target word after receiving additional training, as expected. Condition was

16  not significant, $F_{(1,114)}=0.013$, $p=.909$, $\eta^2<.001$, and neither was the Condition $\times$ Phase

17  interaction, $F_{(1,114)}=3.107$, $p=.081$, $\eta^2=.027$. The Condition $\times$ Phase $\times$ TimeWindow interaction

18  was significant, $F_{(2,114)}=13.481$, $p<.001$, $\eta^2=.106$, reflecting a differential effect of training

19  condition depending on training length and time-point within trial. Neither the 4-way, nor any of

20  the other 3-way interactions were significant.

21

[7] A flexible splitting point that is time-locked to each stimulus takes into account the variability in stimulus duration. This leads to time-windows that are informationally comparable between items. Furthermore, a splitting point roughly close to the middle of the trial (as is the case with the splitting point used) means that the two time windows are comparable in terms of number of data points.
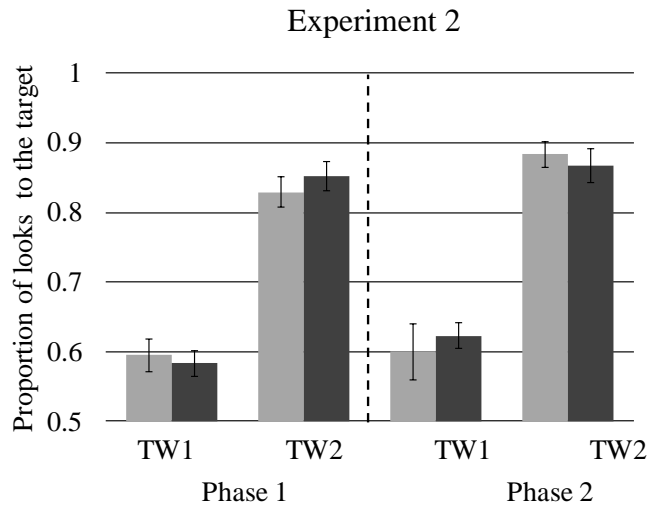
1

### Experiment 1



2

3

### Experiment 2



4

5

### Experiment 3
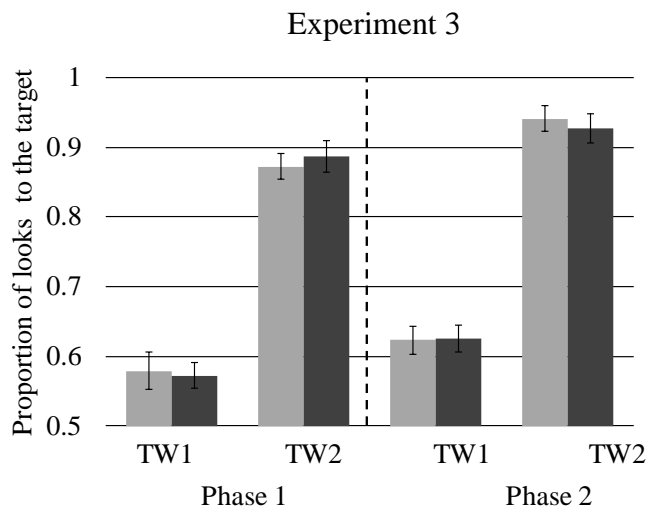


6
7  *Figure 8.* Average proportions of looks to the target per training condition (Perception-
8  Only/Production), testing phase (Phase 1/Phase 2), and time window (TW) within trial (early:
9  TW1/late: TW2) for each of the three experiments. Error bars indicate ±1 within-subject standard

1  error of the mean (Cousineau, 2005; Loftus & Masson, 1994; Morey, 2008). *Note.* An alternative
2  visualization of the data is offered in Figure A1 in the Appendix.
3

4      In following up on the significant Condition × Phase × TimeWindow interaction, post-

5  hoc comparisons (Bonferroni-corrected) showed a significant Condition effect in the late

6  TimeWindow of Phase 2, p=.022, the direction of which indicates more looks to the Perception-

7  Only items late in the trial. In contrast, in the early TimeWindow of Phase 2, the trend was in the

8  opposite direction (i.e., more looks to Production items; p=.059). No simple effect of Production

9  was found in Phase 1.

10     The results of these within-trial analyses show a shift in the direction of the Production

11  effect in Phase 2; specifically, if there is any effect of production early in a trial, this effect seems

12  to be facilitatory; later in the trial, production is clearly detrimental. This within-trial pattern

13  echoes the curve-fitting results: Even though production seems to expedite the activation onset of

14  new words, providing an advantage to produced items early in the trial, items that had only been

15  perceived, without production, enjoy a faster build-up of activation.

16                          **General Discussion**

17     Across three experiments, we examined the role of production on word learning while

18  carefully controlling for other variables such as amount of exposure, talker variability, recall-

19  based facilitation, and attention. Our findings show a robust dissociation: At first production

20  helps, but as learning advances its effect becomes detrimental. This dissociation can potentially

21  reconcile previous results showing both facilitatory (Dodson & Schacter, 2001; Gathercole &

22  Conway, 1988; Hopkins & Edwards, 1972; Hopman & MacDonald, 2018; P. MacDonald &

23  MacLeod, 1998; MacLeod et al., 2010; MacLeod & Bodner, 2017; Zamuner et al., 2016) and

1 detrimental (Baese-Berk, 2019; Baese-Berk & Samuel, under review, 2016; Leach & Samuel,

2 2007; Zamuner et al., 2018) effects of production. Moreover, our findings provide fine-grained

3 timing information about the effect of production on learning – both at the level of training phase

4 (i.e., few versus many training trials) and in terms of real-time processing, at the single-trial

5 scale. In that sense, the present work complements previous work such as the Leach & Samuel's

6 (2007) study that showed larger negative effects of production on perceptual learning but without

7 offering timing information.  As a result, our results offer valuable insights into the mechanisms

8 underlying the seemingly contrasting effects of production.

9 **Early facilitatory effect of production**

10     We found evidence for an early facilitatory effect of production across three experiments.

11 In Experiments 1 and 2, novel words that had been produced during training had a significantly

12 higher activation peak after the first training block (i.e., an effect on the peak); Experiment 3

13 showed the same pattern, but the difference did not reach significance. In other words,

14 production was helpful during the earliest stages of word learning, i.e., during the initial

15 encoding of novel word-forms.

16     In addition, in Experiments 2 and 3, produced items had an advantage at the onset of

17 lexical activation (i.e., an effect on the crossover). The early advantage for production was

18 echoed in our within-trial analyses (the significant Condition $\times$ Phase $\times$ TimeWindow

19 interaction). We suggest that this early activation advantage is a result of production helping the

20 *initial* encoding of novel words. Better lexical configuration of a novel word-form can, in turn,

21 facilitate an earlier onset of activation. Note that in Leach and Samuel's (2007) original contrast

22 between lexical configuration and lexical engagement, an important basis for the distinction was

1    their finding of a facilitatory effect of production on lexical configuration, consistent with this

2    finding.

**Late detrimental effect of production**

4    Evidence for a detrimental effect of production was even more robust; across all

5    experiments and analyses, produced words had a lower activation peak after the second training

6    block. In addition, produced words displayed slower activation build-up (i.e., an effect on the

7    slope) in Experiment 2. This pattern likely indicates that words that were only heard (not

8    produced) were better integrated into the mental lexicon. That is, once a word is well integrated

9    into the system, its recognition reaches a higher level of automatization. As a result, the

10   recognition process moves faster (reflected by the higher slope) and is more effective (reflected

11   by the higher peak).

12   This finding conforms to Leach and Samuel's (2007) concept of lexical engagement, with

13   their results showing better lexical engagement for listeners in a Perception-Only training

14   condition than in a Production condition. This interpretation is also in line with Kapnoula and

15   Samuel (2019), who found that newly learned words were activated faster after participants slept

16   (indicated by the steeper slopes of fixation probability curves). Given the well-documented

17   strengthening role of sleep consolidation in lexical integration, this can be taken as indirect

18   evidence that better integrated words show a more robust (faster and/or higher) pattern of

19   activation.

**Towards a reconciling mechanism**

21   As discussed in the Introduction, adding a production requirement may affect the

22   outcome of word learning in a number of different ways, some of which are not intrinsic to

23   production per se. Our study was designed to control for confounding factors as much as

1    possible, which potentially allows us to identify the mechanism(s) driving any true production

2    effect. For example, participants were always forced to keep the new phonological sequences in

3    their phonological short-term memory until they saw the prompt, thus equalizing the role of such

4    memory effects across conditions. The same aspect of the procedure also controlled for any

5    differences in attention. In addition, no recall was required during training (only immediate

6    repetition), meaning that any production effect could not be driven by the testing effect

7    (Karpicke & Roediger, 2008). Lastly, in Experiments 2 and 3, we controlled for the effect of

8    speaker variability, which is another frequently confounding variable.

9         With these potentially confounding factors controlled, our experiments were designed to

10   examine *dynamic* effects of producing to-be-learned words. That is, we tested how production

11   affects lexical activation in real time using the high temporal resolution of eye tracking, and how

12   this effect may change during the progression of learning.

13        Our results show an early facilitatory effect of production followed by a late detrimental

14   effect. This pattern can reconcile a number of previously reported findings. As discussed earlier,

15   the pattern aligns very well with the theoretical dissociation proposed by Leach and Samuel

16   (2007), according to which lexical configuration precedes lexical engagement. From this

17   perspective, production seems to help early lexical encoding (configuration), but it hurts lexical

18   integration (engagement). Consistent with previous findings (Kapnoula et al., 2015; Kapnoula &

19   McMurray, 2016), we found evidence for both lexical properties being developed within

20   minutes[8] after the onset of learning. This indicates that any internal adjustments that are made to

21   support the development of these two properties (e.g., formation of bottom-up, lateral, and top-

22   down pathways connecting the novel words to other, known representations) may be

---

[8] In Kapnoula and McMurray (2016) and Kapnoula et al. (2015), evidence for integration was found within 15-30 mins after learning onset.

1  theoretically and mechanistically distinct, but their development unfolds in a cascaded and

2  possibly overlapping fashion.

3      Most importantly, our findings and their interpretation within this dual-stage theoretical

4  framework, bring us closer to a comprehensive mechanism of the production effect. For

5  example, as mentioned in the Introduction, one way that production may help word learning is

6  via the addition of articulatory information. Given our results, adding articulatory information

7  should only help with the early encoding of new word-forms; when the learner has very little

8  other information available, every cue can help. This pattern is reminiscent of what has been

9  found in studies looking for motor area activation during speech perception: This can be found,

10  but this is strongly associated with very challenging listening conditions, when other cues are

11  much less accessible (e.g., Nuttall et al., 2016). Our results speak directly to the issue of

12  how/when production may disrupt word learning. There is not a disruptive effect of production

13  on the early encoding of novel word-forms. When production is disruptive, this is associated

14  with later stages, related to lexical integration (e.g., mapping the word-form to its semantic

15  referent).

16  **Limitations and further questions**

17      Although our results provide incisive information about the learning stage in which

18  production is most disruptive, they do not speak to an additional important question: Is the

19  detrimental effect of production in fact a true negative effect, or is the difference between

20  training conditions (here, and in other studies) due to a facilitatory effect of perception that is

21  reduced under production? For example, it may be that production is disruptive in the sense that

22  participants allocate attentional resources towards producing the word and, as a result, have

23  fewer resources available to take advantage of the perceptual input. This would be a "lost

1 opportunity" effect, rather than an active disruption, caused by production. Indeed, this would be

2 in line with findings showing that production becomes detrimental when cognitive load is high –

3 e.g., because the new words are phonologically unfamiliar (Kaushanskaya & Yoo, 2011); or

4 because they are spoken in an unfamiliar accent (Cho & Feldman, 2016); or because the

5 resources themselves are limited due to the participants' young age (López Assef et al., 2021;

6 Zamuner et al., 2018).

7       This idea is also in line with Kapnoula et al. (2015), who report similar degrees of lexical

8 integration of new words independently of whether they were repeated during training or not.

9 That is, production did not seem to affect lexical integration. Furthermore, Baese-Berk and

10 Samuel (2016; under review) have examined this in the domain of learning a new phonetic

11 contrast, and the data are consistent with this (i.e., the idea that some of the detrimental

12 production effects are due to reduced passive exposure). If this is the mechanism, then

13 presumably the disruption could be alleviated if the production requirement were delayed enough

14 for the perceptual processing to finish. Recent work in our laboratory provides evidence that the

15 same pattern holds for learning new words (Kapnoula & Samuel, under review).

16       Another lingering question is why would production help lexical encoding, but hurt

17 further integration? In the Introduction, we present a number of ways in which production may

18 help *or* hurt word learning; however, we did not expect to find evidence for both. That is, our

19 experiments were designed to examine whether production has a positive/negative effect on

20 different aspects and stages of word learning, but they cannot address why production has

21 opposite effects. We speculate that these effects are driven by different mechanisms. A

22 detrimental effect of production on lexical integration could be due to a lost opportunity for

23 additional passive exposure (as argued above); a facilitatory effect on lexical encoding could be

1    due to small differences in attention. Our design minimized differences in attention: participants

2    did not know ahead of time whether they would be asked to repeat the word or not. However, it

3    is conceivable that participants attended more to the phonological structure of a new word once

4    they were asked to repeat it. Such differences in attention could facilitate lexical encoding.

5    Further work is needed to examine the underlying mechanism(s) behind this complex pattern of

6    findings.

7         A potential criticism of this work is that, if ultimately word recognition is highly

8    successful, arguing for a detrimental effect of production may be a misnomer. Indeed, our results

9    do not address the question of whether production affects the probability of correctly recognizing

10   a newly learned word. However, our aim was to examine the effect of production on the quality

11   of the newly learned lexical representations, as reflected by their real-time activation trajectory.

12   In that respect, our results show that production is detrimental when compared to just hearing the

13   new word. Here, we should note that spoken language comprehension likely depends not only on

14   the accuracy of lexical activation, but also on its speed. According to current theories of spoken

15   language comprehension, sentence comprehension is largely based on activation of lexical

16   representations (Altmann & Kamide, 1999; M. C. MacDonald et al., 1994; McRae et al., 1998;

17   Tanenhaus & Trueswell, 1995; Trueswell, 1996). More critically, there is evidence that spoken

18   word recognition happens in parallel to semantic and syntactic processing (Gussow et al., 2019;

19   Yee & Sedivy, 2006). This means that any delays in activating lexical representations can have

20   important downstream consequences at higher levels of processing. Moreover, given the speed

21   with which the speech signal unfolds, it is reasonable to assume that any delays can gradually

22   accumulate in time, which may lead to a growing difficulty in integrating upcoming input. Thus,

23   speed of lexical activation is likely a critical aspect of efficient spoken language comprehension.

1        More broadly, one may wonder about the degree to which our results are relevant to

2    different word learning situations (e.g., in L1 versus L2, or in naturalistic versus classroom-type

3    settings). There are indeed two aspects of our design that perhaps make the task more similar to

4    L1 word learning; first, the novel words were phonologically, phonotactically, and

5    morphologically consistent with the participants' L1 (Spanish); second, we used unfamiliar

6    objects as visual referents. That said, we believe our results capture something fundamental

7    about the cognitive mechanisms of word learning and, in that sense, our findings are relevant to

8    and have implications for word learning in general. Somewhat related to this, even though our

9    results do not directly speak to the question of how production can be best incorporated into

10    word-learning practices in the real world, our findings can certainly be used as a base on which

11    to formulate experimental hypotheses that are more directly relevant to real-world settings. For

12    example, within a second language learning setting, one may predict that delaying the

13    requirement for students to repeat a new word until after they have been exposed to it a few

14    times may lead to more robust learning. Future research in more naturalistic settings can help us

15    test such predictions.

16        Finally, this work examined the role of production on word learning by assessing

17    participants' ability to recognize, rather than *produce* novel words. Our goal was to report results

18    that are directly comparable to previous work, which also used measures of comprehension. In

19    addition, our experimental paradigm (VWP) allowed us to detect fine differences in the dynamic

20    build-up of lexical activation. This would have been quite difficult if we had used production at

21    test. That said, we acknowledge that knowing a word is not limited to recognizing it.

1  **Conclusion and significance**

2  Our findings demonstrate an early facilitatory effect of production, followed by a (more

3  robust) late detrimental effect. Our interpretation of this reversal is that production may facilitate

4  early encoding of new words, but perception is more helpful when it comes to their deeper

5  integration into the mental lexicon.

6  These findings are consistent with a literature that includes both facilitatory (Dodson &

7  Schacter, 2001; Gathercole & Conway, 1988; Hopkins & Edwards, 1972; Hopman &

8  MacDonald, 2018; P. MacDonald & MacLeod, 1998; MacLeod et al., 2010; MacLeod &

9  Bodner, 2017; Zamuner et al., 2016) and detrimental (Baese-Berk, 2019; Baese-Berk & Samuel,

10  under review, 2016; Leach & Samuel, 2007; Zamuner et al., 2018) effects of production,

11  allowing us to suggest a reconciliation of the seemingly contradictory pattern of results. The

12  theoretical framework of our interpretation is based on the two-stage dissociation proposed by

13  Leach and Samuel (2007), according to which different lexical properties (i.e., configuration and

14  engagement) can follow distinct developmental trajectories, shaped by different variables.

15  The results of our three experiments shed light on the journey of novel words into the

16  mental lexicon. As this line of research develops, it has the potential to inform the educational

17  community, clarifying how and when production can be used most effectively to aid novel word

18  learning, and when it should be avoided.

19

1    **References**

2    Allopenna, P. D., Magnuson, J. S., & Tanenhaus, M. K. (1998). Tracking the Time Course of

3        Spoken Word Recognition Using Eye Movements: Evidence for Continuous Mapping

4        Models. *Journal of Memory and Language*, *38*(4), 419–439.

5        https://doi.org/10.1006/jmla.1997.2558

6    Altmann, G. T. M., & Kamide, Y. (1999). Incremental interpretation at verbs: restricting the

7        domain of subsequent reference. *Cognition*, *73*(3), 247–264. https://doi.org/10.1016/S0010-

8        0277(99)00059-1

9    Baddeley, A. D., Papagno, C., & Vallar, G. (1988). When long-term learning depends on short-

10       term storage. *Journal of Memory and Language*, *27*(5), 586–595.

11       https://doi.org/10.1016/0749-596X(88)90028-9

12   Baese-Berk, M. M. (2019). Interactions between speech perception and production during

13       learning of novel phonemic categories. *Attention, Perception, & Psychophysics*, 1–25.

14       https://doi.org/10.3758/s13414-019-01725-4

15   Baese-Berk, M. M., & Samuel, A. G. (under review). *Just give it time: Differential effects of*

16       *disruption and delay on perceptual learning.*

17   Baese-Berk, M. M., & Samuel, A. G. (2016). Listeners beware: Speech production may be bad

18       for learning speech sounds. *Journal of Memory and Language*, *89*, 23–36.

19       https://doi.org/10.1016/J.JML.2015.10.008

20   Boiteau, T. W., Malone, P. S., Peters, S. A., & Almor, A. (2014). Interference between

21       conversation and a concurrent visuomotor task. *Journal of Experimental Psychology.*

*General*, *143*(1), 295. https://doi.org/10.1037/A0031858

Bradlow, A. R., & Hayes, E. (2003). Speaking Clearly for Children With Learning Disabilities. *Article in Journal of Speech Language and Hearing Research*, *46*(1), 80–97. https://doi.org/10.1044/1092-4388(2003/007)

Bradlow, A. R., Torretta, G., & Pisoni, D. (1996). Intelligibility of normal speech I: Global and fine-grained acoustic-phonetic talker characteristics. *Speech Communication*.

Cho, K. W., & Feldman, L. B. (2016). When repeating aloud enhances episodic memory for spoken words: interactions between production- and perception-derived variability. *Journal of Cognitive Psychology*, *28*(6), 673–683. https://doi.org/10.1080/20445911.2016.1182173

Cousineau, D. (2005). Confidence intervals in within-subject designs: A simpler solution to Loftus and Masson's method. *Researchgate.Net*. https://doi.org/10.20982/tqmp.01.1.p042

Creel, S. C., & Tumlin, M. (2011). On-line acoustic and semantic interpretation of talker information. *Journal of Memory and Language*, *65*(3), 264–285.

Dodson, C. S., & Schacter, D. L. (2001). "If I had said it I would have remembered it": Reducing false memories with a distinctiveness heuristic. *Psychonomic Bulletin and Review*, *8*(1), 155–161. https://doi.org/10.3758/BF03196152

Duff, P. (2000). Repetition in foreign language classroom. In *The development of second and foreign language learning through classroom interaction* (pp. 109–138). Lawrence Erlbaum.

Farris-Trimble, A., & McMurray, B. (2013). Test–Retest Reliability of Eye Tracking in the Visual World Paradigm for the Study of Real-Time Spoken Word Recognition. *Journal of*

*Speech, Language, and Hearing Research*, *56*(4), 1328–1345.

Faul, F., Erdfelder, E., Buchner, A., & Lang, A. G. (2009). Statistical power analyses using G*Power 3.1: Tests for correlation and regression analyses. *Behavior Research Methods*, *41*(4), 1149–1160. https://doi.org/10.3758/BRM.41.4.1149

Faul, F., Erdfelder, E., Lang, A. G., & Buchner, A. (2007). G*Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, *39*(2), 175–191. https://doi.org/10.3758/BF03193146

Fernandes, T., Kolinsky, R., & Ventura, P. (2009). The metamorphosis of the statistical segmentation output: Lexicalization during artificial language learning. *Cognition*.

François, C., Cunillera, T., Garcia, E., Laine, M., & Rodriguez-Fornells, A. (2017). Neurophysiological evidence for the interplay of speech segmentation and word-referent mapping during novel word learning. *Neuropsychologia*, *98*, 56–67. https://doi.org/10.1016/j.neuropsychologia.2016.10.006

Gaskell, M. G., & Dumay, N. (2003). Lexical competition and the acquisition of novel words. *Cognition*, *89*(2), 105–132. https://doi.org/10.1016/S0010-0277(03)00070-2

Gathercole, S. E. (2006). Nonword repetition and word learning: The nature of the relationship. *Applied Psycholinguistics*, *27*, 513–543. https://doi.org/10.1017.S0142716406060383

Gathercole, S. E., & Conway, M. A. (1988). Exploring long-term modality effects: Vocalization leads to best retention. *Memory & Cognition*, *16*(2), 110–119. https://doi.org/10.3758/BF03213478

Graf Estes, K., Evans, J. L., Alibali, M. W., & Saffran, J. R. (2007). Can infants map meaning to

newly segmented words? Statistical segmentation and word learning. *Psychological*

*Science*, *18*(3), 254–260. https://doi.org/10.1111/j.1467-9280.2007.01885.x

Gupta, P. (2003). Examining the relationship between word learning, nonword repetition, and

immediate serial recall in adults. *Quarterly Journal of Experimental Psychology Section A:*

*Human Experimental Psychology*, *56 A*(7), 1213–1236.

https://doi.org/10.1080/02724980343000071

Gupta, P. (2008). A computational model of nonword repetition, immediate... - Google Scholar.

In A. Thorn & M. Page (Eds.), *Interactions between short-term and long-term memory in*

*the verbal domain* (pp. 120–147). Psychology Press.

Gupta, P., & Tisdale, J. (2009). Word learning, phonological short-term memory, phonotactic

probability and long-term memory: towards an integrated framework. *Philosophical*

*Transactions of the Royal Society of London. Series B, Biological Sciences*, *364*(1536),

3755–3771. https://doi.org/10.1098/rstb.2009.0132

Gussow, A. E., Kapnoula, E. C., & Molinaro, N. (2019). Any leftovers from a discarded

prediction? Evidence from eye-movements during sentence comprehension. *Language,*

*Cognition and Neuroscience*, 1–18. https://doi.org/10.1080/23273798.2019.1617887

Hay, J., Pelucchi, B., Estes, K., & Saffran, J. (2011). Linking sounds to meanings: Infant

statistical learning in a natural language. *Cognitive Psychology*, *63*(2), 93–106.

Hebb, D. (1961). Distinctive features of learning in the higher animal. *Brain Mechanisms and*

*Learning*, 37–46.

Hickok, G. (2012). Computational neuroanatomy of speech production. *Nature Reviews*

1    *Neuroscience*, *13*(2), 135–145. https://doi.org/10.1038/nrn3158

2    Hickok, G. (2014). The architecture of speech production and the role of the phoneme in speech

3         processing. *Language, Cognition and Neuroscience*, *29*(1), 2–20.

4         https://doi.org/10.1080/01690965.2013.834370

5    Hickok, G., Houde, J., & Rong, F. (2011). Sensorimotor Integration in Speech Processing:

6         Computational Basis and Neural Organization. *Neuron*, *69*(3), 407–422.

7         https://doi.org/10.1016/j.neuron.2011.01.019

8    Höhle, B., Fritzsche, T., Meß, K., Philipp, M., & Gafos, A. (2020). Only the right noise? Effects

9         of phonetic and visual input variability on 14-month-olds' minimal pair word learning.

10        *Developmental Science*, *23*(5). https://doi.org/10.1111/desc.12950

11   Hopkins, R. H., & Edwards, R. E. (1972). Pronunciation effects in recognition memory. *Journal*

12        *of Verbal Learning and Verbal Behavior*, *11*(4), 534–537. https://doi.org/10.1016/S0022-

13        5371(72)80036-7

14   Hopman, E. W. M., & MacDonald, M. C. (2018). Production Practice During Language

15        Learning Improves Comprehension. *Psychological Science*, *29*(6), 961–971.

16        https://doi.org/10.1177/0956797618754486

17   Horst, J. S., & Hout, M. C. (2016). The Novel Object and Unusual Name (NOUN) Database: A

18        collection of novel images for use in experimental research. *Behavior Research Methods*,

19        *48*(4), 1393–1409. https://doi.org/10.3758/s13428-015-0647-3

20   Houston, D. M., & Jusczyk, P. W. (2000). The role of talker-specific information in word

21        segmentation by infants Speech-Language Environments of Children with Cochlear

Implants View project. *Article in Journal of Experimental Psychology Human Perception &*
*Performance*, *26*(5), 1570. https://doi.org/10.1037/0096-1523.26.5.1570

Kadota, S. (2019). *Shadowing as a practice in second language acquisition : connecting inputs*
*and outputs* (S. Kadota (ed.)). Routledge.

Kapnoula, E. C., & McMurray, B. (2016). Newly learned word-forms are abstract and integrated
immediately after acquisition. *Psychonomic Bulletin and Review*, *23*(2), 491–499.
https://doi.org/10.3758/s13423-015-0897-1

Kapnoula, E. C., Packard, S., Gupta, P., & McMurray, B. (2015). Immediate lexical integration
of novel word forms. *Cognition*, *134*, 85–99.
https://doi.org/10.1016/j.cognition.2014.09.007

Kapnoula, E. C., & Samuel, A. G. (under review). *Wait long and prosper! Delaying production*
*during training boosts word learning*.

Kapnoula, E. C., & Samuel, A. G. (2019). Voices in the mental lexicon: Words carry indexical
information that can affect access to their meaning. *Journal of Memory and Language*, *107*,
111–127. https://doi.org/10.1016/J.JML.2019.05.001

Karpicke, J., & Roediger, H. L. (2008). The critical importance of retrieval for learning. *Science*,
*319*(5865), 966–968. https://doi.org/10.1126/science.1152408

Kaushanskaya, M., & Yoo, J. (2011). Rehearsal effects in adult word learning. *Language and*
*Cognitive Processes*, *26*(1), 121–148. https://doi.org/10.1080/01690965.2010.486579

Krashen, S. D. (1985). *The Input Hypothesis: Issues and Implications.* Longman.

Leach, L., & Samuel, A. G. (2007). Lexical configuration and lexical engagement: when adults learn new words. *Cognitive Psychology*, *55*(4), 306–353. https://doi.org/10.1016/j.cogpsych.2007.01.001

Lightbown, P. M., & Spada, N. (2013). *How Languages are Learned*. Oxford University Press.

Loftus, G. R., & Masson, M. E. J. (1994). Using confidence intervals in within-subject designs. *Psychonomic Bulletin & Review*, *1*(4), 476–490. https://doi.org/10.3758/BF03210951

López Assef, B., Desmeules-Trudel, F., Bernard, A., & Zamuner, T. S. (2021). A Shift in the Direction of the Production Effect in Children Aged 2–6 Years. *Child Development*, *92*(6), 2447–2464. https://doi.org/10.1111/cdev.13618

MacDonald, M. C., Pearlmutter, N. J., & Seidenberg, M. S. (1994). The lexical nature of syntactic ambiguity resolution. *Psychological Review*, *101*(4), 676–703. https://doi.org/10.1037//0033-295X.101.4.676

MacDonald, P., & MacLeod, C. (1998). The influence of attention at encoding on direct and indirect remembering. *Acta Psychologica*, *98*(2–3), 291–310.

MacLeod, C. M., & Bodner, G. E. (2017). The Production Effect in Memory. *Current Directions in Psychological Science*, *26*(4), 390–395. https://doi.org/10.1177/0963721417691356

MacLeod, C. M., Gopie, N., Hourihan, K. L., Neary, K. R., & Ozubko, J. D. (2010). The production effect: Delineation of a phenomenon. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *36*(3), 671–685. https://doi.org/10.1037/a0018785

Magnuson, J. S. (2019). Fixations in the visual world paradigm: where, when, why? *Journal of Cultural Cognitive Science*, *3*(2), 113–139. https://doi.org/10.1007/s41809-019-00035-3

Mattys, S. L., & Baddeley, A. (2019). Working memory and second language accent acquisition. *Applied Cognitive Psychology*, *33*(6), 1113–1123. https://doi.org/10.1002/acp.3554

McMurray, B., Clayards, M. A., Tanenhaus, M. K., & Aslin, R. N. (2008). Tracking the time course of phonetic cue integration during spoken word recognition. *Psychonomic Bulletin & Review*, *15*(6), 1064–1071. https://doi.org/10.3758/PBR.15.6.1064

McMurray, B., Kapnoula, E. C. E. C., & Gaskell, M. G. (2016). Learning and integration of new word-forms: Consolidation, pruning and the emergence of automaticity. In M. G. Gaskell & J. Mirković (Eds.), *Speech Perception and Spoken Word Recognition.* (Psychology). https://doi.org/10.4324/9781315772110

McMurray, B., Samelson, V. M., Lee, S. H., & Tomblin, J. B. (2010). Individual differences in online spoken word recognition : Implications for SLI. *Cognitive Psychology*, *60*(1), 1–39. https://doi.org/10.1016/j.cogpsych.2009.06.003

McMurray, B., Tanenhaus, M. K., & Aslin, R. N. (2002). Gradient effects of within-category phonetic variation on lexical access. *Cognition*, *86*(2), B33–B42. https://doi.org/10.1016/S0010-0277(02)00157-9

McRae, K., Spivey-Knowlton, M. J., & Tanenhaus, M. K. (1998). Modeling the Influence of Thematic Fit (and Other Constraints) in On-line Sentence Comprehension. *Journal of Memory and Language*, *38*(3), 283–312. https://doi.org/10.1006/jmla.1997.2543

Mirman, D., Dixon, J. A., & Magnuson, J. S. (2008). Statistical and computational models of the visual world paradigm: Growth curves and individual differences. *Journal of Memory and Language*, *59*(4), 475–494. https://doi.org/10.1016/J.JML.2007.11.006

Mirman, D., Magnuson, J. S., Estes, K. G., & Dixon, J. A. (2008). The link between statistical segmentation and word learning in adults. *Cognition*, *108*(1), 271–280.

Morey, R. D. (2008). Confidence Intervals from Normalized Data: A correction to Cousineau (2005). *Tutorials in Quantitative Methods for Psychology*, *4*(2), 61–64. https://doi.org/10.20982/tqmp.04.2.p061

Nora, A., Renvall, H., Kim, J.-Y., Service, E., & Salmelin, R. (2015). Distinct Effects of Memory Retrieval and Articulatory Preparation when Learning and Accessing New Word Forms. *PLOS ONE*, *10*(5), e0126652. https://doi.org/10.1371/journal.pone.0126652

Norris, D., McQueen, J. M., & Cutler, A. (2003). Perceptual learning in speech. *Cognitive Psychology*, *47*(2), 204–238. https://doi.org/10.1016/S0010-0285(03)00006-9

Nuttall, H. E., Kennedy-Higgins, D., Hogan, J., Devlin, J. T., & Adank, P. (2016). The effect of speech distortion on the excitability of articulatory motor cortex. *NeuroImage*, *128*, 218–226. https://doi.org/10.1016/j.neuroimage.2015.12.038

Ozubko, J. D., & Macleod, C. M. (2010). The Production Effect in Memory: Evidence That Distinctiveness Underlies the Benefit The Production Effect View project. *Article in Journal of Experimental Psychology Learning Memory and Cognition*. https://doi.org/10.1037/a0020604

Ozubko, J. D., Major, J., & Macleod, C. M. (2014). Remembered study mode: Support for the distinctiveness account of the production effect. *Memory*, *22*(5), 509–524. https://doi.org/10.1080/09658211.2013.800554

Page, M., & Norris, D. (1998). The primacy model: a new model of immediate serial recall.

*Psychological Review*, *105*(4), 761–781.

Page, M., & Norris, D. (2008). Is there a common mechanism underlying word-form learning and the Hebb repetition effect? Experimental data and a modelling framework. In A. Thorn & M. Page (Eds.), *Interactions Between Short-Term and Long-Term Memory in the Verbal Domain* (pp. 148–168). Psychology Press.

Page, M., & Norris, D. (2009). *A model linking immediate serial recall, the Hebb repetition effect and the learning of phonological word forms. 364*(1536), 3737–3753. https://doi.org/10.1098/rstb.2009.0173

Payton, K. L., Uchanski, R. M., & Braida, L. D. (1994). Intelligibility of conversational and clear speech in noise and reverberation for listeners with normal and impaired hearing. *Journal of the Acoustical Society of America*, *95*(3), 1581–1592. https://doi.org/10.1121/1.408545

Protopapas, A. (2007). CheckVocal: a program to facilitate checking the accuracy and response time of vocal responses from DMDX. *Behavior Research Methods*, *39*, 859–862. https://doi.org/10.3758/BF03192979

Richtsmeier, P. T., Gerken, L. A., Goffman, L., & Hogan, T. (2009). Statistical frequency in perception affects children's lexical production. *Cognition*, *111*(3), 372–377. https://doi.org/10.1016/j.cognition.2009.02.009

Riley, K. G., & McGregor, K. K. (2012). Noise hampers children's expressive word learning. *Language, Speech, and Hearing Services in Schools*, *43*(3), 325–337. https://doi.org/10.1044/0161-1461(2012/11-0053)

Rodriguez-Fornells, A., Cunillera, T., Mestres-Missé, A., & De Diego-Balaguer, R. (2009).

Neurophysiological mechanisms involved in language learning in adults. In *Philosophical Transactions of the Royal Society B: Biological Sciences* (Vol. 364, Issue 1536, pp. 3711–3735). Royal Society. https://doi.org/10.1098/rstb.2009.0130

Roediger, H. L., & Karpicke, J. D. (2006). Test-enhanced learning: Taking memory tests improves long-term retention. *Psychological Science*, *17*(3), 249–255. https://doi.org/10.1111/j.1467-9280.2006.01693.x

Rost, G. C., & McMurray, B. (2009). Speaker variability augments phonological processing in early word learning. *Developmental Science*, *12*(2), 339–349. https://doi.org/10.1111/j.1467-7687.2008.00786.x

Rost, G. C., & McMurray, B. (2010). Finding the signal by adding noise: The role of noncontrastive phonetic variability in early word learning. *Infancy*, *15*(6), 608–635. https://doi.org/10.1111/j.1532-7078.2010.00033.x

Salverda, A. P., & Tanenhaus, M. K. (2017). The visual world paradigm. In A. M. B. de Groot & P. Hagoort (Eds.), *Research Methods in Psycholinguistics and the Neurobiology of Language*. Wiley.

Samuel, A. G. (1996). Does lexical information influence the perceptual restoration of phonemes? *Journal of Experimental Psychology: General*, *125*(1), 28.

Samuel, A. G., & Larraza, S. (2015). Does listening to non-native speech impair speech perception? *Journal of Memory and Language*, *81*, 51–71.

Scheepers, C., Keller, F., & Lapata, M. (2008). Evidence for Serial Coercion: A Time Course Analysis Using the Visual-World Paradigm. *Cognitive Psychology*, *56*(1), 1–29.

1      https://doi.org/10.1016/j.cogpsych.2006.10.001

2   Seedorff, M., Oleson, J., Cavanaugh, J., & McMurray, B. (2017). Eyetracking Analysis in R. *R*

3      *Package Version 0.1.15*.

4   Seedorff, M., Oleson, J., & McMurray, B. (2018). Detecting when timeseries differ: Using the

5      Bootstrapped Differences of Timeseries (BDOTS) to analyze Visual World Paradigm data

6      (and more). *Journal of Memory and Language*, *102*, 55–67.

7      https://doi.org/10.1016/J.JML.2018.05.004

8   Singh, L. (2008). Influences of high and low variability on infant word recognition. *Cognition*,

9      *106*(2), 833–870.

10  Smiljanić, R., & Bradlow, A. R. (2009). Speaking and hearing clearly: Talker and listener factors

11     in speaking style changes. In *Linguistics and Language Compass* (Vol. 3, Issue 1, pp. 236–

12     264). Blackwell Publishing Inc. https://doi.org/10.1111/j.1749-818X.2008.00112.x

13  Tanenhaus, M. K., Spivey-Knowlton, M. J., Eberhard, K. M., & Sedivy, J. C. (1995). Integration

14     of visual and linguistic information in spoken language comprehension. *Science*, *268*(5217),

15     1632–1634.

16  Tanenhaus, M. K., & Trueswell, J. (1995). *Sentence comprehension.*

17  Thorn, A., & Page, M. (2008). *Interactions Between Short-Term and Long-Term Memory in the*

18     *Verbal Domain*. Psychology Press.

19  Trueswell, J. C. (1996). The Role of Lexical Frequency in Syntactic Ambiguity Resolution.

20     *Journal of Memory and Language*, *35*(4), 566–585. https://doi.org/10.1006/jmla.1996.0030

1    Warren, R. (1970). Perceptual restoration of missing speech sounds. *Science*, *167*(3917), 392–

2        393.

3    Yee, E., & Sedivy, J. C. (2006). Eye movements to pictures reveal transient semantic activation

4        during spoken word recognition. *Journal of Experimental Psychology: Learning, Memory,*

5        *and Cognition*, *32*(1), 1–14. https://doi.org/10.1037/0278-7393.32.1.1

6    Yu, C., & Smith, L. (2007). Rapid word learning under uncertainty via cross-situational

7        statistics. *Psychological Science*.

8    Zamuner, T. S., Morin-Lessard, E., & Strahm, S. (2016). Spoken word recognition of novel

9        words, either produced or only heard during learning. *Journal of Memory and Language*,

10        *89*, 55–67.

11    Zamuner, T. S., Strahm, S., Morin-Lessard, E., & Page, M. (2018). Reverse production effect:

12        children recognize novel words better when they are heard rather than produced.

13        *Developmental Science*, *21*(4), e12636. https://doi.org/10.1111/desc.12636

14

1 **Appendix**

2 Table A1. *Means and standard deviations for accuracy and reaction times (RT) at test by*
3 *Condition, Phase, and Experiment*

| | Accuracy | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Phase 1 | | | Phase 2 | | |
| | *N* | *M* | *SD* | *N* | *M* | *SD* |
| Experiment 1 | | | | | | |
| Perception-Only | 40 | 97.2% | 6.2% | 40 | 99.1% | 3.3% |
| Production | 40 | 97.0% | 9.3% | 40 | 99.2% | 2.5% |
| Experiment 2 | | | | | | |
| Perception-Only | 41 | 94.2% | 10.6% | 41 | 99.4% | 1.9% |
| Production | 41 | 93.1% | 10.3% | 41 | 99.1% | 4.1% |
| Experiment 3 | | | | | | |
| Perception-Only | 41 | 94.2% | 10.1% | 41 | 99.5% | 1.6% |
| Production | 41 | 95.9% | 8.1% | 41 | 99.8% | 1.0% |
| | Reaction times (in ms) | | | | | |
| | Phase 1 | | | Phase 2 | | |
| | *N* | *M* | *SD* | *N* | *M* | *SD* |
| Experiment 1 | | | | | | |
| Perception-Only | 40 | 1,328 | 288 | 40 | 1,097 | 290 |
| Production | 40 | 1,328 | 311 | 40 | 1,088 | 289 |

Experiment 2

| | | | | | | |
|---|---|---|---|---|---|---|
| Perception-Only | 41 | 1,419 | 336 | 41 | 1,154 | 206 |
| Production | 41 | 1,560 | 602 | 41 | 1,189 | 220 |

Experiment 3

| | | | | | | |
|---|---|---|---|---|---|---|
| Perception-Only | 41 | 1,416 | 328 | 41 | 1,122 | 186 |
| Production | 41 | 1,375 | 301 | 41 | 1,127 | 250 |

1
2

1    Table A2. *Full 3×2×2×2 ANOVA results*

| Predictor | Sum of Squares | df | Mean Square | F | p | partial $\eta^2$ |
|---|---|---|---|---|---|---|
| Experiment | 7.263 | 2 | 3.631 | 4.376 | .015 | 0.071 |
| Condition | 0.002 | 1 | 0.002 | 0.013 | .909 | <.01 |
| Phase | 16.095 | 1 | 16.095 | 45.665 | <.001 | 0.286 |
| TimeWindow | 479.539 | 1 | 479.539 | 1044.344 | <.001 | 0.902 |
| Experiment × Condition | 0.045 | 2 | 0.022 | 0.152 | .859 | 0.003 |
| Experiment × Phase | 0.488 | 2 | 0.244 | 0.692 | .503 | 0.012 |
| Experiment × TimeWindow | 5.583 | 2 | 2.791 | 6.079 | .003 | 0.096 |
| Condition × Phase | 0.368 | 1 | 0.368 | 3.107 | .081 | 0.027 |
| Condition × TimeWindow | 0.048 | 1 | 0.048 | 0.497 | .482 | 0.004 |
| Phase × TimeWindow | 7.319 | 1 | 7.319 | 46.086 | <.001 | 0.288 |
| Experiment × Condition × Phase | 0.036 | 2 | 0.018 | 0.150 | .861 | 0.003 |
| Experiment × Condition × TimeWindow | 0.083 | 2 | 0.042 | 0.433 | .650 | 0.008 |
| Condition × Phase × TimeWindow | 1.091 | 1 | 1.091 | 13.481 | <.001 | 0.106 |
| Experiment × Condition × Phase × TimeWindow | 0.017 | 2 | 0.009 | 0.107 | .899 | 0.002 |
| Error | | 114 | | | | |

2
3

69

1    Table A3. *Full results of post-hoc comparisons (Bonferroni-corrected)*

| Phase | TW | Mean Difference (Perception-Only – Production) | SE | p | Lower Bound | Upper Bound |
|---|---|---|---|---|---|---|
| 1 | Early | .017 | .024 | .465 | -.030 | .064 |
| 1 | Late | -.090 | .060 | .131 | -.209 | .028 |
| 2 | Early | -.040 | .021 | .059 | -.082 | .002 |
| 2 | Late | .125 | .054 | .022 | .019 | .232 |

2
3

1
2   *Figure A1*. Average differences of looks to the target (empirical-logit-transformed) per testing
3   phase (Phase 1/Phase 2), and time window (TW) within trial (early: TW1/late: TW2) across
4   experiments. Error bars show 95% confidence intervals of the mean differences. *Note*. This
5   alternative visualization of the ANOVA results presents empirical-logit-transformed data, as
6   used in the statistical analyses.
7