

1 TITLE

2 Conflict monitoring in speech processing: An fMRI study of error detection in speech  
3 production and perception

4

5 ABBREVIATED TITLE

6 Error detection in speech production and perception

7

8 Hanna S. Gauvin<sup>1</sup>, Wouter De Baene<sup>1,2</sup>, Marcel Brass<sup>1</sup> and Robert J. Hartsuiker<sup>1</sup>

9 1 Department of Experimental Psychology, Ghent University, 9000 Ghent, Belgium

10 2 Department of Cognitive Neuropsychology, Tilburg University, 5000 LE Tilburg, The

11 Netherlands

12

13 Correspondence:

14 Hanna S. Gauvin

15 Department of Experimental Psychology

16 Ghent University

17 Henri Dunantlaan 2

18 9000 Ghent, Belgium

19 Hanna.Gauvin@ugent.be

20

21 Conflict of Interest: The authors declare no competing financial interests.

22 Acknowledgements: This research was supported by the Fund for Scientific Research

23 Flanders (FWO), project "Internal self-monitoring: speech perception or forward

24 models?" no. G.0335.11N

25

## Abstract

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

To minimize the number of errors in speech, and thereby facilitate communication, speech is monitored before articulation. It is, however, unclear at which level during speech production monitoring takes place, and what mechanisms are used to detect and correct errors. The present study investigated whether internal verbal monitoring takes place through the speech perception system, as proposed by perception-based theories of speech monitoring, or whether mechanisms independent of perception are applied, as proposed by production-based theories of speech monitoring. With the use of fMRI during a tongue twister task we observed that error detection in internal speech during noise-masked overt speech production and error detection in speech perception both recruit the same neural network, which includes pre-supplementary motor area (pre-SMA), dorsal anterior cingulate cortex (dACC), anterior insula (AI), and inferior frontal gyrus (IFG). Although production and perception recruit similar areas, as proposed by perception-based accounts, we did not find activation in superior temporal areas (which are typically associated with speech perception) during internal speech monitoring in speech production as hypothesized by these accounts. On the contrary, results are highly compatible with a domain general approach to speech monitoring, by which internal speech monitoring takes place through detection of conflict between response options, which is subsequently resolved by a domain general executive center (e.g., the ACC).

48 Keywords: self-monitoring, speech production, speech perception, conflict  
49 monitoring, verbal monitoring, fMRI

50

## 1. Introduction

51 In the domain of language production there is consensus about the existence of an  
52 internal speech monitoring system, which monitors speech before production, in  
53 addition to an external monitoring system (i.e., hearing one's own speech). Evidence  
54 for an internal monitoring system comes from research showing extremely fast self-  
55 corrections for which the external monitoring system would simply be too slow  
56 (Levelt, 1989; Blackmer & Mitton, 1991; Hartsuiker & Kolk 2001), the report of errors  
57 when silently performing a speech task (Oppenheim & Dell, 2008), and the report of  
58 errors when overt speech is masked by loud noise (Lackner & Tuller, 1979; Postma &  
59 Kolk, 1992). However, there is currently no consensus on the underlying nature of  
60 such an internal speech error monitoring mechanism. In a review of verbal  
61 monitoring models, Postma (2000) discusses eleven possible locations during the  
62 process of speaking at which monitoring has been proposed to take place. Most of  
63 the proposed models are directed at monitoring internal speech. Additionally,  
64 external speech can be monitored via perception of the speech and via perception of  
65 the articulators and muscles (proprioceptive feedback) (Abbs & Gracco, 1983; Abbs  
66 et al., 1984; Siegenthaler & Hochberg, 1965).

67 Presently there are roughly three classes of theories on monitoring internal  
68 speech: perception-based accounts (Perceptual Loop Theory, Hartsuiker & Kolk,  
69 2001; Levelt, 1989; Indefrey, 2011), production-based accounts (Local Monitors,  
70 Laver 1980; Conflict Monitors, Nozari et al., 2011), and forward modeling accounts  
71 (e.g., Hickok, 2012; Pickering & Garrod, 2013; Tourville & Guenther, 2011). In the  
72 current study we investigated whether the neural structures involved in verbal  
73 monitoring lend support for any of these three classes of theories. Below we will first

74 outline the theories in more detail, including their neuro-anatomical hypotheses,  
75 after which we outline how the fMRI data can be used to dissociate these theories.

76 Perception-based theories assume that internal speech monitoring takes  
77 place in the speech perception system, during both production and perception.  
78 During production, the phonetic plan is sent directly to the perception system (i.e.,  
79 before articulation) for internal monitoring. Essentially the same monitoring  
80 mechanism would be used for both internal and external monitoring according to  
81 the perception-based theories of monitoring, with internal monitoring using part of  
82 the external monitoring route. Monitoring your own internal and external speech  
83 and monitoring someone else's speech, all rely on the perception system, for which  
84 the superior temporal gyrus is the main neural substrate (e.g. Price, 2012).

85 Production-based theories do not necessarily assume the same monitoring  
86 system for production and perception, and assume that internal monitoring during  
87 production takes place independently of speech perception systems. A recently  
88 proposed production monitoring account uses conflict within the production system  
89 as a basis for monitoring (Nozari et al., 2011). Analogous to domain-general theories  
90 of error detection (e.g., Botvinick et al., 2001; Yeung et al., 2004), monitoring rather  
91 takes place through detection of conflict between response options, which is  
92 subsequently resolved by a domain-general cognitive control unit located in the  
93 Anterior Cingulate Cortex (ACC). Support for these domain-general theories comes  
94 from the vast, and increasing, body of literature in which in response to conflict an  
95 ERN component is found in EEG studies and ACC activation in fMRI studies. Source  
96 localization has traced the ERN component to originate from the ACC (e.g. Van Veen

97 & Carter, 2002; Herrmann et al., 2004). The ERN component and ACC activation are  
98 similar across cognitive domains, suggesting the existence of a domain general  
99 conflict response. For a more detailed overview of the conflict monitoring literature  
100 in relation to the ERN and ACC, see for instance Van Veen and Carter (2006).

101 Another type of production-based monitoring that has been proposed is  
102 monitoring via forward models. Forward modeling accounts of speech monitoring  
103 assume that during production a prediction, or forward model, of the expected  
104 outcome is made. The actual outcome is compared to the predicted outcome, and if  
105 a mismatch between these two is detected, a corrective signal arises. Forward model  
106 theories of speech production are supported by the observation of auditory  
107 response suppression during speech production; based on the prediction of the  
108 sensory feedback of the upcoming event, the sensory cortex is inhibited. When the  
109 sensory feedback is not in accordance with the prediction, an increase in activation is  
110 observed, which might function as a corrective signal (Curio et al., 2000; Heinks-  
111 Maldonado et al., 2005, 2006; Numminen et al., 1999; Eliades & Wang, 2003, 2005).  
112 Direct evidence in support of forward models during speech production comes from  
113 a series of MEG experiments, showing context dependent activation changes in the  
114 auditory cortex in response to imagined speech production (so in the absence of  
115 actual auditory stimulation) at a same time frame as observed after normal speech  
116 production (Tian & Poeppel, 2010, 2013). Most forward model theories rely on  
117 sensory feedback for monitoring. Consequently, internal speech monitoring, which is  
118 investigated in the current study, is outside the scope of these theories. However,  
119 Pickering and Garrod's forward model theory (2013, 2014) does make predictions

120 about monitoring in internal speech production and speech perception. According to  
121 their theory, both during production and perception, predictions are made and  
122 compared to the actual utterance. These comparisons are made in a comparator,  
123 which is a speech-modality (production / perception) independent system. So a  
124 difference between correct and incorrect sentences is expected to lead to  
125 differences in activation in the comparator, which is separate of the perception  
126 system. However, no anatomical predictions are made with respect to this  
127 comparator in Pickering and Garrod's forward model theory.

128         Because production- and perception-based monitoring theories make distinct  
129 predictions about the functional neuroanatomy of speech monitoring, fMRI is a  
130 useful tool to distinguish between these competing theories. Perception-based  
131 monitoring accounts assume that, as the bilateral superior temporal gyri (STG) are  
132 involved in monitoring external speech, internal speech must be monitored via (a  
133 subpart of) the same neuronal structures (Indefrey, 2011). So if (pre-verbal) internal  
134 monitoring is perception-based, we expect superior temporal gyrus (STG) activation  
135 for error detection in both production and perception, even when auditory feedback  
136 is unavailable during production. If monitoring is production-based, however, we  
137 expect to find error monitoring independently of perceptual areas during  
138 production. Production-based monitoring accounts predict activation in areas  
139 associated with subcomponents of the production process, as well as domain-  
140 general areas associated with conflict monitoring in the medial frontal areas, such as  
141 the ACC. In the experiment reported below, we compared internal speech  
142 monitoring during production with external speech monitoring during perception, in  
143 order to investigate whether all monitoring is indeed performed by the perceptual

144 system (as proposed in perception-based theories of monitoring such as the  
145 perceptual loop theory by Levelt, 1989; Indefrey, 2011) or whether monitoring is  
146 performed independently of the perceptual system (as proposed by production-  
147 based theories of monitoring such as the conflict monitoring theory by Nozari et al.,  
148 2011, and the forward model theory by Pickering & Garrod, 2013, 2014).

149         At this moment there are no publications that describe the neuronal  
150 structures involved in internal and external verbal monitoring and their differences.  
151 Only few studies have applied fMRI to investigate internal speech monitoring and  
152 none have compared monitoring in speech production with monitoring in speech  
153 perception. Monitoring of external speech has been investigated predominantly by  
154 manipulating acoustic feedback in the dimensions of frequency or time (McGuire et  
155 al, 1996; Hirano et al, 1997; Fu et al, 2006; Christoffels et al, 2007; Tourville et al,  
156 2008). Perception of altered feedback led to increased activation in the superior  
157 temporal lobe compared to unaltered feedback. Note that these are externally  
158 induced 'errors'; the participant made no error during production, but via  
159 manipulation of the feedback the perception of the speech is changed. There is only  
160 one published fMRI study on error production in language processing that targets  
161 errors made by the producer herself (Abel et al, 2009). In this experiment,  
162 participants overtly named pictures during scanning, and resulting activations during  
163 correct production, incorrect production, and a rest baseline were compared. This  
164 study found increased activations during error production in the ACC, prefrontal and  
165 premotor regions, basal ganglia, thalamus, SMA and precentral gyrus. This  
166 experiment had, however, several limitations: few errors were made, and the



167 reported errors were not very naturalistic as some were merely errors against the  
168 instructed label for each picture (e.g., call this picture ‘flower’ and the participant  
169 responds with ‘sunflower’). There have been no published studies, to the best of our  
170 knowledge, with a direct comparison between fMRI data of speech error detection  
171 in production and perception.

172         The current study therefore aims to investigate the neural underpinnings of  
173 internal verbal monitoring and external verbal monitoring. These data can be used to  
174 distinguish between several highly influential theoretical models of verbal  
175 monitoring, as these theories have neuroanatomically specific predictions. Below we  
176 outline the experimental setup, and we discuss the hypotheses the monitoring  
177 models make with respect to the neuronal functional data.

178         Participants performed a tongue twister task in which they repeated tongue  
179 twister sentences, or listened to a recording of a tongue twister repetition, after  
180 which they judged the repetition on correctness. The percentage of errors in the  
181 perception condition was matched to the number of errors in the production  
182 condition, to allow for a comparison of the areas involved in error detection in both  
183 modalities. In order to test the involvement of the speech perception system in  
184 internal speech monitoring during production, normal feedback was precluded, as  
185 auditory feedback would necessarily involve external monitoring via the speech  
186 perception system. Perceptual-based monitoring is only supported if we find a role  
187 for the perception system in internal speech monitoring.

188         The theories of internal verbal monitoring make the following predictions  
189 with respect to the neural structures involved in verbal monitoring; perception-  
190 based monitoring assumes a major role for the auditory perceptual system, located

191 in the bilateral STG, during both internal and external verbal monitoring. Production-  
192 based monitoring assumes involvement of a domain general monitoring mechanism  
193 located in the ACC, and crucially it assumes no role for the auditory perceptual  
194 system during internal verbal monitoring.

195

## 196 2. Materials and Methods

197

### 198 2.1 Participants

199 Twenty-four participants were recruited from Ghent University, of which 3  
200 were discarded: one due to excessive motion, one because of too many errors  
201 (>80%) and one due to too few errors (<10%). Final analyses included 21 participants  
202 (15 females, 6 males; mean age: 21, ranging from 19 to 30). All reported to be native  
203 speakers of Dutch, have no dyslexia or other speech or language impairments, no  
204 hearing problems, and normal or corrected-to-normal vision. No subject had a  
205 history of neurological, psychiatric, or major medical disorder as assessed by a pre-  
206 scanning questionnaire. All subjects were right handed as assessed by the Edinburgh  
207 Handedness inventory (Oldfield, 1971) ( $n=21$ , EHI score  $M=90.4$ ,  $SD=15.8$ , range =  
208 41.2 to 100, mode = 100). A monetary reward was received for participation. The  
209 study was approved by the local ethical committee of Ghent University's Medical  
210 Department and was conducted in accordance to the Declaration of Helsinki.

211

### 212 2.2. Stimulus material and task design

213 Stimuli were selected on the basis of a pilot study in which 56 tongue twister  
214 sentences were tested. For the production condition sentences were selected that

215 elicited linguistic errors (phonological slips, semantic substitutions, and syntactic  
216 errors). From these, we selected 17 sentences with high error production rates (30%  
217 - 60% of repetitions contained an error). An additional 5 sentences were selected  
218 that were relatively easy (10% - 25% of repetitions contained an error), in order to  
219 prevent discouragement among the participants. Each sentence was presented 3  
220 times per condition. We used tongue twister sentences of the type 'A proper copper  
221 coffee pot', and 'How can a clam cram in a clean cream can?'. A full overview of the  
222 tongue twister sentences is provided in Appendix A. The sentences were selected on  
223 the basis of the errors they elicited, and were not matched for frequency or length<sup>1</sup>.  
224 Audio files for the instruction phase were created in which these sentences were  
225 clearly pronounced at a normal speech rate by a male native speaker of Dutch. These  
226 audio files were presented together with a visual presentation of the tongue twister  
227 at the beginning of each trial.

228 For the perception condition 22 different tongue twister sentences were  
229 selected in order to decrease repetition effects and minimize attention loss. As in the  
230 production condition, each sentence was presented 3 times during the perception  
231 condition. Actual recordings of 4 female participants producing the sentences in the  
232 pilot study, correctly and incorrectly, were used as auditory stimuli. The errors  
233 selected for the perception condition highly resembled those produced in the  
234 production condition, and were all linguistic errors. Pitch was adjusted (increased  
235 with 50 or 20 Hz) for 3 of the 4 participants to facilitate auditory perception in the

---

<sup>1</sup> Note that very long sentences that were difficult to remember were not selected from the pilot study. Only sentences that elicited linguistic errors were selected.

236 scanner. Experiments were created in E-prime 2.0 (Psychology Software Tools, Inc,  
237 Pittsburg, PA).

238 Before entering the scanner participants were briefed on the task. After  
239 entering the scanner the participants again received instructions on the production  
240 task, followed by a familiarization phase and successively the actual experiment.  
241 Participants were instructed to speak normally, while keeping their heads fixed. To  
242 minimize movements, foam pads were placed between the head and head coil. Once  
243 the participant was set up to enter the scanner bore, the participant was asked to  
244 speak, and once again the experimenters stressed to the participants to speak  
245 normally and avoid any head movements, as motion artifacts are often observed  
246 with speech production during acquisition (see below). During the production  
247 condition the experimenter scored the number of incorrectly produced sentences,  
248 which allowed for an error percentage match with the perception condition. After  
249 completing the production condition, participants received instructions for the  
250 perception condition, followed by a familiarization phase and consecutively the  
251 perception condition of the tongue twister task. The total duration of the  
252 experiment was approximately 45 minutes.

253

#### 254 2.2.1 Production Condition

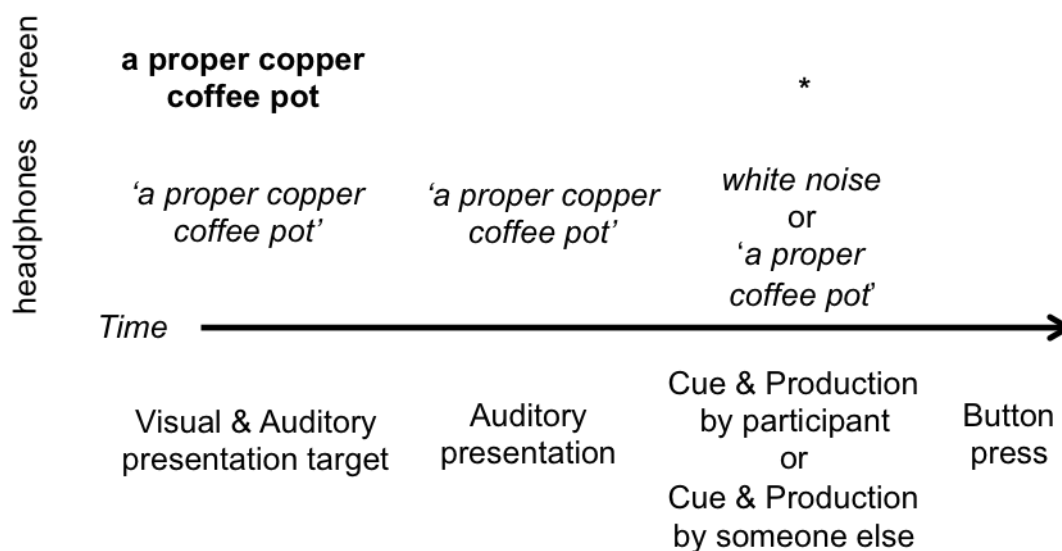
255 Each trial consisted of a visual presentation of the target sentence with a  
256 simultaneous auditory presentation, followed by a blank screen and after 200 ms a  
257 repetition of the auditory presentation. After a pause of 250 ms a visual cue was  
258 presented (\*) to signal to the participant to start producing the target sentence.  
259 After producing the sentence the participant pushed a button to indicate whether

260 the sentence was correct (right hand) or incorrect (left hand). From cue onset until a  
261 correctness judgment was made, after which the cue disappeared, the participant  
262 heard a white noise at maximum volume over the headphones to mask auditory  
263 feedback. An illustration of this task is provided in Figure 1. After a familiarization  
264 phase of 3 trials, three target blocks were presented that each consisted of the 22  
265 tongue twister sentences in random order. Between trials a varying ISI of between  
266 1250 and 5500 ms occurred (mean 2867 ms).

267

### 268 2.2.2 Perception Condition

269 In the perception trials the participants were presented with a visual  
270 presentation of the target sentence with simultaneous auditory presentation,  
271 exactly as in the production condition. After a pause of 200 ms the participants  
272 heard a recording of a person producing the sentence. The participant pushed a  
273 button to indicate whether the sentence was repeated correctly (right hand) or  
274 incorrectly (left hand). An illustration of this task is provided in Figure 1. After a  
275 familiarization phase of 3 trials, three target blocks were presented that each  
276 consisted of the 22 tongue twister sentences in random order. Between trials a  
277 varying ISI of between 1250 and 5500 ms occurred (mean 2867 ms), similar to the  
278 production condition. We constructed 8 versions of the perception condition with  
279 different error rates, ranging from 10% to 45% errors (with 5% intervals), to  
280 approximate the number of errors produced in the production condition.



281

282 Figure 1. Overview of the tongue twister task. The production and perception  
 283 condition only differ in the production part of the sequence; in the production  
 284 condition a white noise is presented over headphones during which the participant  
 285 produces the sentence, while in the perception condition a pre-recorded production  
 286 of the tongue twister sentence is played over headphones.

287

### 288 2.3 Scanning procedure

289 Images were collected with a 3T Magnetom Trio MRI scanner system  
 290 (Siemens Medical Systems, Erlangen, Germany), using a standard 32-channel radio-  
 291 frequency head coil. A 3D high-resolution anatomical image of the whole brain was  
 292 acquired first, for co-registration with the functional images using a T1-weighted 3D  
 293 MPRAGE sequence (TR = 2530 ms, TE = 2.58 ms, TI = 1100 ms, acquisition matrix =  
 294 256 × 256 × 176, sagittal FOV = 220 mm, flip angle = 7°, voxel size = .90 × .86 × .86  
 295 mm<sup>3</sup> (resized to 1 × 1 × 1 mm<sup>3</sup>)). Whole brain functional images were collected  
 296 using a T2\*-weighted EPI sequence, sensitive to BOLD contrast (TR= 2000 ms, TE=28  
 297 ms, image matrix=64×64, FOV=224 mm, flip angle = 80°, slice thickness = 3 mm,

298 distance factor = 17%, voxel size  $3.5 \times 3.5 \times 3.51 \text{ mm}^3$ , 34 axial slices). Specific care  
299 was taken to ensure that frontal areas and (near) complete cerebellum were  
300 included in the imaging volume. A varying number of images were acquired per run  
301 due to the self-paced ending of trials. In the production condition the number of  
302 images per run ranged from 450-528, in the perception condition the number of  
303 images per run ranged from 334-382.

304 Participants went head first and supine into the magnetic bore. They were  
305 instructed to speak normally but to avoid movements of their heads in order to  
306 avoid motion artifacts. Foam pads were placed between the head and head coil to  
307 minimize movement. Auditory stimuli were presented through MR-compatible  
308 headphones with noise-cancellation (OptoACTIVE). An audio recording of the  
309 participant's response was made with an fMRI compatible microphone (OptoACTIVE  
310 FOMRI-III) attached to the headset, which was used to verify the correctness of the  
311 produced sentence. At debriefing participants reported that during production they  
312 were unable to hear themselves speak, confirming that the noise masking of  
313 auditory feedback was successful.

314 While it is generally assumed that overt speech in the scanner will cause large  
315 motion and signal artifacts (see Gracco, Tremblay & Pike, 2005 for an overview) we  
316 did not find this to be the case in our specific set-up. Instead of using a special  
317 scanning procedure (e.g. Eden et al., 1999; Huang et al., 2001; Menenti et al., 2011)  
318 or limit volume acquisition to the time interval after speech production, we applied a  
319 common acquisition procedure. Nevertheless, motions were well within the  
320 boundaries of acceptability (no movement in any direction exceeding the voxel

321 dimensions of 3.5 mm), and no signal artifacts were found. Of the total group only  
322 data of one participant had to be discarded due to excessive motion artifacts.

323

## 324 2.4 Data analysis

### 325 2.4.1 fMRI data pre-processing

326 Data processing and analyses were performed using Matlab and SPM8  
327 software (Wellcome Department of Cognitive Neurology, London, UK). The first nine  
328 scans of all EPI series were excluded from the analysis to minimize T1 relaxation  
329 artifacts and to allow for an optimization of the noise-cancellation. Data processing  
330 started with slice time correction and realignment of the EPI datasets. A mean image  
331 for all EPI volumes was created, to which individual volumes were spatially realigned  
332 by rigid body transformation. The high-resolution structural image was co-registered  
333 with the mean image of the EPI series. The structural image was normalized to the  
334 Montreal Neurological Institute (MNI) template. The normalization parameters were  
335 then applied to the EPI images to ensure an anatomically informed normalization.  
336 Motion parameters were estimated for each session separately. A commonly applied  
337 filter of 8 mm FWHM (full-width at half maximum) was used. The time series data at  
338 each voxel were processed using a high-pass filter with a cut-off of 128 s to remove  
339 low-frequency drifts.

340

### 341 2.4.2 General GLM analyses

342 The subject-level statistical analyses were performed using the general linear  
343 model (GLM). All events of interest were time-locked to the correctness judgments.  
344 We time-locked to judgments rather than to speech errors themselves for several



345 reasons. First, it was not uncommon for participants to produce multiple errors per  
346 sentence. In this case it is unclear which error the activation needs to be time-locked  
347 to. Second, there presumably is high variation in timing between the production of  
348 an error and the detection of that error (e.g. Hartsuiker & Kolk, 2001). So time-  
349 locking to the production is still not time-locking to the error detection. And as the  
350 BOLD response is quite slow and broad (peaks at 5-6 seconds after stimulus onset  
351 and declines slowly until about 10 seconds after stimulus onset), time-locking to the  
352 correctness judgment will still capture relevant activations. For this analysis the  
353 events of interest were Correct trials (where the sentence production was correct)  
354 and Incorrect trials (where the repetition contained an error). Trials where the  
355 participant had given an incorrect judgment formed a separate regressor of no  
356 interest (data loss: 16% in the production condition, 19% in the perception  
357 condition). Vectors containing the event onsets were convolved with the canonical  
358 hemodynamic response function (HRF) to form the main regressors in the design  
359 matrix (the regression model). The vectors were also convolved with the temporal  
360 derivatives and the resulting vectors were entered into the model. In the model, we  
361 also included regressors to account for variance associated with head motion. The  
362 statistical parameter estimates were computed separately for each voxel for all  
363 columns in the design matrix. Separately for the production and perception  
364 condition, one main contrast was calculated for each single subject: erroneous trials  
365 vs. correct trials. These contrasts from the single subject analyses were submitted to  
366 a factorial design with condition (production vs. perception) as factor.

367 Only results significant at the familywise peak-level threshold of  $p < .05$  are  
368 reported. The resulting maps were overlaid onto a structural image of a standard  
369 MNI brain and the coordinates reported correspond to the MNI coordinate system.

370

#### 371 2.4.3 Region of interest analysis

372 To specifically test the involvement of the STS/STG in verbal monitoring, a  
373 region of interest (ROI) analysis was performed for brain regions in the STS/STG that  
374 were previously identified to be involved in verbal monitoring. For ROI analysis  
375 spheres with a radius of 6 mm were created at the peaks of activation clusters with  
376 the use of MarsBar tool for SPM (<http://marsbar.sourceforge.net/>). Resulting  
377 percent signal changes were analyzed in a 2 x 2 (Condition and Accuracy) repeated  
378 measures ANOVA.

379

380

### 3. Results

#### 381 3.1 Behavioral data

382 During scanning of the production trials, the repetitions of the tongue  
383 twisters were recorded. The experimenter later checked these sound files for  
384 correctness of production and judgment. Only the items in which the participant had  
385 correctly identified his or her performance were included in the analysis; the  
386 incorrectly judged items were discarded from all analyses. Overall the participants  
387 repeated 56% of the tongue twisters correctly and produced errors in 28% of the  
388 trials. In the remaining 16% of the trials, the productions were judged incorrectly  
389 (68% misses, 32% false alarms). Frequently produced errors were phonological slips  
390 (of the type '*a proper cropper...*'), word order errors (of the type '*How a clam can*

391 *cram...'*), adjective omissions (of the type '*a proper coffee pot*'), and semantic  
 392 substitutions that specifically seemed to be targeted at a circumvention of the  
 393 troublesome syllables (similar to '*How can a clam cram into a tidy cream can*'). In  
 394 the perception condition participants correctly identified 53% of the items as correct  
 395 and 27% as incorrect repetitions of the tongue twister. In 19% of the trials the  
 396 participants made an incorrect judgment (40% misses, 60% false alarms). The striking  
 397 similarity between the two conditions is the result of online scoring of the  
 398 production trials, to which the perception trials were matched in percentage of  
 399 errors. An overview of accuracy scores and response times measured from cue until  
 400 judgment is provided in Table 1.

401 In the production condition a significant learning effect was observed ( $F(2, 882) =$   
 402  $38.16, p < .001$ ). In the first block 42% of the sentences was produced correctly, in the  
 403 second block 60% of the trials was produced correctly, and in the third block 68% of  
 404 the repetitions was produced correctly.

405 Table 1. *Accuracy and response time in milliseconds for the correct and incorrect*  
 406 *trials in the production and perception condition of the tongue twister task.*

	Production		Perception	
	Score	RT	Score	RT
Correct	M=56%	3693	M=53%	3554
	Range 35%-70%	(SD 1061)	Range 39-74%	(SD 913)
Incorrect	M=28%	4730	M=27%	3641
	Range 15%-55%	(SD 1357)	Range 17-41%	(SD 1134)

407

408 3.2 fMRI data

409 The contrasts made with the fMRI data were the following:

410 1. For each condition separately we contrasted erroneous trials > correct trials

411 Results from contrast 1 were used to make the following contrasts:

412 2. Conjunction analysis comparing similarities between activations during production

413 and perception.

414 3. Disjunction analysis comparing the activations that are independent for the

415 production and perception condition.

416 3a. Production > Perception.

417 3b. Perception > Production.

418

419 3.2.1. Conjunction analysis

420 A conjunction analysis was used to investigate the areas underlying error

421 detection that are common to speech production and speech perception. In this

422 analysis, we tested for a rejection of the conjunction null hypothesis (i.e., only those

423 voxels were reported as active which proved to be significant for speech production

424 and speech perception). The conjunction analyses revealed several clusters that

425 were commonly more active in erroneous compared to correct trials (Table 2; Figure

426 2). Clusters of activation were found in the pre-SMA extending into the dACC, the

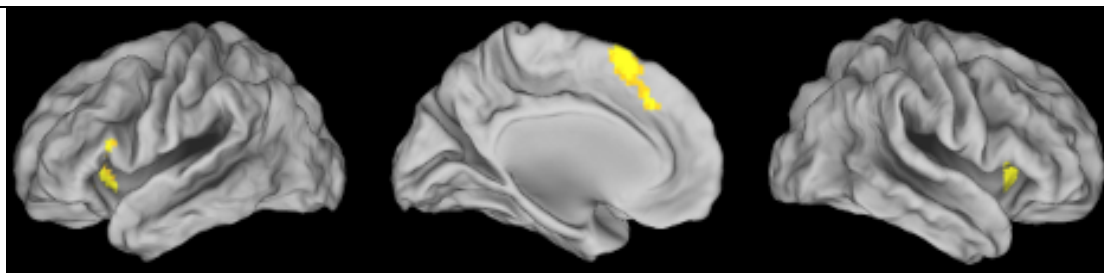
427 left AI and IFG, and the right IFG extending into AI.

428 Table 2.

*Peak Clusters of Activation revealed by Conjunction Error Trials Production and Perception*

Structure	Peak coordinates (MNI)	Z-score	Extent
-----------	------------------------	---------	--------

Pre-Supplementary Motor Area	-6 17 58	5.92	158
Left Insula	-33 20 5	5.95	63
Right Inferior Frontal Pars Triangularis	45 23 1	5.58	62
Left Inferior Frontal Opercularis	-45 20 13	5.05	15



429

430 Figure 2. Activation map averaged across 21 subjects ( $p < .05$ , familywise error  
 431 corrected) of the conjunction analysis error trials production and error trials  
 432 perception.

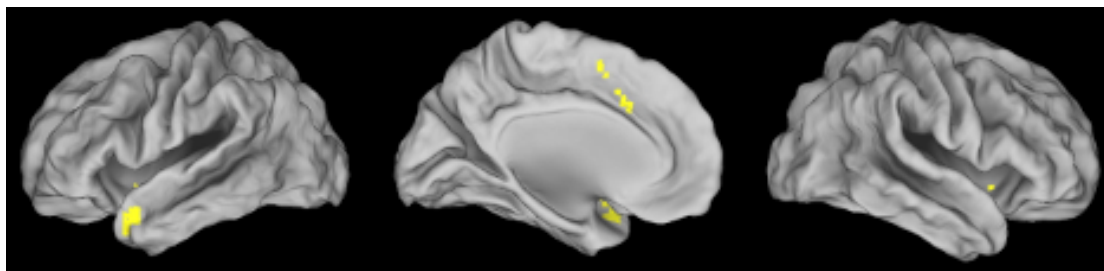
433

### 434 3.2.2 Disjunction analysis

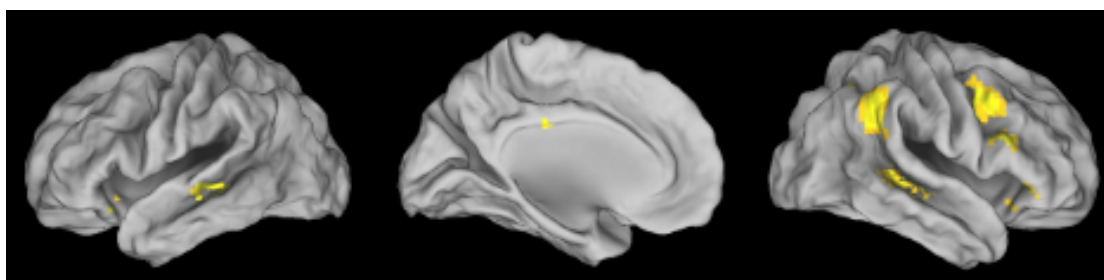
435 To investigate process-specific activations, namely production- and  
 436 perception-specific error detection activation patterns, both an interaction analysis  
 437 and a disjunction analysis can be applied. Both approaches were used to analyze the  
 438 data. As the two analysis methods roughly yielded the same results, we chose to  
 439 report only the results from the disjunction analysis, as they are more  
 440 straightforward to interpret.

441 A disjunction analysis was used to investigate areas active in error detection  
 442 specific for the two modalities, production and perception. Error detection in  
 443 production was masked by error detection in perception to reveal what areas are  
 444 specific for error detection in production. This analysis revealed clusters of activation  
 445 (Table 3, Figure 3) in the left temporal pole, pre-SMA and dACC and BA 48.

446 Error detection in perception was masked by error detection in production to  
 447 reveal areas specific to error detection in perception. This analysis revealed an array  
 448 of clusters, including bilateral posterior superior temporal sulcus / middle temporal  
 449 gyrus (pSTS/MTG), left AI and IFG, right supra marginal gyrus, middle frontal gyrus  
 450 and precentral gyrus, extending into IFG (Table 3, Figure 4).



451  
 452 Figure 3. *Activation map averaged across 21 subjects ( $p < .05$ , familywise error*  
 453 *corrected) of the disjunction analysis error trials production masked by error trials*  
 454 *perception, revealing activation specific for error detection in production.*



455  
 456 Figure 4. *Activation map averaged across 21 subjects ( $p < .05$ , familywise error*  
 457 *corrected) of the disjunction analysis error trials perception masked by error trials*  
 458 *production, revealing activation specific for error detection in perception.*

459

460 Table 3.

*Peak Clusters of Activation revealed by Disjunction Analysis Error Trials Production and Perception*

Structure	Peak coordinates (MNI)	Z-score	Extent
-----------	------------------------	---------	--------

---

*Production Errors – Perception Errors*

Left Temporal Pole	-42 11 -17	5.67	68
ACC	-6 20 34	5.36	18
pre-SMA	-6 8 49	5.36	13
White matter	33 11 -8	5.09	8

*Perception Errors – Production Errors*

Right Middle Frontal Gyrus	45 11 43	6.23	277
Right Middle Temporal Gyrus	54 -37 1	5.75	125
Right Supra Marginal Gyrus	60 -46 31	6.03	173
Left Middle Temporal Gyrus	-57 -28 -5	6.20	56
Right Frontal Inferior Orb	45 35 -5	5.92	18
Right Thalamus	9 -16 10	5.77	10
Right Orbital Inferior Frontal Gyrus	33 23 -14	5.61	10
Corpus Callosum	-3 -25 28	5.06	10
Left Insula	-30 20 -11	5.15	7
Right Middle Frontal Gyrus	30 11 58	4.88	7

---

461

462 3.2.3 ROI analysis of the Superior Temporal Gyrus

463 The perceptual monitoring theories hold that speech monitoring takes place through

464 the speech perception system. Many studies have pointed to the STG as a main locus

465 for speech perception (see Price, 2012) and a possible candidate for perception-

466 based error detection as it has been observed to respond to feedback alterations

467 (McGuire et al, 1996; Hirano et al, 1997; Indefrey &amp; Levelt, 2004; Christoffels et al,

468 2007, 2011; Tourville et al, 2008; Zheng et al, 2010; Takaso et al, 2010). To further  
469 examine the role of STS and STG, the hypothesized locus of the perceptual route for  
470 error detection, additional ROI analyses were conducted. From McGuire,  
471 Silbersweig, and Frith (1996) and Hirano et al. (1997) the clusters that increased for  
472 distorted feedback were selected for ROI analysis, as they are the basis of the  
473 hypothesis for perceptual monitoring through the STS/STG. Nine clusters were  
474 selected, four in the right hemisphere (all STG) and five in the left hemisphere (one  
475 in the STS, four in the STG). In the right hemisphere all selected areas showed a main  
476 effect of modality (all  $p$ 's  $<.005$ ), with higher activation in perception compared to  
477 production. A main effect of accuracy was only significant for one ROI (coordinates:  
478 62 -30 12,  $p < .05$ ), which showed an activation decrease in erroneous trials  
479 compared to correct trials. Significant interactions were observed for three out of  
480 four ROIs (all  $p$ 's  $<.05$ ) (not for 46 -20 4). These interactions were driven by  
481 significant lower activation in erroneous trials compared to correct trials in  
482 production, but not in perception. Results in the left hemisphere gave a more  
483 heterogeneous pattern; all areas showed a significant main effect of modality (all  $p$ 's  
484  $<.05$ ), with four out of five areas showing higher activation for perception compared  
485 to production (coordinates -58 12 4 showed the reverse pattern). With respect to  
486 accuracy an inconsistent and insignificant pattern of activations was observed, with  
487 only a significant main effect in one ROI (coordinates -52 -36 16,  $p < .05$ ), which  
488 showed decreases in erroneous trials compared to correct trials for both production  
489 and perception. A significant interaction was observed in two areas ( $p$ 's  $<.005$ ). This  
490 interaction was driven by significant activation differences between erroneous trials



491 and correct trials (increase in area -58 12 4, and decrease in area -60 -18 4) in  
 492 production, but not in perception.

493         Activation differences between erroneous and correct trials during  
 494 production and perception are presented in table 4. Essentially, these ROI analyses  
 495 show that the bilateral STG are stronger activated during speech perception  
 496 compared to production. With respect to error processing, however, the pattern of  
 497 activations observed was inconsistent with respect to the hypothesis that the STG  
 498 plays a primary role in internal verbal monitoring.

499 Table 4.

*Percentage signal change in bilateral STG in erroneous trials compared to correct trials. Significant signal change is indicated by an asterisk (\*  $p < .05$ , \*\*  $p < .005$ )*

Structure	coordinates (MNI)	Perception	Production
Left STS	-50 -10 0	0.014	0.020
Left STG	-52 -36 16	-0.013	-0.050*
Left STG	-56 -8 0	0.014	-0.013
Left STG	-58 12 4	-0.005	0.109**
Left STG	-60 -18 4	0,031	-0.097**
Right STG	46 -20 4	0.014	-0.015
Right STG	54 -26 8	0.032	-0.070**
Right STG	52 -26 4	0.046*	-0.045*
Right STG	62 -30 12	0.030	-0.108**

500

501

#### 4. Discussion

502           The goal of the current study was to investigate the neuronal structures  
503 underlying internal speech monitoring during production and speech monitoring  
504 during perception, and to use these functional neuroimaging data to distinguish  
505 between current theories of verbal monitoring. Perception-based verbal self-  
506 monitoring theories assume that error detection during speech production and  
507 speech perception both use similar, perceptual routes for error detection.  
508 Production-based theories of self-monitoring do not assume a role for the speech  
509 perception system in internal speech monitoring during production. We observed  
510 that error detection in noise-masked speech production and in speech perception  
511 both recruit the pre-supplementary motor area (pre-SMA), dorsal anterior cingulate  
512 cortex (dACC), bilateral anterior insula (AI), and inferior frontal gyrus (IFG). These  
513 observations suggest that error detection indeed recruits similar neural substrates  
514 and therefore might apply similar mechanisms for monitoring speech during  
515 production and perception. Crucially, no consistent pattern of activation related to  
516 error detection was observed in the bilateral superior temporal sulcus (Hirano et al.  
517 1997; Indefrey & Levelt, 2004; McGuire et al. 1996), suggesting that verbal  
518 monitoring occurs largely independent of speech perception systems.

519           The findings of the activation of a perception-independent monitoring  
520 network, and the inconsistent finding with respect to STG activation, taken together  
521 do not offer support for the perceptual monitoring theories, which assume error  
522 detection in internal speech to take place through speech perception processes  
523 (Hartsuiker & Kolk, 2001; Levelt, 1983, 1989; Indefrey and Levelt, 2004; Indefrey  
524 2011; Hickok 2012), but rather supports a conflict monitoring model of error  
525 detection in speech, as proposed by Nozari et al. (2011). This conflict monitoring

526 theory builds on domain-general theories of error detection and conflict resolution  
527 (e.g., Botvinick et al, 2001; Yeung et al, 2004) and proposes that speech monitoring  
528 takes place by measuring conflict in a processing layer, which is sent to a domain-  
529 general executive center, such as the ACC, which increases control in order to  
530 resolve the conflict. Note, however, that our findings are also compatible with the  
531 forward modeling theory for monitoring as proposed by Pickering and Garrod (2013,  
532 2014), which also assumes a perception independent monitor.

#### 533 4.1 The role of the STG in verbal monitoring

534 If the STS/STG were the main locus for error detection in speech, activation  
535 increases would be expected for erroneous trials compared to correct trials, in both  
536 speech production and perception. Instead we found increased activations in the  
537 STG in production, compared to comprehension, and an inconsistent pattern of  
538 activation with respect to erroneous compared to correct trials. In both hemispheres  
539 one cluster showed a main effect of accuracy, with decreased activation in  
540 erroneous compared to correct trials. Additionally in the right hemisphere we found  
541 an interaction of accuracy and modality, with lower activations in erroneous trials  
542 compared to correct trials in production, but not in perception. This finding is  
543 surprising, and not easy to interpret. At least the finding suggests a role for the right  
544 STG during speech production, related to verbal monitoring. However, we must be  
545 cautious in interpreting this finding, as it is a finding from a post-hoc analysis, and  
546 the direction of the effect does not conform to any of our predictions.

#### 547 4.3 Domain general conflict monitoring

548 The conflict monitoring literature supports an explanation of the current  
549 findings within a framework of a domain general monitoring mechanism. The

550 structures found to be active in monitoring during speech perception and internal  
551 speech monitoring during speech production (the pre-SMA, ACC, IFG, and AI) are all  
552 regions that have been related to conflict processing in numerous tasks that require  
553 conflict resolution. The same network has been found to be active in both error  
554 making and error observation in the action domain: error detection increased  
555 activity in the ACC, SMA, pre-SMA, and AI (Newman-Norlund et al. 2009; Desmet et  
556 al. 2013, Monfardini et al. 2013). In the literature this network is also described as  
557 the cingulo-opercular network, which has been related to task maintenance (e.g.  
558 Dosenbach et al., 2008). The pre-SMA and ACC play a critical role in performance  
559 monitoring and adjustment of cognitive control (e.g., Botvinick et al., 2001;  
560 Ullsperger & Von Cramon, 2006; Bonini et al., 2014). The ACC has consistently been  
561 found to be activated after response conflict detection, errors, and unfavorable  
562 outcomes (see Ridderinkhof, 2004 for an overview). Also the dorsal ACC has been  
563 localized as the primary generator of the ERN component (e.g., Van Veen & Carter,  
564 2002; Herrmann et al., 2004). The IFG / AI has also frequently been observed in  
565 cognitive control tasks and tasks engaging attentional processes (e.g., Craig, 2010),  
566 and is hypothesized to be responsible for signaling awareness and in regulating  
567 response selection (see Tops & Boksem, 2011 for an overview). Increased right IFG  
568 activation is often observed in tasks involving stopping one's actions, including  
569 stopping speech (Xue et al., 2008). Increased right IFG activation was also observed  
570 in preparation of word pairs that were primed to lead to embarrassing vs. neutral  
571 speech errors, showing its involvement in increased control during language  
572 processing (Severens et al., 2011). Together these areas form a domain-general  
573 network for conflict resolution.

#### 574 4.4 Process specific activations

575           Apart from the domain-general activations, as observed in the conjunction  
576 analysis, error detection in speech perception and production showed process-  
577 specific activations. Self-monitoring of internal speech during noise-masked speech  
578 production recruited the left temporal pole and pre-SMA and ACC. The pre-SMA is  
579 known to have a somatotopic organization (Chainay et al, 2004; Alario et al, 2006),  
580 resulting in process-specific activations. Left temporal pole activations are observed  
581 in tasks requiring the composition of sentence meaning, and more specifically in the  
582 processing of syntactic structure (Vandenberghe, 2002; Grodzinsky & Friederici  
583 2006; Humphries, 2006

584           Error detection in speech perception revealed process-specific activations in  
585 a few clusters in the left hemisphere, and a more extensive pattern of activation  
586 clusters in the right hemisphere. Left hemisphere activations include anterior insula  
587 and posterior middle temporal gyrus. The left insula has interestingly been  
588 demonstrated to play a crucial role in phonological retrieval and articulation (Shafto  
589 et al., 2010). Activations in the pSTS/MTG are observed bilaterally in response to  
590 (noisy) auditory stimuli (Bates, 2003; Boatman 2004; Fu et al. 2006), and in  
591 integration of auditory and visual information (Beauchamp et al., 2004). Left MTG  
592 has also been linked to semantic processing (e.g. Demonet et al., 1992, 1994;  
593 Vandenberghe et al, 1996; Stromswold et al, 1996; Binder et al, 2009; Diaz and  
594 McCarthy, 2009; but see Price, 2012). In the right hemisphere large clusters are  
595 observed in the posterior middle frontal gyrus, precentral gyrus, in the  
596 supramarginal gyrus, in the IFG/AI, and in the pSTS/MTG. The supramarginal gyrus is  
597 involved in phonological perception and decision making (Hartwigsen et al. 2010;

598 Buchsbaum et al. 2008; McDermott et al. 2003; Price et al. 1997) although it typically  
599 does not show up in speech comprehension tasks (Hickok and Poeppel, 2007;  
600 Rauschecker and Scott, 2009).

#### 601 4.5 Similar findings in monitoring language processing

602 A highly similar pattern of activation for error perception processing was  
603 found in a study in which participant detected semantic errors during reading  
604 (Raposo & Marques, 2013). Compared to correct sentences, sentences with  
605 semantic anomalies increased attention in the right precentral gyrus, right marginal  
606 gyrus, and the ACC. The same areas are observed to be increased in activation in the  
607 perception condition of the current experiment. The fact that monitoring in this  
608 different modality, namely reading, shows similar results further supports a domain  
609 general monitoring mechanism.

610 The current findings are also in line with preceding research into language  
611 control and altered feedback monitoring, which consistently reported activations in  
612 the ACC, SMA and frontal areas (e.g. Fu et al, 2006; Christoffels et al, 2007; Tourville  
613 et al, 2008; Piai et al. 2013). One interesting difference between the before-  
614 mentioned studies of Fu et al. (2006) and Christoffels et al. (2007) into feedback  
615 monitoring and our findings is that we did not find increased activations in the  
616 cerebellum. These cerebellar activations during feedback processing have also been  
617 related to error detection in perception-based models, as it is hypothesized to drive  
618 corrective motor commands to the motor cortex after receiving input from  
619 somatosensory and auditory areas (Ito, 2008; Tourville & Guenther, 2011; Hickok,  
620 2012). While the studies above specifically looked at the effect of manipulating  
621 external feedback, we have excluded external feedback by noise masking. This hints

622 that the role of the cerebellum might be more closely related to external feedback  
623 instead of monitoring proper.

624 In line with our findings are recent studies in which fMRI was used to study  
625 conflict resolution in language processing. Wittfoth et al (2009) investigated  
626 emotional conflict processing in speech perception, and Piai et al. (2013)  
627 investigated attentional conflict in language and non-language processing.  
628 Processing of emotional conflicting information (e.g., a semantically positive  
629 sentence with a negative prosody) also showed an increase in BOLD response in the  
630 posterior medial prefrontal cortex extending into ACC, bilateral insula and IFG,  
631 posterior cingulate and inferior parietal lobule. Processing of attentional conflict in a  
632 Stroop Task (color word is printed in an incongruent ink color), a Picture-Word  
633 Interference Task (picture and distractor are semantically related), and a Simon Task  
634 (press a left or right button to a visual stimulus presented on the opposite side) all  
635 elicited ACC activation. So what we observe in speech error detection are activations  
636 consistent with a domain-general error detection mechanism, through performance  
637 monitoring and adjustment of cognitive control.

638 The finding of the current study, that of a conflict monitoring system which  
639 operates during both production and perception, and which has been observed to  
640 perform the same task in non-linguistic processes is important for three reasons.  
641 First of all, these results have provided a preliminary answer to the question  
642 whether verbal monitoring in production is perception-based. Clearly verbal  
643 monitoring can occur largely independent of perception systems, and is therefore  
644 production-based.

645           Second, as the network described here for verbal monitoring already has  
646 been studied much more extensively in relationship to conflict monitoring, we can  
647 now further investigate whether conflict monitoring mechanisms can apply similarly  
648 to verbal monitoring. This could hugely increase our understanding of verbal  
649 monitoring, and how this could lead to monitoring deficits, which presumably  
650 underlie speech pathologies such as stuttering, and auditory verbal hallucinations  
651 such as observed in schizophrenia.

652           Third, most current theories of production-based monitoring are limited to  
653 speech production. The current findings provide insight into how verbal monitoring  
654 might occur during speech perception; namely highly similar as during production.

#### 655 4.6 Limitations

656           The current study has some limitations, of which a few pertain to the use of  
657 noise masking. The first issue regarding the presentation of noise masking during  
658 production is that it might have induced activations (e.g. Scott & McGettigan, 2013;  
659 Scott et al., 2004) and increased the cognitive load for the participants. However, as  
660 this would have equally affected the erroneous and the correct trials, which we  
661 contrasted to see the neural basis of error detection in verbal monitoring, the  
662 activations we report in the current paper are not noise-induced activations. If  
663 indeed the presence of noise did increase the cognitive load, it might have resulted  
664 in the production of more errors, which would have been beneficial for the current  
665 study.

666           A second comment related to noise masking during production is that  
667 proprioception and bone conduction of the produced speech cannot be excluded as  
668 a monitoring channel. Lackner and Tuller (1979) hypothesized that word selection



669 errors could be detected on the basis of tactile feedback. However, a more recent  
670 study by Postma and Noordanus (1996) contradicts this claim. In their study errors  
671 were reported during four production conditions: silent, mouthed, noise-masked  
672 and normal feedback. The number of reported errors were the same for the first  
673 three conditions, but increased in the fourth. Only the feedback from the external  
674 channel after production provides additional information for error detection on top  
675 of internal channel monitoring. If proprioception and bone conduction were  
676 channels by which monitoring can take place on top of internal speech, one would  
677 expect to see an increase in number of errors reported in the mouthed  
678 (proprioception) and noise-masked condition (proprioception and bone conduction)  
679 compared to the silent condition. But since proprioception and bone conduction did  
680 not contribute to the detection of more errors compared to the silent speech task,  
681 we cannot assume these channels to be of significant value for monitoring.

682         Despite these limitations resulting from noise masking during speech  
683 production, we opted for noise-masked feedback. By noise masking the overt speech  
684 with headphones, the participant could not hear his or her auditory feedback and  
685 would thus have to monitor their internal speech, and the experimenter could use  
686 the produced overt speech to verify the correctness of the repetition. Another  
687 benefit is that by having the participant produce overt speech, unlike covert speech,  
688 it is certain beyond doubt that the speech plan is fully formed (Barch et al., 1999;  
689 Huang et al., 2001; Gracco et al., 2005).

690         The use of a button press response for error detection can also be seen as a  
691 limitation, as it makes the task somewhat less naturalistic, and focuses the attention  
692 of the participants on error detection. The button press was included in the

693 paradigm to measure whether the participant was aware of the error or not. People  
694 do not correct all their speech errors (e.g. Nootboom, 1980), but it is unclear  
695 whether uncorrected errors are ones the producer was unaware of, or ones where  
696 the producer was aware of the error but did not bother to correct it (Berg, 1986). So  
697 the only way to be sure that a participant was aware of the error was to directly ask  
698 the participants. Also, if large numbers of both conscious and unconscious errors had  
699 been made, it would have been interesting to investigate whether a difference exists  
700 in brain activations between conscious and unconscious error production. However,  
701 too few unconscious errors were produced to make this comparison.

702         A further limitation of this study is that there was no counterbalancing of the  
703 order of the production and the perception condition; each participant first  
704 performed the production condition and then the perception condition. Although  
705 this lack of counterbalancing may have disadvantages, we felt these were  
706 outweighed greatly by the advantage of being able to match the error percentages  
707 in perception to that in production. This is of course only possible with a fixed order  
708 of the conditions, and allows for a direct comparison between production and  
709 perception. An unbalanced distribution of error percentages in the production and  
710 perception condition would severely impair the validity of a comparison between  
711 error detection in the production and perception condition.

712

713         In summary, our results suggest that error detection in speech processing  
714 takes place through a domain-general conflict monitoring system, which comprises  
715 the dorsal anterior cingulate cortex, supplementary motor area, bilateral anterior  
716 insula, and inferior frontal gyrus. This network, which has been consistently

717 observed in non-linguistic conflict, is recruited for both speech perception and  
718 speech production. The lack of evidence for the involvement of the superior  
719 temporal gyrus does not offer support for perceptual theories of error monitoring.  
720 The involvement of the conflict-monitoring network rather argues for a conflict  
721 monitoring account of error detection in speech.  
722

- 723 References
- 724 Abel S, Dressel K, Kümmerer D, Saur D, Mader D, Weiller C, Huber W (2009) Correct  
725 and erroneous picture naming responses in healthy subjects. *Neuroscience Letters*  
726 463:167-171.
- 727 Abbs JH, Gracco VL (1983). Sensorimotor actions in the control multi-joint speech  
728 gestures. *Trends in Neuroscience* 6:391-395.
- 729 Abbs JH, Gracco VL, Cole KJ (1984). Control of multimovement coordination:  
730 sensorimotor mechanisms in speech motor programming. *Journal of Motor Behavior*  
731 16:195-231.
- 732 Alario F-X, Chainay H, Lehericy S, Cohen L (2006) The role of the supplementary  
733 motor area (SMA) in word production. *Brain Research* 1076:129-143.
- 734 Barch DM, Carter CS, Braver TS, Sabb FW, Noll DC, Cohen JC (1999) Overt verbal  
735 responding during fMRI scanning: empirical investigations of problems and potential  
736 solutions. *Neuroimage* 10:642-657.
- 737 Bates E, Wilson SM, Saygin AP, Dick F, Sereno MI, Knight RT, Dronkers NF (2003)  
738 Voxel-based lesion-symptom mapping. *Nature Neuroscience* 6:448–450.
- 739 Beauchamp MS, Lee KE, Argall BD, Martin A (2004) Integration of auditory and visual  
740 information about objects in superior temporal sulcus. *Neuron* 41: 809-823
- 741 Berg T (1986) The aftermath of error occurrence: Psycholinguistic evidence from cut-  
742 offs. *Language and Communication* 6:195–213.
- 743 Binder JR, Desai RH, Graves WW, Conant LL (2009) Where is the semantic  
744 system? A critical review and meta-analysis of 120 functional neuroimaging studies.  
745 *Cerebral Cortex* 19:2767–2796.

- 746 Blackmer ER, Mitton JL (1991) Theories of monitoring and the timing of repairs in  
747 spontaneous speech. *Cognition* 39:173-194.
- 748 Boatman D (2004) Cortical bases of speech perception: evidence from functional  
749 lesion studies. *Cognition* 92:47–65.
- 750 Bonini F, Burle B, Liégeois-Chauvel C, Régis J, Chauvel P, Vidal F (2014) Action  
751 monitoring and medial frontal cortex: leading role of supplementary motor area.  
752 *Science* 343:888-891.
- 753 Botvinick MM, Braver TS, Barch DM, Carter CS, Cohen JD (2001) Conflict Monitoring  
754 and Cognitive Control. *Psychological Review* 108:624-652.
- 755 Buchsbaum BR, D’Esposito M (2008) The search for the phonological store: from  
756 loop to convolution. *Journal of Cognitive Neuroscience* 20:762–778
- 757 Chainay H, Krainik A, Tanguy ML, Gerardin E, Le Bihan D, Lehericy S. (2004) Foot, face  
758 and hand representation in the human supplementary motor area. *Neuroreport*  
759 15:765–769.
- 760 Christoffels IK, Formisano E, Schiller NO (2007) Neural correlates of verbal feedback  
761 processing: an fMRI study employing overt speech. *Human Brain Mapping* 28:868–  
762 879.
- 763 Christoffels IK, Van De Ven V, Waldorp LJ, Formisano E, Schiller NO (2011) The  
764 sensory consequences of speaking: parametric neural cancellation during speech in  
765 auditory cortex. *PLoS One*: 6.
- 766 Craig AD (2010) Once an island, now the focus of attention. *Brain Structure and*  
767 *Function* 214:395–396.
- 768 Curio G, Neuloh G, Numminen J, Jousmaki V, Hari R. (2000) Speaking modifies voice-  
769 evoked activity in the human auditory cortex. *Human Brain Mapping*, 9, 183–191.

- 770 Demonet JF, Chollet F, Ramsay S, Cardebat D, Nespoulous JL, Wise R, Rascol A,  
771 Frackowiak R (1992) The anatomy of phonological and semantic processing in  
772 normal subjects. *Brain* 115:1753–1768.
- 773 Demonet JF, Price C, Wise R, Frackowiak RS (1994) Differential activation of right and  
774 left posterior sylvian regions by semantic and phonological tasks: a positron-  
775 emission tomography study in normal human subjects. *Neuroscience Letters*  
776 182:25–28.
- 777 Desmet C, Deschrijver E, Brass M (2013) How social is error observation? The neural  
778 mechanisms underlying the observation of human and machine errors. *Social*  
779 *Cognitive and Affective Neuroscience* 9:427-435.
- 780 Diaz MT, McCarthy G (2009) A comparison of brain activity evoked by single content  
781 and function words: an fMRI investigation of implicit word processing. *Brain*  
782 *Research* 1282:38–49.
- 783 Dosenbach NU, Fair DA, Cohen AL, Schlaggar BL, Petersen SE (2008) A dual-networks  
784 architecture of top-down control. *Trends in cognitive sciences* 12:99-105.
- 785 Eden GE, Joseph JE, Brown HE, Brown CP, Zeffiro TA (1999) Utilizing hemodynamic  
786 delay and dispersion to detect fMRI signal change without auditory interference: the  
787 behavior interleaved gradients technique. *Magnetic Resonance in Medicine* 41:13-  
788 20.
- 789 Eliades SJ, Wang X (2003) Sensory-motor interaction in the primate auditory cortex  
790 during self-initiated vocalizations. *Journal of Neurophysiology* 89:2194–2207.
- 791 Eliades SJ, Wang X (2005) Dynamics of auditory-vocal interaction in monkey auditory  
792 cortex. *Cerebral Cortex* 15:1510–1523

793 Fu CH, Vythelingum GN, Brammer MJ, Williams SC, Amaro Jr E, Andrew CM, Yaguez  
794 L, van Haren NE, Matsumoto K, McGuire PK (2006) An fMRI study of verbal self-  
795 monitoring: neural correlates of auditory verbal feedback. *Cerebral Cortex* 16:969–  
796 977.

797 Gracco VL, Tremblay P, Pike B (2005) Imaging speech production using fMRI.  
798 *Neuroimage* 26:294–301.

799 Grodzinsky Y, Friederici AD (2006) Neuroimaging of syntax and syntactic processing.  
800 *Current Opinion in Neurobiology* 16:240–246.

801 Hartwigsen G, Baumgaertner A, Price CJ, Koehnke M, Ulmer S, Siebner HR (2010)  
802 Phonological decisions require both the left and right supramarginal gyri.  
803 *Proceedings of the National Academy of Sciences U.S.A.* 107:16494–16499.

804 Hartsuiker RJ, Kolk HHJ (2001) Error monitoring in speech production: A  
805 computational test of the perceptual loop theory. *Cognitive Psychology* 42:113-157.

806 Heinks-Maldonado TH, Mathalon DH, Gray M, Ford JM (2005) Fine-tuning of auditory  
807 cortex during speech production. *Psychophysiology* 42:180–190.

808 Heinks-Maldonado TH, Nagarajan SS, Houde JF (2006) Magnetoencephalographic  
809 evidence for a precise forward model in speech production. *Neuroreport* 17:1375–  
810 1379.

811 Herrmann MJ, Rommler J, Ehlis AC, Heidrich A, Fallgatter AJ (2004) Source  
812 localization (LORETA) of the error-related-negativity (ERN/Ne) and positivity (Pe).  
813 *Cognitive Brain Research* 20:294–299

814 Hickok G, Poeppel D (2007) The cortical organization of speech processing. *Nature*  
815 *Reviews Neuroscience* 8:393–402

- 816 Hickok G (2012) Computational neuroanatomy of speech production. *Nature*  
817 *Reviews Neuroscience* 13:135-145
- 818 Hirano S, Kojima H, Naito Y, Honjo I, Kamoto Y, Okazawa H, Ishizu K, Yonekura Y,  
819 Nagahama Y, Fukuyama H, Konishi J (1997) Cortical processing mechanism for  
820 vocalization with auditory verbal feedback. *Neuroreport* 8:2379–2382.
- 821 Huang J, Carr TH, Cao Y (2001) Comparing cortical activations for silent and overt  
822 speech using event-related fMRI. *Human Brain Mapping* 15:39–53.
- 823 Humphries C, Binder JR, Medler DA, Liebenthal E (2006) Syntactic and semantic  
824 modulation of neural activity during auditory sentence comprehension. *Journal of*  
825 *Cognitive Neuroscience* 18:665–679.
- 826 Indefrey P (2011) The spatial and temporal signatures of word production  
827 components: a critical update. *Frontiers in Psychology* 2:255.
- 828 Indefrey P, Levelt WJM (2004) The spatial and temporal signatures of word  
829 production components. *Cognition* 92:101–144.
- 830 Ito M (2008) Control of mental activities by internal models in the cerebellum.  
831 *Nature Reviews Neuroscience* 9:304–313.
- 832 Lackner JR, Tuller BH (1979) Role of efference monitoring in the detection of self-  
833 produced speech errors. In: *Sentence processing* (Cooper WE, Walker ECT, ed), pp.  
834 281–294. Hillsdale, N.J.: Erlbaum.
- 835 Laver J (1980) Monitoring systems in the neurolinguistic control of speech  
836 production. In: *Errors in linguistic performance: Slips of the tongue, ear, pen, and*  
837 *hand* (Fromkin VA, ed), pp. 287-305. New York: Academic Press.
- 838 Levelt WJM (1983) Monitoring and self-repair in speech. *Cognition* 14:41–104.



839 Levelt WJM (1989) *Speaking: From intention to articulation*. Cambridge, MA: MIT  
840 Press.

841 Matlab and SPM8 software Wellcome Department of Cognitive Neurology, London,  
842 UK.

843 McDermott KB, Petersen SE, Watson JM, Ojemann JG (2003) A procedure for  
844 identifying regions preferentially activated by attention to semantic and  
845 phonological relations using functional magnetic resonance imaging.  
846 *Neuropsychologia* 41:293–303

847 McGuire PK, Silbersweig DA, Frith CD (1996) Functional neuroanatomy of verbal self-  
848 monitoring. *Brain* 119:907–917.

849 Menenti L, Gierhan SME, Segaert K, Hagoort P (2011) Shared language: Overlap and  
850 segregation of the neuronal infrastructure for speaking and listening revealed by  
851 fMRI. *Psychological Science* 22:1173–82.

852 Monfardini E, Gazzola V, Boussaoud D, Brovelli A, Keysers C, Wicker B (2013)  
853 Vicarious neural processing of outcomes during observational learning. *PLOS One* 8.

854 Newman-Norlund RD, Ganesh S, van Schie HT, De Bruijn ERA, Bekkering H (2009)  
855 Self-identification and empathy modulate error-related brain activity during the  
856 observation of penalty shots between friend and foe. *Social Cognitive and Affective*  
857 *Neuroscience* 4:10-22.

858 Nootboom SG (1980) *Speaking and unspeaking: detection and correction of*  
859 *phonological and lexical errors in spontaneous speech* In: *Errors in linguistic*  
860 *performance: Slips of the tongue, ear, pen, and hand* (Fromkin VA, ed), pp. 87-95.  
861 New York: Academic Press.

- 862 Nozari N, Dell GS, Schwartz MF (2011) Is comprehension necessary for error  
863 detection? A conflict-based account of monitoring in speech production. *Cognitive*  
864 *Psychology* 63:1-33.
- 865 Numminen J, Salmelin R, Hari R (1999) Subject's own speech reduces reactivity of the  
866 human auditory cortex. *Neuroscience Letters* 265:119–122.
- 867 Oppenheim GM, Dell GS (2008) Inner speech slips exhibit lexical bias, but not the  
868 phonemic similarity effect. *Cognition* 106:528-537.
- 869 Piai V, Roelofs A, Acheson DJ, Takashima A (2013) Attention for speaking: Neural  
870 substrates of general and specific mechanisms for monitoring and control. *Frontiers*  
871 *in Human Neuroscience* 7:832.
- 872 Pickering MJ, Garrod S (2013) An integrated theory of language production and  
873 comprehension. *Behavioral and Brain Sciences* 36:329-392.
- 874 Pickering MJ, Garrod S (2014) Self- Other and joint monitoring using forward models.  
875 *Frontiers in Human Neuroscience* 8:132-143.
- 876 Postma A, Kolk HHJ (1992) The effects of noise masking and required accuracy on  
877 speech errors, disfluencies, and self-repairs. *Journal of Speech and Hearing Research*  
878 35:537–544.
- 879 Postma A, Noordanus C (1996) Production and detection of speech errors in silent,  
880 mouthed, noise-masked, and normal auditory feedback speech. *Language and*  
881 *Speech* 39: 375–392.
- 882 Postma A (2000) Detection of errors during speech production: a review of speech  
883 monitoring models. *Cognition* 77:97–132.
- 884 Price CJ, Moore CJ, Humphreys GW, Wise RJS (1997) Segregating semantic from  
885 phonological processes during reading. *Journal of Cognitive Neuroscience* 9:727–733

- 886 Price CJ (2012) The anatomy of language: a review of 100 fMRI studies published in  
887 2009. *Annals of the New York Academy of Sciences* 1191:62–88
- 888 Raposo A Marques JF (2013) The contribution of fronto-parietal regions to sentence  
889 comprehension: Insights from the Moses illusion. *NeuroImage* 83:431-437.
- 890 Rauschecker JP, Scott SK (2009) Maps and streams in the auditory cortex: nonhuman  
891 primates illuminate human speech processing. *Nature Neuroscience* 12:718–724.
- 892 Ridderinkhof RK, Van den Wildenberg WPM, Segalowitzd SJ, Cartere CS (2004)  
893 Neurocognitive mechanisms of cognitive control: The role of prefrontal cortex in  
894 action selection, response inhibition, performance monitoring, and reward-based  
895 learning. *Brain and Cognition* 56:29–140.
- 896 Scott S K, McGettigan C (2013) The neural processing of masked speech. *Hearing*  
897 *research*, 303:58-66.
- 898 Scott SK, Rosen S, Wickham L, Wise RJS (2004) A positron emission tomography  
899 study of the neural basis of informational and energetic masking efforts in speech  
900 perception. *The Journal of the Acoustical Society of America* 115:813-821.
- 901 Shafto MA, Stamatakis EA, Tam PP, Tyler LK (2010) Word retrieval failures in old age:  
902 the relationship between structure and function. *Journal of Cognitive Neuroscience*  
903 22:1530-1540.
- 904 Severens E, Kühn S, Hartsuiker RJ, Brass M (2011). Functional mechanisms involved  
905 in the internal inhibition of taboo words. *Social cognitive and affective neuroscience*,  
906 nsr030.
- 907 Siegenthaler BM, Hochberg I (1965). Reaction time of the tongue to auditory and  
908 tactile stimula- tion. *Perceptual and Motor Skills* 21:387-393.

- 909 Stromswold K, Caplan D, Alpert N, Rauch S (1996) Localization of syntactic  
910 comprehension by positron emission tomography. *Brain and Language* 52:452–473.
- 911 Takaso H, Eisner F, Wise RJ, Scott SK (2010) The effect of delayed auditory feedback  
912 on activity in the temporal lobe while speaking: a positron emission tomography  
913 study. *Journal of Speech Language and Hearing Research* 53:226–236.
- 914 Tian X, Poeppel D (2010) Mental imagery of speech and movement implicates the  
915 dynamics of internal forward models. *Frontiers in Psychology* 1: 166.
- 916 Tian X, Poeppel D (2013) The effect of imagination on stimulation: the functional  
917 specificity of efference copies in speech processing. *Journal of cognitive*  
918 *neuroscience* 25: 1020-1036.
- 919 Tops M, Boksem MA (2011) A potential role of the inferior frontal gyrus and anterior  
920 insula in cognitive control, brain rhythms, and event-related potentials. *Frontiers in*  
921 *Psychology* 2:330.
- 922 Tourville JA, Reilly KJ, Guenther FH (2008) Neural mechanisms underlying auditory  
923 feedback control of speech. *Neuroimage* 39:1429–1443.
- 924 Tourville JT, Guenther FH (2011) The DIVA model: A neural theory of speech  
925 acquisition and production. *Language and Cognitive Processes* 25:952-981.
- 926 Ullsperger M, von Cramon DY (2006) The role of intact frontostriatal circuits in error  
927 processing. *Journal of Cognitive Neuroscience*, 18:651-664.
- 928 Van Veen V, Carter CS (2006) Error detection, correction, and prevention in the  
929 brain: a brief review of data and theories. *Clinical EEG and neuroscience*, 37:330-335.
- 930 Van Veen V, Carter CS (2002) The timing of action-monitoring processes in the  
931 anterior cingulate cortex. *Journal of cognitive neuroscience*, 14:593-602.

- 932 Vandenberghe R, Price C, Wise R, Josephs O, Frackowiak RS (1996) Functional  
933 anatomy of a common semantic system for words and pictures. *Nature* 383:254–  
934 256.
- 935 Vandenberghe R, Nobre AC, Price CJ (2002) The response of left temporal cortex to  
936 sentences. *Journal of Cognitive Neuroscience* 14:550–560.
- 937 Wittfoth M, Schröder C, Schardt DM, Dengler R, Heinze HJ, Kotz SA (2009). On  
938 emotional conflict: interference resolution of happy and angry prosody reveals  
939 valence-specific effects. *Cerebral Cortex*, 20:383-392.
- 940 Xue G, Aron AR, Poldrack RA (2008). Common neural substrates for inhibition of  
941 spoken and manual responses. *Cerebral Cortex* 18:1923-1932.
- 942 Yeung N, Botvinick MM, Cohen JD (2004) The neural basis of error detection: conflict  
943 monitoring and the error-related negativity. *Psychological Review* 11:931–959.
- 944 Zheng ZZ, Munhall KG, Johnsrude IS (2010) Functional overlap between regions  
945 involved in speech perception and in monitoring one's own voice during speech  
946 production. *Journal of Cognitive Neuroscience* 22:1770-1781.
- 947

948 Appendix A

949 Dutch tongue twister sentences with rough translations to English in italics.

950 Tongue twister used in the production condition.

951 1. Ruud rups raspt rap rode ronde radijsjes.

952 *Cathy caterpillar quicky grates round red radishes.*

953 2. De meid sneed zeven scheve sneden brood.

954 *The maid cut seven skew slices of bread.*

955 3. Als apen apen naapen, apen apen apen na.

956 *When monkeys mimic monkeys, monkeys mimic monkeys.*

957 4. Wiske mixt whisky in de whisky mixer.

958 *Wilma mixes whisky in the whisky mixer.*

959 5. Gijs grijpt de grijsgrauwe gans graag gauw.

960 *Gordon gladly grabs the grey goose swiftly.*

961 6. Baardige artsen helpen aarzelende bedelaars.

962 *Bearded doctors help hesitant beggars.*

963 7. Als een potvis in een pispot pist, zit de pispot vol met potvispis.

964 *If a sperm whale pisses in a pissjar, the pissjar is filled with sperm whale piss.*

965 8. Een pet met een platte klep is een plattekleppet.

966 *A cap with a flat flap is a flat flap cap.*

967 9. Vaders vader vond vier vuile vesten van vier vuile venten.

968 *Father's father found four filthy cardigans of four filthy blokes.*

969 10. Sluwe feministen foeteren op flemende sloeries.

970 *Sly feminists grumble about flanneling floozies.*

971 11. Jeukt jouw jeukende neus zoals mijn jeukende neus jeukt?

- 972 *Does your itchy nose itch like my itchy nose itches?*
- 973 12. De koetsier poetst de postkoets met postkoetspoets.
- 974 *The coachman polishes the coach with coach polish.*
- 975 13. Pappa pakt de platte blauwe bakpan.
- 976 *Daddy grabs the flat blue frying pan.*
- 977 14. Pseudo-psychologen sporten als speren.
- 978 *Pseudo-psychologists sport like crazy.*
- 979 15. Aaibare kraaien leggen kale kraaie-eieren.
- 980 *Cuddly crows lay bald crow eggs.*
- 981 16. Ping pingpongde de pingpongbal naar Pong.
- 982 *Ping ping-ponged the ping pong ball to Pong.*
- 983 17. Krakende krekels trippelen op tegels.
- 984 *Creaking crickets patter on the tiling.*
- 985 18. De kat krabt de krullen van de trap.
- 986 *The cat scratches shavings of the stairs.*
- 987 19. Knappe kappers kappen knap.
- 988 *Handsome hairdressers cut handsomely.*
- 989 20. Achtentachtig achterdochtige doktersdochters.
- 990 *Eighty-eight suspicious doctor's daughters.*
- 991 21. Zeven zotten zullen zes zomerse zondagen zwemmen zonder zwembroek.
- 992 *Seven fools will swim six Sundays without swimming trunks.*
- 993 22. Piet's priesterpij is piepklein.
- 994 *Pete's priests frock is very tiny.*
- 995

- 996 Tongue Twisters used in the perception condition
- 997 1. Trillend trippelde tante Tiny tandloos naar de treiterende tandarts toe.  
 998 *Aunt Tilly tremblingly toddled toothless to the harassing dentist.*
- 999 2. Tijdens de afwas viel de asbak in de afwasbak.  
 1000 *During the wash up the ashtray fell into the kitchen sink.*
- 1001 3. Liesje leerde Lotje lopen langs de lange Lindenlaan.  
 1002 *Lacey learned Laney how to walk along the long Linden lane.*
- 1003 4. Knappe slakken snakken naar slappe sla.  
 1004 *Pretty snakes yearn for limp lettuce.*
- 1005 5. Toen Lotje niet wou lopen, liet Liefje Lotje staan.  
 1006 *As Lany would not walk, Lacy left Lany.*
- 1007 6. Vissende vissers die vissen naar vissen, maar vissende vissers die vangen vaak  
 1008 bot.  
 1009 *Fishing fishermen fish for fish, but fishing fishermen often catch zilch.*
- 1010 7. Dikke drilboren drillen door dikke deuren.  
 1011 *Large drills drill trough thick doors.*
- 1012 8. Kriegelig kocht Krelis kilo's kruimige krieltjes.  
 1013 *Grumpily Gary bought kilo's of floury spuds.*
- 1014 9. De dunne dokter duwde de dikke dame door de draaiende draaideur.  
 1015 *The thin doctor pushed the fat lady through the spinning revolving door.*
- 1016 10. Trollen rollebollen als dollen in de drollen.  
 1017 *Trolls horse around in the turds like crazy.*
- 1018 11. Ezels eten netels niet en netels eten ezels niet.  
 1019 *Donkeys don't eat nettles, and nettles don't eat donkeys.*



- 1020 12. De pasgewassen was was pas gewassen nadat de pasgewassen was gewassen  
1021 was.  
1022 *The freshly washed laundry was only washed after the freshly washed laundry*  
1023 *was washed.*
- 1024 13. De magere marktskraamvrouw kookte veel makreel.  
1025 *The skinny stall woman cooked lots of mackerel.*
- 1026 14. De toetsenist test het toetsenbord.  
1027 *The keyboardist tests the keyboard.*
- 1028 15. De grommende beer bromt beestachtig geestig.  
1029 *The growling bear grumbles mightily funny.*
- 1030 16. Kniezende kneuzen kiezen kale keukens.  
1031 *Moping misfits choose bare kitchens.*
- 1032 17. Babbelende baby's dromen van dommelende bosduifjes.  
1033 *Babbling babies dream of dozy wild pigeons.*
- 1034 18. Pinnige dikke piloten drinken prille pils.  
1035 *Stingy fat pilots drink early bears.*
- 1036 19. Nukkige nuchtere Nellie is niet nuttig.  
1037 *Crancky sober Nelly is not usefull.*
- 1038 20. Gerooide woudreuzen groeien in mooie wouden.  
1039 *Cleared wood giants grow in pretty woods.*
- 1040 21. Slome slavinnen lopen in sombere lompen.  
1041 *Slow slaves walk around in dreary duds.*
- 1042 22. De stille prinses at knisperende spritsen.  
1043 *The princess ate crackling cookies.*