# Very short utterances
## and timing in turn-taking

**Mattias Heldner[1,2], Jens Edlund[1], Anna Hjalmarsson[1], Kornel Laskowski[1]**
**[1]Speech, Music and Hearing, KTH, Stockholm, Sweden**
**[2]New job! Department of Linguistics, Stockholm University, Stockholm, Sweden**

**Is there a difference in the timing of very short utterances, such as backchannels, compared to longer utterances? (i.e. timing relative to neighboring utterances)**
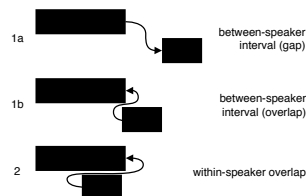
**How does the inclusion/exclusion of between-speaker intervals (i.e. gaps, overlaps, no-gap-no-overlaps) adjacent to very short utterances affect the distribution of such intervals?**

## Speech material

A subset of the Spontal corpus, about 8 hours of spontaneous two-party face-to-face conversations in Swedish recorded with close talk microphones.

## Procedure

1. Speech activity detection (100 ms frame step, 200 ms frame size) produced a segmentation into TALKSPURTS ≥ 200 ms; and PAUSES ≥ 200 ms for each speaker.
2. The TALKSPURTS were subdivided into very short utterances (VSUs) and their complement (NONVSUs) based on their duration: VSUs = TALKSPURTS ≥ 200 ms and ≤ 1000 ms; NONVSUs = TALKSPURTS ≥ 1100 ms.
3. The TALKSPURTS and PAUSES of the individual speakers were combined to identify intervals of SINGLE-SPEAKER SPEECH for each speaker, JOINT SILENCE and JOINT SPEECH.
4. Sequences of such intervals were used to identify BETWEEN-SPEAKER INTERVALS = intervals of JOINT SILENCE or JOINT SPEECH preceded and followed by SINGLE-SPEAKER SPEECH from different speakers; and WITHIN-SPEAKER OVERLAPS = JOINT SPEECH preceded and followed by SINGLE-SPEAKER SPEECH from the same speaker.
5. Durations of BETWEEN-SPEAKER INTERVALS and WITHIN-SPEAKER OVERLAPS were calculated.



6. The VSU/NONVSU distinction was used to identify four types of between-speaker intervals: (i) VSU-VSU; (ii) VSU-NONVSU; (iii) NONVSU-VSU; and (iv) NONVSU-NONVSU.
7. Durations of BETWEEN-SPEAKER INTERVALS were related to detection thresholds for gaps and overlaps (Heldner, 2011): GAP = JOINT SILENCE ≥ 200 ms; OVERLAP = JOINT SPEECH ≥ 200 ms; NO-GAP-NO-OVERLAP = JOINT SPEECH or JOINT SILENCE ≤ 100 ms.

## Timing differences

Table 1. Descriptive statistics for the durations (in ms) of the four types of between-speaker intervals

|  | Mean | SD | N |
|---|---|---|---|
| VSU-VSU | 281 | 599 | 432 |
| VSU-NONVSU | 287 | 630 | 1618 |
| NONVSU-VSU | 203 | 480 | 1621 |
| NONVSU-NONVSU | −36 | 821 | 2835 |
| Total | 125 | 709 | 6506 |

Table 2. Distribution of perceived speaker change categories for the four types of between-speaker intervals

|  | Perceived as | % |
|---|---|---|
| VSU-VSU | OVERLAP | 11.8 |
|  | NO-GAP-NO-OVERLAP | 36.1 |
|  | GAP | 52.1 |
| VSU-NONVSU | OVERLAP | 13.8 |
|  | NO-GAP-NO-OVERLAP | 34.1 |
|  | GAP | 52.1 |
| NONVSU-VSU | OVERLAP | 12.3 |
|  | NO-GAP-NO-OVERLAP | 37.5 |
|  | GAP | 50.2 |
| NONVSU-NONVSU | OVERLAP | 37.5 |
|  | NO-GAP-NO-OVERLAP | 25.7 |
|  | GAP | 36.8 |

- Short utterances, such as backchannels, are timed differently relative to neighboring utterances than longer utterances:
- Between-speaker intervals delineated by VSUs on either or both sides have relatively fewer and shorter OVERLAPS and more NO-GAP-NO-OVERLAPS than NONVSU-NONVSUs
- Conversely, NONVSU-NONVSUs have relatively more and longer OVERLAPS; and fewer GAPS and NO-GAP-NO-OVERLAPS than the other types of between-speaker intervals
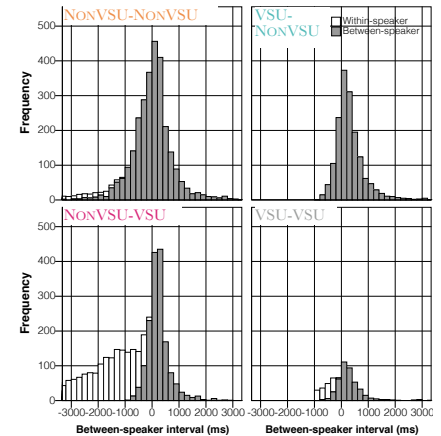


Figure 1. Histograms of between-speaker intervals (grey bars) in the four types of between-speaker intervals. The bin size is 200 ms. The white bars show the interval from the onset of WITHIN-SPEAKER OVERLAP to the offset of the nearest single-speaker speech forward in time.

## Distributions

- While these timing differences may to some extent be an artifact of the definition of VSUs in terms of duration, the choice to include or exclude such short utterances from distributions of between-speaker interval durations affects those distributions considerably:
- Excluding short utterances results in lower central tendency, which has been interpreted as more precise timing, but also in higher variance and relatively fewer NO-GAP-NO-OVERLAPS, which may well be interpreted as less precise timing
- If some proportion of the WITHIN-SPEAKER OVERLAPS had indeed been intended as Overlaps,

## Conclusions

It is important to keep track of whether short utterances, such as backchannels, are included or not when interpreting and comparing between-speaker intervals.

## Acknowledgements