



Københavns Universitet



## The Logic of comparability

Gregersen, Frans; Barner-Rasmussen, Michael

*Published in:*

Corpus Linguistics and Linguistic Theory

*DOI:*

[10.1515/CLLT.2011.002](https://doi.org/10.1515/CLLT.2011.002)

*Publication date:*

2011

*Document Version*

Early version, also known as pre-print

*Citation for published version (APA):*

Gregersen, F., & Barner-Rasmussen, M. (2011). The Logic of comparability: On genres and phonetic variation in a project on language change in real time . *Corpus Linguistics and Linguistic Theory*, 7(1), 7-36.  
<https://doi.org/10.1515/CLLT.2011.002>

# The Logic of comparability: On genres and phonetic variation in a project on language change in real time

FRANS GREGERSEN and MICHAEL BARNER-RASMUSSEN<sup>1</sup>

## *Abstract*

*This paper is based on data from the LANCHART (Language Change in Real Time) corpus. LANCHART is 'data-based' in that transcripts are orthographically normalized and both transcripts and annotations are time stamped and stored in a database. The corpus is structured according to generation, gender, geography, class, and time of recording. The socio-linguistic issue treated is the relationship between genre, as defined in the context of a so-called Discourse Context Analysis, and a particularly frequent and significant phonetic variable, viz. the (æ) variable in Modern Danish. Through repeated searches in the corpus we show that genres do have an effect on phonetic variation and that they frequently develop in real time, at least with regard to the patterning of the selected phonetic variable. The results are discussed in the context of intra-individual variation within the sociolinguistic interview, i.e., the time-honored style problem.*

*Keywords:* sociolinguistics, genre, pragmatic variation and phonetic variation, Danish, language change in real time, data-based transcript model

## **1. Introduction**

Historically speaking, corpus linguistics and sociolinguistics together make up an empirical response to the traditional focus on introspective evidence (Labov 1972a). In this respect, sociolinguistics and corpus linguistics are alike. They base themselves squarely and solely on observational evidence (Labov 1984). But there are also important differences between the two approaches. Sociolinguistics from the outset has focused on spoken language as the relevant setting for the study of language variation and change. Corpus linguists, however, have on the whole used their ready access to large, and rapidly increasing,

corpora of written language (though Biber's work is the obvious exception to this general claim having made a point of contrasting evidence from such corpora with speech data; cf. Biber 1988, 1995, Biber et al. 1998: 106ff, Biber & Conrad 2009). Only recently do corpus linguistics and sociolinguistics approach a situation where the massive sociolinguistic investment in building large digitized corpora of spoken language makes it fruitful, indeed necessary, to use corpus-linguistic methods (i.e. time stamping, tool building of search machines and repeated searches) to address sociolinguistic problems.

## 2. Problems of comparability

Observational data almost invariably are available as recordings. Recordings are compared, but that presupposes that they are comparable. The essential problem of comparability is in our view to control intra-individual variation *within* recordings in order to get at inter-individual differences *between* recordings, both those from different persons and those from the same person but from another point in time. This problem used to be couched in the discourse of style. Intra-individual variation was seen as dependent on various contexts (Labov 1966, 1972) and there was extended discussion on how to be sure that what was compared was really data from specific styles only. Style was the controlling factor (Labov 2001). Since Eckert (2003) and Coupland (2007) style has, however, not been the same, and we refrain here from discussing comparability in style terms (but cf. Section 8 below). Instead we focus on *the conditions for inferring change from a sample or corpus of sociolinguistic interviews recorded at two different points in real time*.

In empirical studies of language change based on spoken corpora, change is either inferred from a sample stratified for age in so-called apparent time (cf. the most recent discussion in José 2010), or studied on the basis of recordings/studies at two (or more) separate points in time, known as change in "real time".<sup>2</sup> Obviously, there has to be some distance in time between the two recordings for us to be able to detect change.<sup>3</sup> We adopt the convention of abbreviating the original (first) study S1, while the repetition of S1, i.e. the equivalent new study, is designated as (an instance of an) S2 (Gregersen 2009).

The specific features of the S1 characteristically limit the S2 possibilities. We want S2 to conform to S1 so that they are comparable. Comparability here means to control for all relevant factors except time so that we are certain that time is the essential difference between the data sets.

One barrier to comparability might be *field methods*, as different methods may lead to characteristically different data. For instance, one tradition, originating with Labov, has relied on the sociolinguistic interview combined with a fixed set of thematically structured questions to generate the necessary amount

of good data (Labov 1966, 1984). Another tradition has preferred to arrange a group discussion, where the interviewer no longer has to occupy a privileged place of interaction but acts more like a facilitator (Labov's New York study is a case in point, cf. Labov 1966 [2006]: 59). Or, the field worker may simply be absent from the group recording, in which case the interviewer is not able to control thematic development at all. In Nordic sociolinguistics this latter, very influential, paradigm may be called the Thelander-Gumperz method (Gregersen 2009: 11–13).

In the first type, theme seems to be the vital factor in securing comparability, hence the focus on the so-called danger of death stories (see below) or the 'first fight' stories. But we ask the reader to bear in mind that it is, in fact, not only the theme which is comparable, but often also the characteristic structure of the speech, i.e. the *genre*. In the second and third type of field methods, interaction and acquaintance are seen as the decisive factors: with only a facilitator or no field worker present, what is there to prevent informants from talking freely to persons they know or interact with daily?<sup>4</sup>

Field methods may matter, but so do *designs*. The first choice in making S2 conform to S1 is the replication. The solution of replicating exactly in S2 what went on in the S1 is conveniently exemplified by Fowler's replication of Labov's classic department store study (Labov 1994: 87ff) and the replication by Pope, Meyerhoff and Ladd (2007) of the equally classic Martha's Vineyard study (Pope, Meyerhoff and Ladd 2007, cf. also Blake and Josey 2003). The Fowler replication of the department store study also demonstrates why it is often impossible to undertake a perfect replication. In the S1 there were three department stores: Saks, Macy's and Klein. In 1986 Klein was no more and so Fowler had to find a substitute. The real world may be a real obstacle to replicating a study in real time.

Two further barriers to comparability relate to the changing social (and potentially sociolinguistic) meanings of aspects of the interview process, such as topic or interviewer. Perhaps the *topics* used as generators of narratives might no longer perform the same function today. To take the most obvious case: if it is really true that New York is now much safer than it was 50 years ago, the so-called danger of death question might not work as well as a generator of narratives now as it did back then. When asked "Have you ever been in a situation where you thought there was a serious danger of your being killed? That you thought to yourself 'This is it?'" (Labov 1966 (2006): 70) many informants might – just as most Danes – answer with a plain 'no',<sup>5</sup> and the interviewer would have to go on without personal narratives of 'danger of death'.<sup>6</sup> If s/he found another 'rich' subject which would trigger personal narratives full of emotion and told with involvement, s/he would have solved this specific problem of comparability only because he or she had *changed* the S2 successfully to conform functionally to S1. This shows that a careful consideration of

function within the speech event may lead to changes in design between the S1 and the S2.

Consider also the notion of the *interviewer*. Is this particular role a historically variable one if we look, for example, at a time distance of 20 years? That depends on (the development of) the position of interviewers in public discourse and public awareness. The notion of an interviewer is tightly bound to the role of the mass media, and if it was a privilege to be interviewed by the mass media back in the 1960s when Labov carried out his first interviews, we may understand why interviewees at that time treated the event of an interview as an occasion for putting on a formal style (Labov 1966 [2006]: 59). If, however, interviews have now become standard practice in and outside of the mass media, the problem might no longer be to overcome a barrier between a casual and a formal style but rather to get the formal style at all. Interview style would no longer by default be formal style.

We shall not in this paper discuss these problems of comparability further, though they might all be illustrated by our data. Rather we focus here on the problem of controlling *intra-individual variation*. In any given S1 recording there will always be some intra-individual variation. But how do we control for this when we carry out the one and only important comparison, that is, between the S1 and the S2? To be more precise: which parts of the S1 may be compared to which parts of the S2?

This involves first finding a practical solution to problems inherent in the characteristics of the linguistic levels: at the level of phonetics, it is simply not practically possible to code all occurrences of the variables, since there are so many of them. This goes in particular for the most frequent ones, which may for this very reason be supposed to be important to perception. This raises two pertinent questions: if we cannot code all occurrences, *which ones should we code*, what selection should we make? This we call the *selection* problem. And related to this: *how many instances* of each variable for each recording? Or rather: how few can we code and still be reasonably sure that we have a sample that is a true subsample of the total set of instances? We refer to this as the *probe* problem. In the sections below we give a brief description of the LANCHART study of language change in the Danish speech community in real time (Sections 3.1 through 4.1) in order to give the necessary preconditions for our answers to these two questions (Section 4.2).

### 3. The LANCHART project

#### 3.1 *Aims and basic terminology*

In 2005 the LANCHART Centre was established at the University of Copenhagen.<sup>7</sup> The purpose of the centre was and is to study variation and change

both at the level of the individual and at the level of the speech community in real time (Sankoff 2005). An essential part of the plan was to digitize and transcribe both previous and planned recordings and to conserve them as transcription files linked to the recording as well as as time-stamped, meta-dated database entries (cf. Kendall 2008). At present (June 2010), the LANCHART corpus contains 452 transcribed sound files with around 5.4 million words out of a total of 1814 recordings. The transcribed corpus thus includes 25% of the total number of recordings. The phonetically and pragmatically coded material (see below) is moreover a fraction of this, although it comprises more than 200,000 manual annotations.<sup>8</sup>

A fair number of early sociolinguistic and dialectological studies analyzed the speech of a stratified section of informants at various places in the Danish speech community during the late 1970s through the 1980s. These studies constitute the LANCHART S1s. Since S2 recordings have been staged from 2005 and onwards, there is in this case a time distance of around twenty years between the S1s and the S2s. The earlier studies were stratified for age and we took care to group all the informants in ‘generations’ so that we could control for differences in effect between informants who were young at the time of the first recording and thus middle-aged in the S2 setting (generation 2) on the one hand, and on the other hand informants who were already middle-aged at the time of the S1 and hence 20 years older (some of them reaching the status of ‘old’ in the S2). A further problem arose here in that the generations to be unified are defined on the basis of chronological age whereas another equally crucial fact is ‘age at time of recording’.<sup>9</sup> The core corpus consists of recordings of informant samples stratified for age, gender, and social class from:

- Vinderup in Western Jutland (S1 from 1978);
- Odder in Eastern Jutland (S1 from 1986–89);
- Næstved, Southern Zealand (two different S1s, both of them from 1986ff);
- and Copenhagen, Central Zealand (S1 from 1986–88).

The core of the LANCHART study is a panel of 42 informants from Copenhagen, 19 from Næstved, 24 from Odder, and 18 from Vinderup, who have all been recorded twice. In this paper, however, the Copenhagen, Næstved, and Odder informants (and not the Vinderup informants, cf. note 9) make up the core sample of 85 informants recorded twice within 20 years.<sup>10</sup>

### 3.2 *Speaker variables*

Four speaker variables are controlled: *gender*, *age* defined as a broader section of time, a time slice (Eckert 1997) i.e. ‘generation’, *social class*, and *geographical location*<sup>11</sup> (Gregersen 2009).

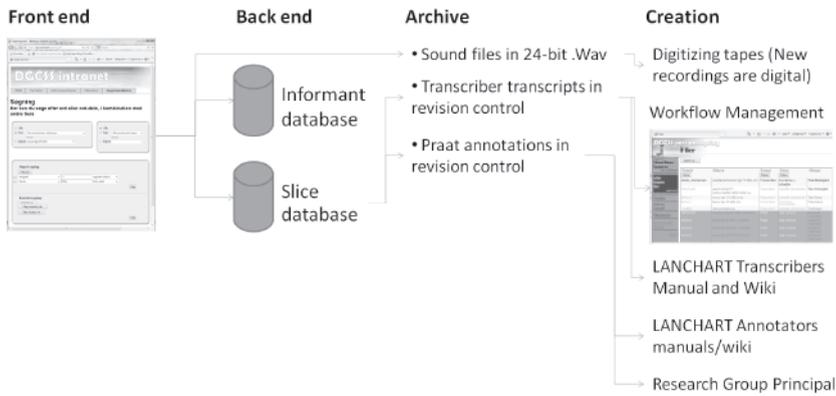
12 *F. Gregersen and M. Barner-Rasmussen*

Figure 1. *The elements of the LANCHART Corpus: a web front end connecting to a database/file system back end that in turn gets its data from the archive of transcriptions, annotations and digitized files.*

### 3.3 Creating the corpus

Kendall (2007, 2008) calls for “sociolinguists to explicitly focus on the *storage, management* and *preservation* of their data”, initially to make sociolinguistic data more like data as we know it from other areas. Ideally, a variable tabulation such as the ones presented in this paper should be directly connected to the source sound-files via time-codes and other meta-data so that other researchers may verify and critique the findings. Kendall further proposes that new types of analysis and new ways of working with data will result from such a focus. In this we wholeheartedly agree.

The LANCHART project is the largest of its kind in Denmark and from the outset it was clear that the endeavor would be impossible without the integration of information technology. To be able to prepare ultimately 1,800+ recordings for the corpus it was decided to use the programs Transcriber<sup>12</sup> for transcriptions and Praat<sup>13</sup> for annotations and to manage the data files created here using revision control software. The process of managing the transcription process was automated using home-grown workflow management software. The vast majority of recordings were made using state-of-the-art (at the time) equipment, e.g. DPA 4088 F lavalier microphones and Revox A77 or Marantz CP430 recorders. Interviews on tape were digitized in 24-bit lossless WAV and S2 recordings made using digital recorders at the same resolution.

In transcribing, special care was taken to ensure the highest possible level of uniformity in order to enable search and comparison across interviews. A detailed manual was prepared and transcribers trained in its use. Word spellings follow the strict orthographical norm institutionalized by the Danish Language

Board in its orthographical dictionary (*Retskrivningsordbogen*) and all transcripts are proof-read by our most experienced transcribers. Additional transcription features, e.g. pauses and false starts, have been kept to an absolute minimum. Transcripts are aligned with sound files at the utterance level. This means that the alignment is not sufficiently precise to carry out automatic analysis of e.g. vowel formants, but precise enough to identify where a word occurs.

The choice of Praat for annotations was based on the Praat data model that allows for any number of ‘annotation tiers’, each tier spanning the length of the sound file and containing any number of time-stamped annotations. Each type of coding can then be manipulated and stored in a separate tier. This construct is equivalent to a table of “core data elements for a data-based transcript” (Kendall 2007) where each element is a cell in the table and contains a speaker reference, the textual representation of the markup, and a start and an end time. Manual annotations are notoriously hard to reproduce (McEnery 2006), so all phonetic mark up is done twice and proof-read by a third person to minimize inter-individual variation in the mark up.

All information about speakers (informants and field workers) is stored in a separate database. This ensures that informant data can be updated, for instance with information from the recordings themselves, without any change of transcripts.

#### **4. Comparability and the discourse context analysis**

Both S1 recordings and S2 ones have thus been made uniform in a number of respects. But the problem of comparing only similar types of intra-individual linguistic behavior cannot be solved by uniformity of transcripts alone. In order to control this type of variation in a data-based corpus, we developed a six-dimensional description of each interview or group recording called the Discourse Context Analysis (abbreviated DCA) (Gregersen, Beck Nielsen & Thøgersen 2009). A framework such as the DCA enables us to detect comparable passages in recordings which originally were made for different purposes and/or used different methods, so that our analyses compare only truly comparable sections from S1 and S2. Think of the DCA as equivalent to normalization of acoustic measurements of phonetic variables<sup>14</sup>: Normalization has to make measurements so much alike that they may be compared. When comparing we can then see the real differences.

In the DCA the six dimensions are the following:

- Type of Speech event (S);
- Activity type (A);

14 *F. Gregersen and M. Barner-Rasmussen*

- Type of Interaction structure (I);
- Type of Macro speech act (M);
- Type of speech Genre (G);
- and Enunciation (U, for Danish ‘udsigelse’).

The *Type of speech event* (S) dimension refers to the relationship between field worker and informant<sup>15</sup>, and to the various types of recording – single person interview, group interview (i.e., with an interviewer), and group discussion (i.e., without any interviewer). *Activities* (A) of various *types* are performed in the recordings, such as collecting data on the informant’s social background or that of other potential informants, conversation with the informant, or collecting language attitude data. The *Interactional* (I) dimension distinguishes five subcategories:

- I4<sup>16</sup>: departure from the regular question-answer format and the interviewer’s more or less pronounced and obvious monopoly on strong initiatives, i.e. more symmetrical interaction;
- I5: informant initiative with response by interviewer, i.e. reversal of interaction roles;
- I6: fight for the floor;
- I7: informant initiative with response by other informant;
- and I8: monologue.

The dimensions of Macro Speech Act (M) and Genre (G) will be detailed in Section 4.1 below. The dimension of *Enunciation* (U) refers to the fact that informants sometimes do not claim full responsibility for the words they utter but rather quote someone else – or repeat what the field worker has told them to say. In these cases we make a note of this in the annotation tier, so that we do not base our analysis of the informants’ phonetic patterns on e.g. passages which contain quotations.

All DCA coding is performed on the basis of the transcripts in order to avoid the vicious circle of being influenced by precisely those phonetic phenomena which we aim to explain. Two coders are involved. One coder first codes the whole recording either for Activity type, Interaction and Genre, or for Type of Speech Event, Macro Speech Act and Enunciation. A second coder then enters the codes into Praat, simultaneously checking that the codes are consistent with the annotator’s manual (which regrettably is only available in Danish for the moment). In order to prevent misunderstandings, we use CAPs in the following for all categories within a given dimension.

The DCA is relevant in and of itself and may e.g. be used to characterize the effects of different field methods on the data (Gregersen, Beck Nielsen and Thøgersen 2009) but it also informs the phonetic analysis: only passages delimited as belonging to the Macro Speech Act of EXCHANGE OF INFORMA-

TION (see below) were coded for phonetic variation. This Macro Speech Act was present in all recordings and thus maximized the data for analysis while still controlling for discourse context.

#### 4.1 *Genres (and Macro Speech Acts) in the DCA*

Although the concept of genre figures prominently in early sociolinguistic linguistic theory, e.g. as the final letter in the Hymes SPEAKING formula for the speech event (for a brief but very instructive review of the various traditions, see the entry 'genre' in Swann et al. 2004), genre is a notoriously difficult notion within sociolinguistics. The meaning attached to the term within the field of Systemic Functional Linguistics (Martin 1997, Martin & Rose 2003) differs from e.g. Labov (2001) where narratives are rather seen as one of many contexts giving rise to informal speech (cf. Labov 1966 [2006]: 70f). And corpus linguists such as Biber & Conrad (2009: 16) see Genre as "usually once-occurring in the text, in a particular place in the text". Finally, the literary traditions are solely concerned with written genres (Frow 2006).

In the LANCHART study we are concerned with recordings of speech. The informants all participate in a speech event which we may categorize as some sort of *conversation* in the sense of Biber & Conrad 2009 where this is seen as the basic interpersonal spoken *register*, or more specifically *a sociolinguistic interview*. The DCA Genre dimension works on passages *within* these interviews. The category of genre, in the DCA conceptualization, thus demands the application of structural criteria in order to capture units which are equivalent to, or even identical with, everyday or phenomenologically pertinent categories.<sup>17</sup> The eight categories of the Genre classification are: NARRATIVES, GENERAL ACCOUNTS, SPECIFIC ACCOUNTS, SOAP BOX, GOSSIP, CONFIDENCES, REFLECTIONS and JOKES. Of these eight categories, we will be especially concerned with *the narrative field* as outlined below, viz. NARRATIVES, GENERAL ACCOUNTS and SPECIFIC ACCOUNTS.

We wanted to be able both to compare our results to those of Labov in the area of narrative studies, and to broaden the notion of narrative. Hence, we sliced the narrative field into three categories, with our genre NARRATIVE (abbreviated Gna) being identical to the quintessential Labovian narrative, i.e. an account of a past event which the narrator has experienced him- or herself and actively taken part in and where the narrator views the past event as reportable and thus expects to make a point or entertain the audience with an account. A Labovian narrative also has a definite conclusion in response to the implicit 'so what?' that is the eternal threat for any story-teller.

Supplementing the Labovian NARRATIVE, the two genres SPECIFIC ACCOUNT (Gsr) and GENERAL ACCOUNT (Ggr) diverge in that the first one

does not concern reportable events (but otherwise conforms to narrative rules in telling about past events, most often in the temporal order in which the original events occurred) while the second, GENERAL ACCOUNT, diverges from the Labovian format in that events may be reportable or not, but they are told not as specific occurrences but as generic events, that is, as something that used to happen and most often would have this or that event structure. Thus all the three genres are narrative in kind, but only one of them is identical to the Labovian narrative, viz. NARRATIVE. Together they fill out the narrative field in the recordings.

The difference between the dimensions of Genre and Macro Speech Acts resides first in the way the analysis is applied to the data: The Macro Speech Act analysis is a coding of the sum total of exchanges in the speech event. All interaction which is there in the recording, and hence in the transcription, has to be classified as belonging to one of five categories: EXCHANGE OF INFORMATION, EXCHANGE OF ATTITUDES, EXCHANGE OF EMOTIONS, SPEECH ACCOMPANYING ACTION, and EXCHANGE OF FICTION. Since all exchange has to be classified as belonging to one and only one of these five categories, they are necessarily broad and not readily characterized by formal features. In this respect they differ from the genres where formal features are of the essence.

The DCA was construed in a continuous dialogue with a so-called ‘exploratory corpus’ of 21 recordings selected so as to maximize the contrasts within the sum total of S1s and S2s. The exploratory corpus is phonetically (and grammatically) coded from start to finish so that it is possible to conduct searches which detail the use of genres for those few informants who were included in that corpus. In total it contains 311,311 words, 40,308 manual annotations and 933,933 machine made annotations. Please bear in mind that the procedure used for the genre analysis automatically divides the transcripts into two main sections, viz. passages coded as belonging to a genre and passages which are outside any genre.

#### 4.2 *Our solution to the selection and probe problems*

At the end of Section 2 above, we distinguished two problems, *the selection problem*, viz. which occurrences to code, and *the probe problem*, viz. how many of them. Our solution to the selection problem was to choose only passages which were labeled as EXCHANGE OF INFORMATION by the Macro Speech Act analysis for phonetic analysis. This particular solution exploits the complementarity in the DCA of the dimensions of Macro Speech Act and Genre: It follows from the principles of the Discourse Context Analysis that Genre and Macro Speech Act are independent dimensions. Hence, different

subcategories of Genre may very well occur within a single Macro Speech Act, and in fact they do. In this case, various genres occur within the Macro Speech Act of EXCHANGE OF INFORMATION. The basic idea of this paper is thus to use the Macro Speech Act EXCHANGE OF INFORMATION as our first constraint on the comparison. Within this subset of the data we contrast the genres within the narrative field with the sum total of passages *outside any genre* in order to discuss whether there is any difference which is related to one or more Genre subcategories. We refine the searches successively in order to hone in on the difference between subcorpora and finish with the individual informant as a possible source of variation. For this purpose we use the exploratory files as a testing ground.

## 5. The sociolinguistic variable (æ)

To illustrate our comparison of genres with passages not coded as belonging to any genre, we require a sociolinguistic variable which has been demonstrated to be both socially and stylistically sensitive, and which is frequent enough to reveal systematic variation.

The most thoroughly investigated cases of variation in the sociolinguistic literature on Danish are the changes which have occurred in the realization of the short /a/ phoneme throughout the 20<sup>th</sup> century (see Brink & Lund 1975; Normann Jørgensen 1980; Holmberg 1991; Gregersen, Maegaard & Phrao 2009). Brink and Lund (1975) and Normann Jørgensen (1980) convincingly demonstrate that this variable was extremely important socially. In contrast, Holmberg (1991) shows that only a decade later, the variable was *not* socially determined but still stylistically active. This variable is thus perfect for an exploratory study.

A number of sources, primarily the thorough investigation by Brink and Lund (1975), indicate that Danish had only one short /a/ quality in the beginning of the 19<sup>th</sup> century (Gregersen 2009a). This was the 'normal' unrounded low vowel found in most European languages. Two splits, however, created the situation we investigate and which is here abbreviated as 'the (æ) variable'. The first split separated contexts before coronal consonants and syllable boundaries from all others. In this particular context the short /a/ was raised and realized as [æ], creating a textbook example of complementary distribution: One short /a/ phoneme manifested differently in complementarily distributed environments, viz. as [æ] and [ɑ].

A further split created the sociolinguistic differentiation between on the one hand a raised (æ), realized as [ɛ], which in the beginning of the 20<sup>th</sup> century was used primarily by the working class and by men, and on the other hand the

(æ) proper, i.e. realized as [æ]. Following Normann Jørgensen (1980) and Gregersen, Maegaard and Pharao (2009), we may reconstruct the development as follows. The working class male pronunciation caught on and was at the end of the 1970s close to becoming universal, when the change reversed and the values of raised (æ) fell to around 10% in Copenhagen. What caused this reversal, or even whether the reversal is in some measure an artefact related to our ways of judging segmental vowel quality auditorily, is still a topic for discussion, but there is no denying the fact that what used to be a working class (predominantly male) stereotype has now become a much less clear signal precisely because the use is not universal.

In this respect we may witness a development which is particularly interesting because the level of awareness for this particular feature of phonetic variation in Danish was and is so high that there is even a popular label for the raised variant of (æ), i.e. the so-called ‘flat a’. The historical nature of the variable as a WC shibboleth is at odds with the present distribution of variants, in which the raised variant is in no way characteristic of the working class pronunciation of the short (æ). The present day social meaning of the raised (æ) will be discussed in the remaining sections of this paper from the viewpoint of the dimension of Genre while still having an open eye towards possible influences from speaker variables such as class, geography and generation.

## **6. The LANCHART Work Bench: The Search Engine**

This leads us to briefly introduce and discuss the front-end part of the LANCHART corpus of data-based transcripts and annotations. Based on the Praat-file archive a searchable database of all core data elements (we call those slices hereafter) in all tiers in all transcripts is maintained based on the most recent revision of the archives (see illustration in Section 3.3). As the corpus grows and changes, a search may yield different results at different times, but the LANCHART Search Engine can reproduce these results by populating the back end database with the archive files as they were when the original search was carried out.

The front end user interface allows selection of subsets of the entire corpus based on study(ies), generation(s), speaker(s) etc. This has been very valuable because of the time requirements to produce new annotations. A small(er) set of interviews may be selected for annotation and subsequently searched to allow for quicker, more efficient hypothesis testing. This is the way the DCA has been planned and executed, and this paper represents a significant validation of its usefulness. We believe it is a central benefit of data-basing transcripts and annotations.

The screenshot shows a search interface with three main sections:

- Filer:** Contains radio buttons for 'Alle', 'Sæt', and 'Eneelt'. Below them are dropdown menus for 'Alle prioriterede filer', 'bysoc1-gf-ABK', and an empty dropdown.
- Informanter:** Contains radio buttons for 'Alle', 'Sæt', and 'Eneelt'. Below them are dropdown menus for 'Alle prioriterede telere' and an empty dropdown.
- Søgekriterier:** A larger section with a 'tilføj søgekriterier' button. It contains dropdown menus for 'fjern ortografi', 'fjern grammatik', 'man', 'GA', 'Hele ordet', and 'Del af ordet'. To the right of these are navigation icons: '+/-', '1', '2', '3', '4', '5', '+/-'.

At the bottom left of the 'Søgekriterier' section is a 'Søg' button.

Figure 2. Search criteria can look in adjacent or further removed slices from the 'base slice'. In the example shown, only places in the interview where the grammatical code 'GA' is found one or two slices removed from any occurrence of 'man' will be found. This is indicated by the tick in boxes 1, 2 and 3 above. Also note the selection options for interviews ('Filer'): all files, a predefined set of files or a single file, and informants (again: all informants, a set of informants or a single informant).

Searches may look for matches with the whole contents of a slice, a part of same, and even matching regular expressions; for example, searching for all cases and versions of a lemma can be done in the same search. Searches are also context sensitive, so one may look for concordances between the contents of multiple tiers at the same position or at a remove (see figure 2 below).

In the present paper we use this functionality to search for a phonetic variable both inside and outside the genres mentioned. Incidentally, this is a task enabled by this way of working with sociolinguistic data in that we did not plan for this particular set of searches when we prepared the DCA – data-basing allows for and encourages re-using data.

Results can be output either as simple counts of the number of occurrences found, as KWIC concordances (using html-tables), or as comma-separated value files containing the names of the files matched, all background information for the participants, as well as the content of every slice in every tier of the annotation file. This output may be imported into any number of post-processing applications for data analysis, variable tabulations and so on. Because of the ease of performing such searches in the data-based corpus, inquiries can start out as general searches and, through perusing the results, become more and more specific in an ongoing dialogue with the corpus. We can 'have a look'.

It is also noteworthy that since the search engine is using standard three-tier application design principles, the underlying database is directly queryable (using standard SQL queries), so that complicated searches can be performed even when the search engine's interface does not allow this (yet). The architecture allows for separate development efforts on front end, back end and archives so that new features can be added continually without 'breaking the app'.

## 7. Results

### 7.1 Do the LANCHART genres change in real time?

As detailed in Section 4.1 above, we search for the named genres NARRATIVES, GENERAL ACCOUNTS, and SPECIFIC ACCOUNTS, which we aggregate (using a regular expression in the search engine) as a single entity, the *narrative field*, and for the results of the codings of the ( $\text{æ}$ ) variable. We discard all double codings within the dimension of Genre and all realizations of the variable other than a straight raised ( $\text{æ}$ ) and a straight non-raised ( $\text{æ}$ ), i.e. we discard all tokens where the coders could not decide whether a given instance of ( $\text{æ}$ ) was raised or not. They are infrequent in the data and we could not see any pattern in the coding that would be an indication of an intermediate value having psychological, and hence also sociolinguistic, reality<sup>18</sup>. Consequently, documented in the tables below are relationships between auditorily coded phonetic variants within unambiguously assigned genre passages and the various speaker variables. Since we are concerned with change, we give the S1 results and the S2 results separately (tables 1 and 2 respectively).

In the first column we identify each project by location and generation (note that Copenhagen = Bysoc, The Copenhagen Urban Sociolinguistics Project). Thus Næstved 1 refers to the generation of informants from Næstved born between 1943 and 1963, and Odder 1 to the group of informants from Odder born between 1943 and 1963, while Odder 2 and Bysoc 2 indicate the younger generation of informants, born between 1963 and 1973. We detail the percentage of raised ( $\text{æ}$ )'s and give the total number of instances.

In Table 1, we note that only one of the data sets shows significant differences between the narrative field and passages outside genre, viz. Næstved 1. We also note significantly different generational patterns in the Odder data set,

Table 1. *Raised ( $\text{æ}$ ) percentages for the narrative field (NARRATIVES, GENERAL ACCOUNTS, and SPECIFIC ACCOUNTS aggregated) on the one hand, and in passages outside Genres on the other. This table is concerned with the S1 recordings only (1978ff); 85 recordings in total.*

Project and generation	Genres within the narrative field		Passages outside Genres		Level of significance
	Raised ( $\text{æ}$ ) percentage	N	Raised ( $\text{æ}$ ) percentage	N	
Bysoc 1	9%	346	9%	668	0.8983 ns
Bysoc 2	4%	162	7%	510	0.1885 ns
Næstved 1	25%	253	17%	519	0.0156*
Odder 1	18%	143	24%	232	0.2066 ns
Odder 2	0%	68	1%	278	1 ns <sup>19</sup>
TOTAL	13%	972	11%	2227	0.111 ns

Table 2. *Raised (æ) percentages for the narrative field (NARRATIVES, GENERAL ACCOUNTS, and SPECIFIC ACCOUNTS aggregated) on the one hand, and in passages outside Genres on the other. This table is concerned with the S2 recordings only (2005ff); 85 informants.*

Project and generation	Genres within the narrative field		Passages outside Genres		Level of significance
	Raised (æ) percentage	N	Raised (æ) percentage	N	
Bysoc 1	12%	284	9%	709	0.0758 ns
Bysoc 2	9%	231	7%	494	0.4005 ns
Næstved 1	15%	241	14%	594	0.4771 ns
Odder 1	15%	130	11%	298	0.3024 ns
Odder 2	4%	115	1%	351	0.1319 ns <sup>20</sup>
TOTAL	12%	1001	9%	2446	0.0091**

both in aggregated genres and in passages outside genre (both are significant at  $p < 0.0001^{***}$ , Fisher's exact test).

In Table 2, the picture is different. In the S2 recordings, narrative fields contain significantly more raised (æ) than passages outside genre. Furthermore, the distribution itself is new: the total for the S1 recordings is significantly different from the total for the S2 recordings with regard to passages outside genre ( $p = 0.0063^{**}$  Chi Square), the direction of change going toward fewer raised variants in real time. The rates of (æ) raising in the narrative field in the S1 recordings are not, however, significantly different from the same passages in the S2 recordings (see figure 3). This means that the pattern has reached some kind of predictability, which is indeed the case: all the data sets in the S2 recordings manifest more raised (æ) within the narrative field than in passages outside any genre. In Figure 3, this is illustrated by comparing the distance between the columns for the narrative field and passages outside Genres for the S1 and the S2 data sets.

On the whole, the S2 recordings feature less raised (æ)s. But the pattern of use for the (æ) variable has changed less within the narrative field than in the passages outside Genres making the difference significant in the S2 recordings. We take this as corroborating evidence that we have indeed captured a pattern of genre differences, viz. between the aggregated narrative genres and passages outside Genres.

Obviously aggregated measures in themselves do not convince us that all genres are sociolinguistically real in the sense that they are reflected in a specific pattern of the phonetic variable (æ). Perhaps the aggregated measure obscures rather than reveals the pattern. We have to look at the differences *within* the narrative field, i.e. the difference between each narrative genre.

We now take a look at the aggregated numbers for each genre separately in the S1 and the S2 recordings (compare Tables 3 and 4).

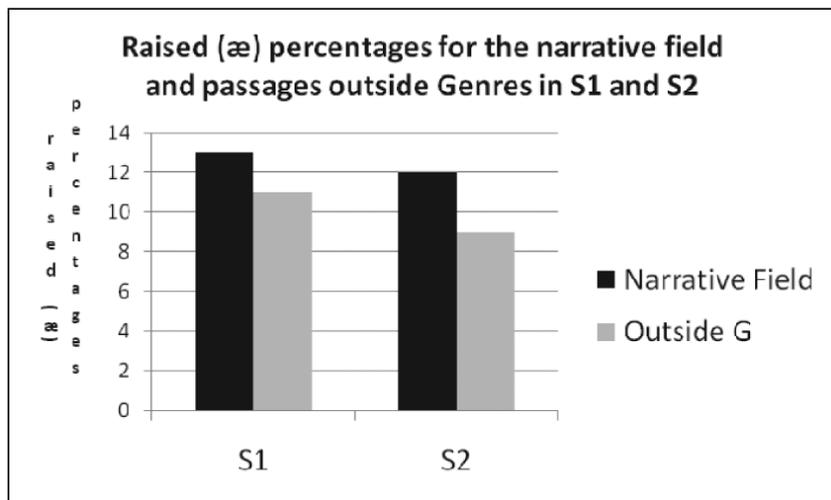


Figure 3. Raised (æ) percentages in the S1 and S2 recordings for the narrative field on the one hand and passages outside Genres on the other. The difference between the two columns in the S1 data set is not significant (Chi square:  $p = 0.111$ ) whereas the difference between the two columns in the S2 data set is significant at  $p = 0.0091^{**}$ .  $N$  is 3199 for S1 and 3447 for S2.

Table 3. Separate figures for raised (æ) percentages in the three narrative genres and passages outside Genres, only S1 recordings

data set	NARRATIVES		GENERAL ACCOUNT		SPECIFIC ACCOUNT		Outside G
	raised (æ)%	N	raised (æ)%	N	raised (æ)%	N	raised (æ)%
Bysoc 1	4%	57	10%	128	10%	161	9%
Bysoc 2	0%	21	5%	98	5%	43	7%
Næstved 1	25%	44	22%	107	27%	102	17%
Odder 1	50%	10	28%	47	9%	86	24%
Odder 2	0%	0	0%	46	0%	22	1%
TOTAL	14%	132	13%	426	13%	414	11%

Looking at the figures for the five data sets and the three genres separately we have the chance of judging whether the three narrative genres are geographically specific or generation specific, and whether they change over time or are stable. First, we observe a geographical difference: The figures follow the general trend in the LANCHART study. Raised (æ) is more frequent in Næstved 1 and Odder 1 than in the Copenhagen data set. Across generations, young informants in generation 2 have fewer raised variants than the older

Table 4. *Separate figures for raised (æ) percentages in the three narrative genres and passages outside Genres, only S2 recordings*

data set	NARRATIVES		GENERAL ACCOUNT		SPECIFIC ACCOUNT		Outside G
	raised (æ)%	N	raised (æ)%	N	raised (æ)%	N	raised (æ)%
Bysoc 1	12%	41	10%	123	11%	106	9%
Bysoc 2	11%	45	3%	63	11%	133	7%
Næstved 1	25%	8	18%	110	11%	133	14%
Odder 1	19%	37	4%	27	17%	66	11%
Odder 2	0%	11	10%	40	2%	64	1%
TOTAL	14%	132	11%	363	11%	502	9%

ones of generation 1, with the exception of Odder for GENERAL ACCOUNTS in the S2 recordings.

One genre seems to have become quite stable: In the S2 recordings the genre SPECIFIC ACCOUNT has figures of 11% except in Odder. This is very different from the percentages found in the S1 recordings. Perhaps this is an effect of differences between the S1 and the S2 in Odder. In both the S1 Odder data sets, SPECIFIC ACCOUNTS have significantly different values than the comparable figures for passages outside any genre, but the direction of divergence is different. In the older generation (Odder 1) there is *more* raised (æ) in SPECIFIC ACCOUNTS ( $p = 0.0272^*$ ). In the younger generation (Odder 2) there is significantly *less* ( $p = 0.0042^{**}$ ), and the level of raised (æ) in the SPECIFIC ACCOUNTS also diverges from that in the other genres (the difference between GENERAL ACCOUNTS and SPECIFIC ACCOUNTS is significant at  $p = 0.0055^{**}$ ). This asymmetry may result from the fieldwork technique used in the S1 Odder recordings, where two field workers interviewed one informant. Presumably the presence of two field workers works differently with the younger informants than with the middle-aged. Or it may simply be the case that SPECIFIC ACCOUNTS have different functions for the younger informants. The effect has disappeared in the S2 Odder recordings, where another fieldwork technique with only one interviewer was used throughout.

The genre GENERAL ACCOUNT seems to be at the other end of the stability continuum, with figures ranging over a much wider area, yet the total ends up at precisely the same figure as SPECIFIC ACCOUNT. In the S2 recordings, the genre GENERAL ACCOUNT has changed places with SPECIFIC ACCOUNT, so to speak, in that the level of raised (æ) is significantly lower than it was in the S1 recordings ( $p = 0.0129^*$ ). We shall come back to this instability below.

The genre NARRATIVE alone is responsible for the significant difference between S1 and S2 recordings in Tables 1 and 2, since the rate of raised (æ) is

completely stable (when looking at the TOTALs) for the NARRATIVE, while the rates of raised (æ) are dropping both for the two other narrative genres and for the passages outside any genre. Recall that the difference between the two figures for passages outside Genres is significant at  $p = 0.0063^{**}$ . The difference between the figures for NARRATIVES and those for outside Genres are significant at  $p = 0.0487^*$ . The genre NARRATIVE, i.e. the Labovian personal narrative, is apparently different enough from the two other narrative genres to distinguish itself *within* the narrative field by staying at the same level of (æ) raising.

We conclude that there is some evidence that genres as defined in the DCA are responsible for significant differences in the quantitative patterning of the (æ) variable. Figures for specific genres can tell us more about which passages have a divergent pattern of raised (æ) use in contrast to passages outside any Genre. Also we seem to find evidence that a specific field work technique might have an effect on the use of genres.

In this sense the genres delimited by the DCA are sociolinguistically real.

### *7.2 Are genres 'sociolinguistically real' at the level of the individual informant in a single recording?*

In this section we investigate whether genres are 'real' in another sense. We focus on the individual and ask whether the passages assigned to specific genres are significantly different as to (æ) variation in the course of one session with one informant. All factors have been kept constant except genre. For this purpose, we exploit the exploratory sample, where recordings were coded from start to finish, giving us an opportunity to contrast passages belonging to different genres.

Though the (æ) variable has been selected for this investigation precisely because of its high frequency, it is not the case that there are instances of (æ) in all the various genres which the coding apparatus allows. But a number of genres that we have not mentioned before will be represented in the tables, such as SOAP BOX (abbreviated below as Gsb) (Labov 2001: 91), consisting of general remarks of a political nature in a broad sense ranging from local school politics to opinions about presidential elections, GOSSIP (Gsl), and REFLECTIONS (Gre), consisting of expressed thoughts on how the informant's life has developed in general.

For one of the Middle Class men we have three recordings in the exploratory sample: an S1 group interview together with the informant's father and brother, the S1 single person interview, and the corresponding S2 interview. In this case there is no relevant difference between the field methods employed in the S1 (Gregersen, Albris and Pedersen 1991) and the S2 recordings.

The relevant figures are given in the three tables 5–7 below. In all cases the informant is a Middle Class male Copenhagenener. The various Genres are ab-

Table 5. *Raised (æ) percentages for passages inside Genres and outside for a Middle Class male Copenhagener in the S1 group conversation.*

Genre	raised (æ)%	N
Gna	26%	118
Ggr	20%	25
Gsr	23%	13
Gsb	23%	31
Gsl	24%	33
INSIDE G	25%	220
Outside G	26%	272
<i>TOTAL</i>	25%	492

Table 6. *Raised (æ) percentages for passages inside Genres and outside for the same Middle Class male Copenhagener in the S1 sociolinguistic interview.*

Genre	raised (æ)%	N
Gna	42%	26
Ggr	8%	34
Gsr	30%	53
Gsb	44%	23
Gsl	33%	3
Gre	44%	9
INSIDE G	30%	148
Outside G	28%	458
<i>TOTAL</i>	29%	606

Table 7. *Raised (æ) percentages for passages inside Genres and outside for the same Middle Class male Copenhagener in the S2 sociolinguistic interview.*

Genre	raised (æ)%	N
Gna	30%	54
Ggr	40%	70
Gsr	31%	35
Gsb	33%	42
Gre	31%	13
INSIDE G	34%	214
Outside G	28%	593
<i>TOTAL</i>	29%	807

breviated as follows: Gna = NARRATIVE; Ggr = GENERAL ACCOUNT; Gsr = SPECIFIC ACCOUNT; Gsb = SOAP BOX; Gsl = GOSSIP. INSIDE G contains the sum total of all passages labeled as a Genre while Outside G gives the sum total of all passages not assigned to any Genre.

It is fairly remarkable how stable this informant is. He seems not to distinguish too much, if at all, between genres. His figures for (æ) raising are stable, thus lending some credence to the much debated doctrine of Labov's that we may indeed find that the vernacular used in the sphere of intimacy is more systematic than the linguistic behavior we tap into when we tape.

In Table 6, the figures for (æ) raising in the genre GENERAL ACCOUNT diverges drastically from the figures of the other genres. The difference between GENERAL ACCOUNT and NARRATIVES is significant at  $p = 0.0045^{**}$  (Fisher's exact test) and the difference between GENERAL ACCOUNT and SOAP BOX is significant at  $p = 0.0035^{**}$  (Fisher's exact test). In other words, this interview recording features a rather systematic distribution of raising, too. Although the percentages are in general higher than those of the group discussion, it is only the extremely low figure for the genre GENERAL ACCOUNT which destroys the generalization that the raised (æ) is a characteristic of passages labeled as belonging to a genre and thus may be seen as some sort of signal for genre. What is there about GENERAL ACCOUNT that is distinct?

While GENERAL ACCOUNT manifested the *lowest* raised (æ) percentage of all genres in both S1 recordings, especially in the S1 interview, it is the genre with the *highest* percentage of raised (æ) in the S2 interview. This difference across recordings is significant at  $p = 0.0011^{**}$ . Why?

There are three possibilities. First it could be a fact about the phonetic variable. In the S1 interview what is significant is the absence of raised variants in passages labeled as GENERAL ACCOUNT. But perhaps the non-raised variant is no longer used to make the genre of GENERAL ACCOUNT distinct. Other variables could have taken its place. Second it could be a fact about the genre. We may say that the genre GENERAL ACCOUNT has lost its special status with respect to this variable and now patterns like the other passages belonging to a genre. The high rate of raised (æ)s serves to signal a passage as belonging to a genre instead of being outside. It might also be the case that this is a fact about the individual. This is indeed made likely by a look at the figures for the Copenhagen data set in general. When we look at the figures in Table 3 in Section 7.1 above, we find that the raised (æ) percentages are much lower than the rates we find with this informant. And there is no special status for GENERAL ACCOUNT.

Even stronger evidence for the reality of genres versus passages outside any genre may be had from the other informants in the exploratory sample, in particular a married couple from Næstved 1 who participated both in the S1 and the S2 recordings. For the four possible cases (two informants in S1 and S2 recordings), the rates of raised (æ) inside passages belonging to a genre and in passages outside Genres are significantly different in three:

- WC male S1: difference between INSIDE G and Outside G;  $p = 0.0644$  ns
- WC male S2: difference between INSIDE G and Outside G;  $p = 0.0468^*$

- MC female S1: difference between INSIDE G and Outside G;  $p = 0.0212^*$
- MC female S2: difference between INSIDE G and Outside G;  $p = 0.0005^{***}$

In sum, the figures lend some credence to our claim that genres are also real at the level of the individual informant. The categories seem to capture at least some of the intra-individual variation present in a recording. But if so, they are not necessarily stable in time – at least not with regard to this variable. There is one possible generalization, though: The pattern is in all cases, except one, that the rates of raised ( $\text{æ}$ ) outside Genres are *lower* than INSIDE GENRE. Thus we may conclude that the raised ( $\text{æ}$ ) is most often used to signal, or is triggered by, a passage belonging to a genre, any genre.

### 7.3 *Are the probes used in the main study representative for the whole recording?*

In this section we shall use the exploratory sample to study the relationship between the various genres, the use of raised ( $\text{æ}$ ), and the individual recordings. As mentioned above (Section 4.2) our solution to the *selection problem* involved coding at least 40 instances of any phonetic variable within the Macro Speech Act of EXCHANGE OF INFORMATION. This constitutes what we call *the probe* for all informants. The results in Section 7.1 are based on this probe.

But since the exploratory sample was coded from start to finish, it is possible for us to contrast the probe taken for each of the 5 individuals who are both in the exploratory and the total sample, with the countless occurrences in the recording from which it was taken in the first place. For three of these individuals, a MC male from Copenhagen Bysoc 1 (*informant 1*), a WC male from Næstved 1 (*informant 2*), and a MC female from Næstved 1 (*informant 3*) (the latter two are a married couple), there is both an S1 and an S2 recording; whereas for another Næstved 1 married couple, both of them from the WC (*informant 4* being the female while the male is *informant 5*), there is only an S2 recording. Our wish is for the probe figures in table 7 below to be equal to or at least not significantly different from those of the *TOTAL*. In that case we may trust the probe to be a faithful picture of the informant's production of the ( $\text{æ}$ ) variable.

The result, evident from Table 8 is that in 2 of the 8 possible cases, the probe is significantly different from the exploratory coding. In the first case, the WC male from Næstved 1, we can see that the figure is too high, thus giving us the impression that the person in question has changed quite drastically in real time. The crucial fact, however, is that he has actually changed in real time, although not that much. The difference between the two totals is significant at  $p = 0.0149^*$ . The second case is the WC male from Næstved 1, who only

28 *F. Gregersen and M. Barner-Rasmussen*Table 8. *Comparison of the figures for the exploratory sample (coded from start to finish) and the probes used for the phonetic study of change in real time for the 5 informants and 8 recordings where this is possible.*

RECORDING	Informant	Genre	raised (æ)%	N	
S1	1	INSIDE G	30%	148	$p = 0.2638$ ns
		Outside G	28%	458	
		TOTAL	29%	606	
		Probe	21%	42	
S2	1	INSIDE G	34%	214	$p = 0.0712$ ns
		Outside G	28%	593	
		TOTAL	29%	806	
		Probe	41%	51	
S1	2	INSIDE G	9%	109	$p = 0.4493$ ns
		Outside G	18%	82	
		TOTAL	13%	191	
		Probe	17%	46	
S2	2	INSIDE G	27%	158	$p = 0.0001$ ***
		Outside G	19%	515	
		TOTAL	21%	673	
		Probe	48%	46	
S1	3	INSIDE G	39%	188	$p = 0.7716$ ns
		Outside G	26%	109	
		TOTAL	34%	297	
		Probe	36%	47	
S2	3	INSIDE G	29%	129	$p = 0.5238$ ns
		Outside G	15%	376	
		TOTAL	18%	505	
		Probe	23%	40	
S2	4	INSIDE G	27%	115	$p = 0.0522$ ns
		Outside G	27%	268	
		TOTAL	27%	383	
		Probe	14%	49	
S2	5	INSIDE G	9%	299	Fisher: $p = 0.0251$ *
		Outside G	9%	256	
		TOTAL	9%	555	
		Probe	0%	45	

participates in the S2 recordings (informant 5) where the probe underestimates his number of raised (æ)s.

Looking at the figures for genre passages (INSIDE G) vs. figures for Outside G we note a curious pattern in the four Næstved recordings which feature informants 2 and 3: In the S1 recordings the husband and wife diverge so that the husband (a WC male) has very few instances of raised (æ) in his genre passages whereas the female part (an MC female) has a lot more raising INSIDE

Gs than outside. In the S2 recordings, however, the pattern is completely different: The male has the 'normal' pattern of more raising INSIDE GENRES than outside as does the female, although she has less raising altogether. They converge with time. The probe captures this nicely for the female but it fails miserably in the case of the male. In general the probe seems to be appropriately placed between the values for passages labeled as belong to a genre and passages not belonging to any. We do want the probe to capture the fact that not all passages belong to a genre.

We conclude that the probes are valid indicators of the total sample in 6 out of 8 cases. In the two cases where the probes differ significantly from the total, they go in the right direction: Informant 2 actually has a significantly new distribution in the S2 recordings – it is just not that different. And informant 5 actually has a very low amount of raised (æ)s, but he has some.

## 8. Discussion

The genre 'narrative' figures prominently in the time honored discussion of intra-individual variation, once known as the discussion of 'style' or 'stylistic shifting'. The *locus classicus* for this discussion is of course originally Labov (1966), which however reached a wider audience as Chapter 3 of Labov (1972). In this chapter, Labov first formulated the observer's paradox and then pointed out ways to circumvent it. Among the contexts which an interviewer might use to elicit a style that is close to the vernacular, while not being identical with it, is the 'danger of death' question (see Section 2 above). Note that Labov is concerned with the field methods appropriate to elicit variation within a single setting, namely the sociolinguistic interview, in order to get at the informant's casual speech; and with methods to determine quantitatively whether or not we had succeeded in doing so.

As we noted above, this type of style discussion was more or less superseded when Eckert (2001) and Coupland (2001, 2007), coming from different corners of the field, introduced a view of stylistic practices which focuses on the individual's use of stylistic resources at his or her disposal to create a specific persona in a particular context. The relationship between these two takes on style seems to be that the latter presupposes the former: an interpretation of a particular pattern as projecting a certain persona must at least to some extent rely on the pattern manifested being actually typically associated with persons who act or typify the *persona*<sup>21</sup> (cf. Macaulay 2009, Bell 1984). And it works better with specific forms than with quantitative patterns of variants. In this view of the discussion, the Labovian take focuses on quantitative patterns of variants as evidence of intra-individual variation, while the Eckert-Coupland view instead focuses qualitatively on the social meaning that is evoked in the

listener/analyst by the informant's use of a particular stereotypical variant. In doing so, the latter view ascribes intention to the speaker to signal or project a specific identity (Schilling-Estes 2004: 378). Perhaps the distinction between a deterministic and a voluntaristic approach to linguistic variation (Gregersen 2005) may shed some light on this issue. We see the Scylla of determinism as lurking behind the Labovian take (Schilling-Estes 2004: 383) but the Charybdis of voluntarism as equally threatening for the Eckert-Couplandian view.

In our view, the introduction of the Discourse Context Analysis, and in particular the introduction of a level of analysis based on the notion of genre, allows us to see intra-individual variation as partly determined by variation in the ascription of passages during an unclear speech event like 'the sociolinguistic interview' to various structured sub-events, namely genres. Genres, as Bakhtin put it, "are the drive belts from the history of society to the history of language" (Bakhtin 1986: 65) and thus have a connection 'upwards' toward spheres of society. On the other hand, speech genres are sub-routines embedded in a flow of conversation which serves to establish the base line from which the pattern in the genre may, and will often, diverge (consider the notion of passages 'outside genres' above). Genres are intentional acts. Thus it is actually the speaker's intention to signal that a specific instance of a particular genre is about to become invoked. This is done thematically and structurally but it seems that frequency of particular phonetic variables such as the raised (æ) is also used as a cue – or as an effect<sup>22</sup>.

The advantage of the notion of style was that it was convenient shorthand to discuss the notion of good data and the 'vernacular' which in Labov's classic formulation was seen as more systematic than the linguistic behavior encountered with an open microphone, hence the observer's paradox. The notion of style has turned out to be less than convenient, though, in discussions of what happens in other mundane situations with language behavior, in particular distributions of variants. On the one hand, the notion of context invokes notions of norms, especially when used to characterize reading passages and the reading aloud of word lists in contrast to casual speech. Here societal pressures are mediated through education in ways which must be culturally sensitive in order to do justice to differences in teaching practice, amount of schooling etc. (Baugh 2001). On the other hand, the difference between casual and non-casual speech is believed to be derived from norming behavior of in principle the same type, here, however, as mediated by the monitor. The monitor is either on, meaning attention is being paid to speech, or less on (or even off), meaning that attention presumably is elsewhere, in any case not on speech (Schilling-Estes 2004: 378f).

But this reduction of the intra-individual variation to a single dimension derivable from attention to speech does not tally with the fact that e.g. narratives often are told with as much attention to speech as reading aloud (Butters 2000).

On the whole, we shall not succeed in doing justice to intra-individual variation if we stick to one dimension only. With the DCA we can characterize any number of combinations. Here we have only focused on genre because it is a comparatively simple notion, relatively easy to operationalize, but it is essential to the whole conception presented here that it is only one dimension out of six.

Nevertheless, we believe that the notion of genre as embedded in the Discourse Context Analysis makes it possible to characterize an important factor in the range of intra-individual variation so typical of sociolinguistic data. With the DCA we may capture the many ways in which speech in the sociolinguistic interview situation emulates speech in situations not readily available to sociolinguists.

Finally, we hope to have documented that a type of analysis such as the Discourse Context Analysis is essential to consider in order to accomplish comparability across data sets or corpora. Comparability cannot be had at the level of entire speech events such as the sociolinguistic interview – even if it were possible to code all instances of frequent phonetic variables (which it is not). Furthermore, a specific part of the DCA, the Genre analysis, may be seen as placed between the in principle deterministic conception of context-as-determinants-of-intra-individual-variation and the in principle voluntaristic projection of a persona by the use of specific (patterns of) variants. If we have contributed to uniting the quantitative analysis of patterns of phonetic variants with the qualitative interpretation of a specific instance of use, we have accomplished our goal.

## **Bionotes**

Frans Gregersen is professor of Danish language at the University of Copenhagen and director of the DNRF LANCHART Centre. His research focuses on sociolinguistic issues in the Danish speech community in general, ranging from language politics to the analysis of phonetic innovations. Central is the problem of linguistic change and its embedding in the history of the community. His recent publications include the collection of papers from the LANCHART centre which marks volume 41 of the International Journal *Acta Linguistica Hafniensia* (2009) and a collection of his own papers called *Københavnsk sociolingvistik* (Oslo: Novus 2009) (most of them are in Danish but some are in English). He has edited volume 32, 1 of *Nordic Journal of Linguistics* focussing on sociolinguistics (with Unn Røyneland) and has just delivered the introduction to a collection of papers from the ICLaVE 5 in Copenhagen to John Benjamins for publishing later this year. E-mail: fg@hum.ku.dk

Michael Barner-Rasmussen is IT officer and as such in charge of all dedicated programmes used at the DNRF LANCHART Centre. Recently, he has

completed the new search engine format and is collaborating with the centre staff on developing it further. He is also active in the Danish CLARIN project and has collaborated with the Department of Scandinavian Research on developing a base of Runic inscriptions (<<http://runer.ku.dk>>) and on creating a digital atlas of the Danish historical-administrative geography (<<http://www.digdag.dk/index.php?lang=en>>). His publications focus on Digital Humanities and include a 2009 paper, “The LANCHART Search Engine—Making important progress in data and data archiving reuse”, presented at the Digital Humanities 2009 Conference. E-mail: [mbr@hum.ku.dk](mailto:mbr@hum.ku.dk)

## Notes

1. This paper would not have been half as readable without the generous help of the two editors, Tyler Kendall and Gerard Van Herk who made us clarify a lot of issues. We also gratefully acknowledge help from two anonymous reviewers. For what remains on the pages after this process we remain fully responsible.
2. This crucially depends on whether we study language acquisition (in which case any longitudinal study qualifies as a study of real time (acquisitional) change) or rather language change in communities. In the latter case we study (a number of) adults who have passed any so-called critical age limit for acquisition. Here we follow standard practice in only calling the latter studies studies of ‘change in real time’.
3. How large the time distance between an S1 and an S2 has to be, is a matter for discussion. In panel studies of young informants who seem to go through several significant ‘ages’ in rapid succession, i.e. developing from youngster to student and on to e.g. married couples, the time distance between an S1 and an S2 may be rather short as demonstrated by Hernes (2007).
4. The obvious answer is the microphone: the observer is present in the guise of a mic and the observer’s paradox is still valid.
5. This is the main reason why we have not worked with a fixed set of themes. Rich topics vary with geography and class culture even within such a homogeneous community as Denmark.
6. One of the two anonymous referees has kindly directed our attention to Butters (2000). In this interesting paper Ronald Butters argues that the meeting between an interviewer and an informant may be framed in terms of visitor-guest or intruder-private citizen dimensions in a way that prohibits danger of death questions as a possible rich subject. The question may simply call forth too strong memories for the informant to cope with and/or the interviewer may be ill at ease so that the whole passage becomes awkward. This is a long discussion, cf. Wolfson (1976), Møller (1991), but partly different from the one we wish to contribute to here. What we argue is that culture and history may make specific questions less useful than they were once. What Butters and Wolfson argue is that interviewers should be aware of the fact that they may be trespassing, or rightly or wrongly feel that they are trespassing, on ground where angels fear to tread.
7. The LANCHART centre is financed by the Danish National Research Foundation through a grant 2005–2013 to the first author. We gratefully acknowledge the support of the DNRF. The LANCHART acronym stands for Language CHange in Real Time, cf. <<http://www.lanchart.dk/>> (in English) and <<http://www.dgcss.dk/>> (in Danish).
8. There are 13.4 million time-stamped annotations in the corpus besides the 5+ million words.

9. This is actually only relevant for the Vinderup study since the generation 1 informants in Vinderup firstly were not necessarily past the critical age when they were recorded the first time and secondly were recorded around 30 years later in the S2.
10. As a way of getting both age stratification and data for future studies, we have included recordings of a third generation at all four sites; namely, those who were in the eighth or ninth grade in public comprehensive schools at the time of the S2 re-recordings. These generation 3 data will not be reported on here since they are not relevant for the panel study (but more so for the trend study).
11. The four sites include one metropolis (Copenhagen), one major city under the dominance of the metropolis (Næstved), and two smaller provincial cities: one that is very close to a larger conurbation into which it is being partly integrated (Odder, near Aarhus), and one that is more isolated and located in a traditionally dialect-speaking area (Vinderup). This geographical dimension in the corpus makes it possible to test models of diffusion from a norm center, in this case Copenhagen. That Copenhagen, and only Copenhagen, functions as the norm center for Danish speech has been documented in language attitude tests carried out at all the LANCHART sites (Kristiansen 2009).
12. Downloadable at <<http://sourceforge.net/projects/trans/files/>>.
13. Download and info at <<http://www.fon.hum.uva.nl/praat/>>.
14. cf. <[http://ncslaap.lib.ncsu.edu/tools/norm/about\\_normalization1.php](http://ncslaap.lib.ncsu.edu/tools/norm/about_normalization1.php)>
15. They are: unknown to each other, known to each other by a friend-of-a-friend (Milroy 1980) or known directly to each other.
16. The categories of I1 through 3 were discarded for lack of consistency in the annotations.
17. We have had to solve the problem of whether one of these categories may be embedded in another (e.g. a JOKE inside a NARRATIVE (in that case the joke is coded as occurring between the two parts of the narrative)) or whether they are mutually exclusive. The answer is that we have operationalized them as mutually exclusive after a tryout period where we allowed multiple codings. The unfortunate result was a proliferation of categories which were of no significance quantitatively and hence the principle had to be that the categories have to be coded *as if* they were mutually exclusive. This did not cause any problems.
18. In contrast to another phonetic variable, viz. the raising of epsilon before the velar nasal, where the codings of intermediate variants turned out to be generation specific.
19. In this case Fisher's exact test was used because of the low number of cases.
20. In this case Chi square was used with Yates' correction.
21. This depends on whether the *persona* has any relation to the speaker variables we use in the quantitative studies.
22. The discussion of which way the causal relation operates is raised by Finegan and Biber (2001).

## References

- Albris, Jon. 1991. Style analysis. In Frans Gregersen & Inge Lise Pedersen (eds.), *The Copenhagen Study in urban sociolinguistics, Part 1*, 45–106. Copenhagen: C. A. Reitzel.
- Bakhtin, M. M. 1986. *Speech genres and other late essays*. Translated by Vern W. McGee, edited by Caryl Emerson and Michael Holquist. Austin: University of Texas Press.
- Basbøll, Hans. 2005. *The phonology of Danish*. Cambridge: Cambridge University Press.
- Baugh, John. 2001. A dissection of style-shifting. In Penelope Eckert & John R. Rickford (eds.), *Style and sociolinguistic variation*, 109–118. Cambridge: Cambridge University Press.
- Biber, Douglas. 1988. *Variation across speech and writing*. Cambridge: Cambridge University Press.

- Biber, Douglas. 1995. *Dimensions of register variation. A cross linguistic comparison*. Cambridge: Cambridge University Press.
- Biber, Douglas, Susan Conrad & Randi Reppen. 1998. *Corpus Linguistics. Investigating Language Structure and Use*. Cambridge: Cambridge University Press.
- Biber, Douglas & Susan Conrad. 2009. *Register, Genre and Style* (Cambridge Textbooks in Linguistics). Cambridge: Cambridge University Press.
- Blake, Renée & Meredith Josey. 2003. The /ay/ diphthong in a Martha's Vineyard community: What can we say 40 years after Labov? *Language in Society* 32(4). 451–484.
- Brink, Lars and Jørn Lund. 1975. *Dansk Rigmål I–II* [Danish Standard Language]. København: Gyldendal.
- Butters, Ronald R. 2000. Conversational anomalies in eliciting danger-of-death narratives. *Southern Journal of Linguistics* 24(1). 69–81.
- Coupland, Nikolas. 2001. Language, situation and the relational self: theorizing dialect-style in sociolinguistics. In Penelope Eckert & John R. Rickford (eds.), *Style and Sociolinguistic Variation*, 185–210. Cambridge: Cambridge University Press.
- Coupland, Nikolas. 2007. *Style: Language variation and identity*. Cambridge: Cambridge University Press.
- Eckert, Penelope. 1997. Age as a sociolinguistic variable. Florian Coulmas (ed.), *The Handbook of Sociolinguistics*, 151–167. London: Blackwell.
- Eckert, Penelope. 2001. Style and social meaning. In Penelope Eckert & John R. Rickford (eds.), *Style and sociolinguistic variation*, 119–126. Cambridge: Cambridge University Press.
- Finegan, Edward & Douglas Biber. 2001. Register variation and social dialect variation: The Register Axiom. In Penelope Eckert & John R. Rickford (eds.), *Style and Sociolinguistic Variation*, 235–267. Cambridge: Cambridge University Press.
- Frow, John. 2006. *Genre. The new critical idiom*. Oxford: Routledge.
- Gregersen, Frans. 2005. Mellem sociolingvistisk determinisme og voluntarisme [Between sociolinguistic determinism and voluntarism]. In Gudbrand Alhaug, Endre Mørck & Aud-Kirsti Pedersen (red.), *Mot rikare mål å trå*. Festskrift til Tove Bull, 229–240. Oslo: Novus forlag.
- Gregersen, Frans. 2009. The data and design of the LANCHART study. *Acta Linguistica Hafniensia* 41. 3–29.
- Gregersen, Frans. 2009a. Hvad ved vi – om det såkaldt “flade” a? [What do we know about the so-called ‘flat’ (æ)?] In Ken Farø et al. (red.), *Sprogvidenskab i glimt. 70 tekster om sprog i teori og praksis*, 17–25. Odense: Syddansk Universitetsforlag.
- Gregersen, Frans, Jon Albris & Inge Lise Pedersen. 1991. Data and design of the Copenhagen study. In Frans Gregersen and Inge Lise Pedersen (eds.), *The Copenhagen Study in Urban Sociolinguistics, Part 1*, 5–43. Copenhagen: C. A. Reitzel.
- Gregersen, Frans, Søren Beck Nielsen & Jacob Thøgersen. 2009. Stepping into the same river twice: on the discourse context analysis in the LANCHART project. *Acta Linguistica Hafniensia* 41. 30–63.
- Gregersen, Frans, Marie Maegaard & Nicolai Pharao. 2009. The long and short of (æ)-variation in Danish – a panel study of short (æ)-variants in Danish in real time. *Acta Linguistica Hafniensia* 41. 64–82.
- Hernes, Reidunn. 2006. *Talemål i endring? Ein longitudinell studie av talemåsutvikling og språkleg røyndomsoppfatning hjå ungdomar i Os*. [Spoken language in change? A longitudinal study of young speakers from OS's spoken language development and perception of reality], Dr. art.-avhandling, Bergen: Nordisk Institutt, Universitetet i Bergen.
- Holmberg, Henrik. 1991. The sociophonetics of some vowel variables in Copenhagen speech. In Frans Gregersen & Inge Lise Pedersen (eds.), *The Copenhagen Study in Urban Sociolinguistics, Part 1*, 107–239. Copenhagen: C.A. Reitzel.

- José, Brian. 2010. The Apparent-Time construct and stable variation. Final /z/ devoicing in North-western Indiana. *Journal of Sociolinguistics*, 14(1). 34–59.
- Jørgensen, J. Normann. 1980. Det flade a vil sejre ['The raised (æ) will be victorious'], *SAML* 7. 67–124. København: Institut for Anvendt og Matematisk Lingvistik.
- Kendall, Tyler & Amanda French. 2006. Digital Audio Archives, Computer-Enhanced Transcripts, and New Methods in Sociolinguistic Analysis. *DIGITAL HUMANITIES 2006 The First ADHO International Conference Abstracts*. 110–112. <<http://www.allc-ach2006.colloques.paris-sorbonne.fr/DHs.pdf> 110-112>.
- Kendall, Tyler. 2007. Enhancing Sociolinguistic Data Collections: The North Carolina Sociolinguistic Archive and Analysis Project. *Penn Working Papers in Linguistics* 13(2). 15–26.
- Kendall, Tyler. 2008. On the history and future of sociolinguistic Data. *Language and Linguistics Compass* 2(2). 332–351.
- Kensing, Finn, Jesper Simonsen & Keld Bødker. 1998. MUST – a Method for Participatory Design. *Human-Computer Interaction*, 13(2). 167–198.
- Kristiansen, Tore. 2009. The Macro-level social meanings of late-modern Danish accents. *Acta Linguistica Hafniensia* 41. 167–192.
- Labov, William. 1966. *The social stratification of English in New York City*. Washington: Center for Applied Linguistics.
- Labov, William. 1972. *Sociolinguistic patterns*. Philadelphia: University of Pennsylvania Press.
- Labov, William. 1972a. Some principles of linguistic methodology, *Language in Society*, 1. 97–120.
- Labov, William. 1984. Field Methods of the project on linguistic change and variation. In John Baugh & Joel Sherzer (eds.), *Language in Use*, 28–53. Englewood Cliffs: Prentice-Hall.
- Labov, William. 1994. *Principles of Linguistic Change. Volume 1: Internal Factors*. Oxford: Blackwell.
- Labov, William. 2001. The anatomy of Style Shifting. In Penelope Eckert & John R. Rickford (eds.), *Style and sociolinguistic variation*, 85–108. Cambridge: Cambridge University Press.
- Labov, William. 1966 (2006). *The social stratification of English in New York City. Second edition*. Cambridge: Cambridge University Press.
- Macaulay, Ronald. 2009. Review of Nikolas Coupland. 2007. Style: Language variation and identity, *Language in Society* 38(1). 119–122.
- Martin, Jim R. 1997. Analyzing genre: functional parameters. In Frances Christie & J. R. Martin (eds.), *Genre and institutions, Social processes in the workplace and school* (Open Linguistics Series), 3–39. New York: Cassell.
- Martin, Jim R. & David Rose. 2003. *Working with Discourse. Meaning beyond the clause* (Open Linguistics Series). London/New York: Continuum.
- McEnery, Tony, Richard Xiao & Yukio Tono. 2006. *Corpus based language studies: An advanced resource book*. (Routledge Applied Linguistics Series). London: Routledge/Taylor and Francis.
- Milroy, Lesley. 1980. *Language and Social Networks* (Language in Society series 2). Oxford/New York: Blackwell.
- Møller, Erik. 1991. Narratives in the sociolinguistic interview. In Frans Gregersen & Inge Lise Pedersen (eds.), *The Copenhagen Study in Urban Sociolinguistics, Part 2*, 241–335. Copenhagen: C.A. Reitzel.
- Pope, Jennifer, Miriam Meyerhoff & D. Robert Ladd. 2007. Forty years of language change on Martha's Vineyard. *Language* 83(3). 615–627.
- Sankoff, Gillian. 2005. Cross-sectional and longitudinal studies in sociolinguistics. In Ulrich Ammon, Norbert Dittmar, Klaus J. Mattheier & Peter Trudgill (eds.), *Sociolinguistics/ Soziolinguistik: An international Handbook of the Science of Language and Society*, 2nd edition. 1003–1012, Berlin: Mouton de Gruyter.

36 *F. Gregersen and M. Barner-Rasmussen*

- Schilling-Estes, Natalie. 2004. Investigating Stylistic Variation. In J. K. Chambers, Peter Trudgill & Natalie Schilling-Estes (eds.), *The Handbook of Language Variation and Change*, 375–401. Oxford: Blackwell.
- Swann, Joan, Ana Deumerts, Theresa Lillis & Rajend Meshtrie. 2004. *A dictionary of sociolinguistics*. Edinburgh: Edinburgh University Press.
- Wolfson, Nessa. 1976. Speech Events and Natural Speech: some Implications for Sociolinguistic Methodology, *Language in Society* 5. 189–209.