# CORD
## Corpus Resource Database

http://www.helsinki.fi/varieng/CoRD

## General introduction

As the number of electronic corpora increases every year, it is becoming ever more challenging for linguists to be aware of, let alone familiar with the corpora available at any given time. The members of the Research Unit for Variation, Contacts, and Change in English (VARIENG) believe passionately that corpus linguistic research requires not only familiarity with the latest methodologies, but also with the details of each corpus. To that end, VARIENG has launched a free, open-access resource to facilitate the sharing of information on available linguistic corpora. In the first instance, we endeavor to collect information on all English language corpora.

Compiling a corpus takes a long time. It is also hard work. The last thing one wants to do is compile a specialist corpus only to discover that a colleague did the same thing a few years earlier and would have been happy to share it with you. As CoRD grows, chances of accidentally overlapping compilations are diminished.

### Authoritative information

All CoRD entries are written by the compilers of the respective corpus. No guess work, no interpretation. When you read a CoRD entry, you can trust the information to be accurate. Each CoRD entry comes with full details of when and by whom the information was submitted. A feedback form is made available for CoRD users to report any inaccuracies, and these can be communicated to the authors for clarification or explanation.

### Up-to-date and accurate

CoRD entries can be updated to reflect the latest news about ongoing corpus projects. CoRD can also be used to make lists of errata available to corpus users. Stratifications originally envisioned may change due to new scholarship, and newly discovered information may change the classification of a particular text. Through the pages of CoRD, compilers can communicate the latest data to corpus users. Naturally, a full record will be kept of the update history, and CoRD pages always indicate when and by whom the page was last updated.
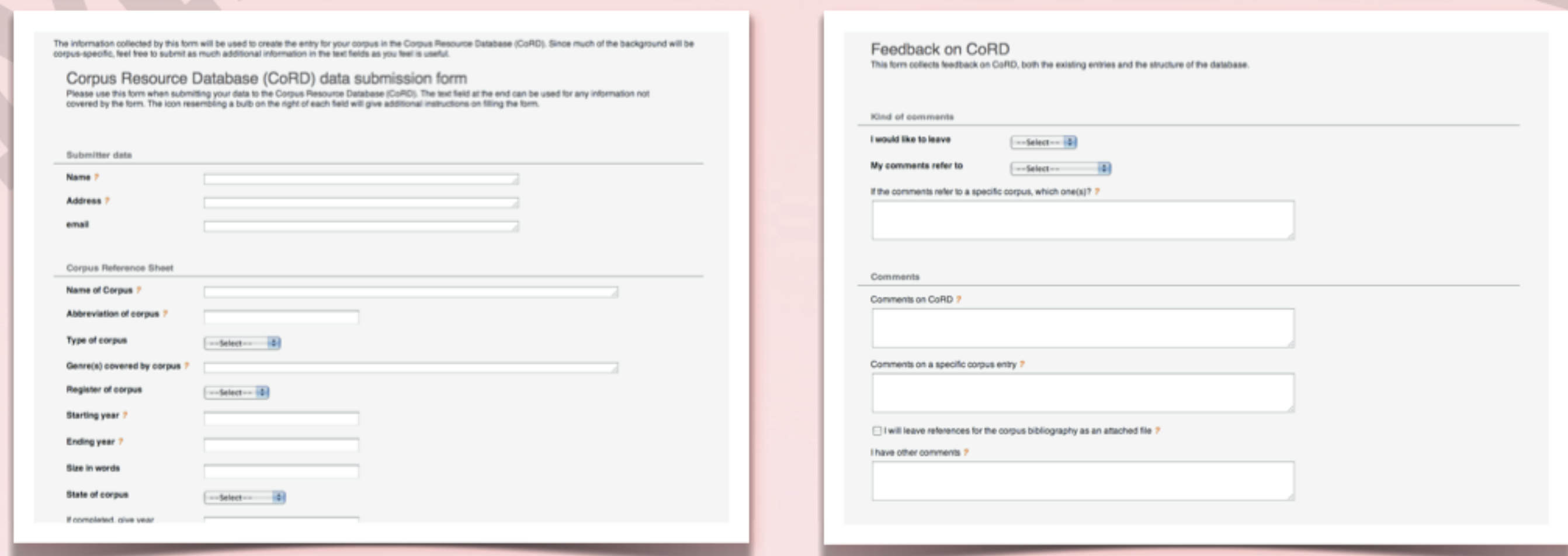


### Submitting a corpus to CoRD

Corpus information can be submitted in one of two ways. The simplest option is to use the electronic form, available on the CoRD website. We take the information and manually convert it into HTML following the guidelines of the CoRD template.



Alternatively, if your corpus information is more complicated, or if you have specific requests about the layout and presentation, you can get in touch with us and we'll be happy to work with you in designing the CoRD entry. Naturally, you can add to, alter, or remove any information at any time.

## What is CoRD for?

### Reference information

A common problem in corpus linguistic scholarship is deciding on the correct reference line when citing a particular corpus. CoRD provides the recommended reference line for each corpus, written by the compilers themselves.

### A quick source of information

Nearly all computers today are connected to the Internet. Through CoRD, users of corpora are one mouse click away from checking the basic data of a corpus such as its wordcount or that of its subcorpora, the annotation scheme, or perhaps the file format used in the distribution version.

### Giving credit where credit is due

One of the shortcomings of conventional citation formats is that often only the leaders of research projects are mentioned by name. All the others, many of whom may have contributed extensive amounts of time and labour, are relegated to a footnote. CoRD allows corpus compilers to acknowledge all team members. Because CoRD entries are free-form and open-ended, the compilers can use as much space as necessary to explain the roles and responsibilities of each team member.

Today, funding bodies are also increasingly interested in public acknowledgment of their contributions to research. Long-running projects, as well as those with many project members, may have benefited from funding by a number of different organizations. Through CoRD, corpus compilation projects can credit sources of funding in as much detail as necessary.

### Availability information

Some corpora are made available through open access academic repositories, while others are obtainable by subscription, membership in an organization, or directly from the compilers. CoRD tells you who to contact.

### Teaching

Teachers of corpus linguistics can use CoRD in a number of different ways. In addition to referring students to CoRD for basic information on corpora, CoRD entries can serve as introductory material to the process of corpus compilation. Narrative accounts, such as those provided by the compilers of the Helsinki Corpus and the Helsinki Corpus of British English Dialects, can help current and future corpus compilers understand both the challenges and joys of compiling corpora.

### Recording the living history of corpus linguistics

Although a relatively young discipline, corpus linguistics has in actual fact been around for quite some time. Over the years, corpus compilation methods have undergone considerable changes. CoRD is an online resource for recording and sharing not only hard facts, but also the stories and experiences of the compilers. Photographs, audio recordings and video clips can be used to liven up the presentation even further.

## The Structure of a CoRD entry

CoRD entries are essentially open-ended in structure. Because the information comes directly from the compilers, it is as accurate as possible. While we suggest a basic two-fold structure for most entries, the details specific to each corpus will ultimately dictate how the information is arranged and presented.

In principle, CoRD entries comprise two main sections: *Basic Information* and *Background Information*. The *front page* of every CoRD entry gives basic information on the timespan and size of the corpus as well as on the genres and text types covered. Other information on the front page includes the names of the compilers and other team members, and the recommended reference line. Additional information on funding and details related to academic affiliations can also be included. Information on availability comes next, followed by up-to-date contact information. Naturally, hyperlinks to the compilers' own website or to other resources can be included as well.

Information on the structure of the corpus also goes under *Basic Information*. Charts and tables can be used to illustrate the organization of the corpus. This sub-section details the principles of temporal, sociolinguistic and genre-related stratification systems employed. VARIENG researchers believe that it is extremely important not only to know what kind of stratifications are used, but also what principles lie behind them. The stratification can be explained in detail, up to and including bibliographical information of text sources.

Any annotation schemes used can be explained and detailed. This will allow corpus users to quickly consult the annotations used in one or more corpora, thus saving plenty of valuable time and effort otherwise used in chasing down such information. Importantly, the compilers have as much space as they like to truly explain the reasons behind the choices they made. Facsimile images of early material can be included to illustrate issues arising in the editing of manuscripts and early printed material, while audio files can be used when describing the transcription of spoken data.

The *Background* section of a CoRD entry affords the compiler(s) the opportunity to narrate the details of the compilation process. This can include information both on the progression of the compilation itself, and on the human side of working on corpora. Cooperation with specialist consultants, such as historians or literary scholars, can be discussed. Compilers are encouraged to submit any and all multimedia content relevant to the process of corpus compilation. Photographs of equipment used and scanned images of notebook and journal entries can be included as historiographical evidence of the compiling process.

Finally, CoRD is envisioned to provide an up-to-date *bibliography* of research conducted using a particular corpus. This would not only allow members of the corpus linguistic community to know of and potentially contact colleagues interested in similar topics, but also provide corpus compilers with a single point of reference when evaluating the fruits of their labour. Like the corpus information, all bibliographical information will be submitted by the authors of the articles.

## CoRD today

### What's on CoRD right now?

CoRD includes the detailed descriptions of five corpora. In addition to the *Helsinki Corpus*, CoRD entries covers three corpora compiled in Helsinki: the *Corpus of British English Dialects*, The *Corpus of Early English Correspondence* (CEEC), and the *Corpus of Early English Medical Writing* (CEEM), the latter two with detailed sections on their subcorpora. In addition, CoRD includes entries on the *SCOTS* and *NECTE* corpora.

### Multimedia on CoRD

Although still in an early stage, CoRD already includes a number of innovative features. In addition to standard visual aids like charts and maps, several entries feature photographs of the compilers and of the equipment used during the 1970's, 80's, and 90's. The entry on the *British English Dialects* corpus features audio clips to illustrate the recordings made in the 1970's and 80's on which the transcriptions are based. Entries on all corpora compiled in Helsinki come with a video presentation by the compilers: Matti Rissanen on *Helsinki Corpus*, Terttu Nevalainen on *CEEC*, Irma Taavitsainen on *CEEM*, and Kirsti Peitsara on the *British English Dialects corpus*.

### Example: pages of The *Helsinki Corpus of British English Dialects* on CoRD



CoRD entry for the *Helsinki Corpus of British English Dialects* was designed by Anna-Liisa Vasko and Simo Ahava.

Copyright for original photography and audio samples held by project members.

## VARIENG

The Research Unit for Variation, Contacts, and Change in English (VARIENG) is a Center of Excellence funded by the Academy of Finland. VARIENG functions at the Universities of Helsinki and Jyväskylä, in affiliation with the Departments of English and Modern Languages, respectively.

### CoRD team

The CoRD concept was envisioned by Terttu Nevalainen, Jukka Tyrkkö and Minna Palander-Collin. The basic structure for CoRD entries was developed by Jukka Tyrkkö, Matti Rissanen, Terttu Nevalainen, Irma Taavitsainen and Arja Nurmi, who also designed the electronic submission and feedback forms. The template for the CoRD website was designed by Tanja Säily and Jukka Tyrkkö, the present CoRD coordinator.

### Contact information

To contact VARIENG, see http://www.helsinki.fi/varieng/