# Qualitative data analysis software: The state of the art

*Susanne Friese**

The aim of this paper is to provide an overview of the 'state of the art' of QDA or CAQDAS software. As the range and number of packages have increased over the years, I needed to make a decision about how I wanted to approach it. As the recent hype on 'big' and social media data has also left its mark on the type of functionality we see emerge in CAQDAS packages, and with this on the need for different types of analysis tools, I decided to use Kahneman's ideas about slow and fast thinking as a framework (Kahneman, 2011). Slow thinking in the context of CAQDAS is related to researcher driven analysis and fast thinking to tool and data driven analysis. The paper is divided into two parts. In the first part I describe trends and new developments and in the second part I offer a critical appraisal. I assume that readers are familiar with the basic functionalities of at least one of the CAQDAS packages.

## Qualitative Data Analysis Software or CAQDAS?

In most of my writings, I point out that I prefer the acronym CAQDAS (computer-aided qualitative data analysis software) as compared to QDA software as it more accurately describes what these programs are doing: aiding analysis rather than analyzing qualitative data. Looking at the new developments and the directions some CAQDAS packages are pursuing, I may need to revise my yearlong persistence to use the term CAQDAS rather than QDA, at least for some packages. The companies developing the programs react to the type of data that has become available in recent years and that is of potential interest for the users of their software. If we go back a couple of years, this started with the widespread availability of PDF documents. One program after the other offered support for this format. In recent years, the option of importing data from reference managers for the purpose of conducting a literature review or meta-analysis was added. Another recent new feature is the import of survey data for the analysis of open-ended questions, because web-based surveys have become more common. In the latest versions of a number of packages, social media data have entered the scene. NVivo has built a tool to capture websites, Facebook, Twitter, and YouTube data. MAXQDA 12 offers Twitter import, ATLAS.ti 8 Twitter and Facebook import. Along with those tools automated analysis options are implemented. Data are automatically coded on the basis of hashtags, and information about retweets, 'mentions', and geolocation are used to generate social network maps. The Plus version of NVivo 11 adds another

* Dr. Susanne Friese is a senior research partner at Max Planck Institute for the Study of Religious and Ethnic Diversity, Göttingen, Germany. Email: Friese@mmg.mpg.de.

automated feature by promising automated insights through thematic and sentiment coding. Below I will discuss and evaluate these features in more detail. The inclusion of such features shows a potential shift from CAQDAS to QDA software, from researcher driven qualitative data analysis to software driven analysis of qualitative data. As Van Dijck (this issue) points out in her introductory article, both will be needed, also in the future. Big data analysis will not displace the necessity for close reading of data, but we cannot use the same tools for both tasks. Are the new tools that CAQDAS offers responding to this need? And if so, are they fit for the task of handling big data? From a user point of view, is it desirable to have software that incorporates the full array of tools ranging from supporting interpretive approaches to gaining quick insights from big data? Or would a division of labor make more sense in terms of using specialized programs for close reading of data on the one hand and distant reading on the other? In order to find answers to these questions, I found Bosch's reference to Kahneman's book on fast and slow thinking very appealing as it provides a suitable framework for analyzing what is happening in terms of software development for CAQDAS (Bosch, this issue).

*Fast and slow thinking*
Kahneman (2011: 23) describes these two modes of thinking (fast and slow) as 'System 1' and 'System 2':

> "System 1 runs automatically and System 2 is normally in a comfortable low-effort mode ... System 1 continuously generates suggestions for System 2: impressions, intuitions, intentions, and feelings. ... If all goes smoothly, which is most of the time, System 2 adopts the suggestions of System 1 with little or no modifications. ... When System 1 runs into difficulty, it calls on System 2 to support more detailed and specific processing that may solve the problem of the moment."

As more and more data become available, it is impossible for System 2 to analyze them all. Therefore, new tools are developed that appeal more to System 1's way of thinking, fast and intuitive. Based on many experimental studies, Kahneman found out that System 1 decisions, judgments, and evaluations are in many situations good enough, but are easily influenced by illusions and biases. Even though these flaws exist, we need System 1 type thinking, because System 2 is lazy, needs more incentives to begin to work, takes more effort and energy, and is much slower. In facing big data, we probably won't get far relying on System 2 and we need more tools that assist System 1, or tools that prepare data in a way that allows System 2 to make sense of it. In the following, I will review the trends and new developments of CAQDAS tools in light of their support for System 1 and System 2 ways of thinking and discuss the pros and cons.

## Trends and new developments

### Types of data
In recent years, CAQDAS has supported many more data types and file formats. Due to space constraints, they cannot all be discussed in detail here. In the next section, however, I touch upon social media data.

### Analytic tools
In terms of analytic tools, we see an increase in numbers and options, and the tools have become more sophisticated. While the tools available in the initial years of CAQDAS were typically addressing System 2's way of thinking, most of the new tools appeal to System 1, but not all. Examples of System 2 tools are simple and complex code queries using Boolean, semantic, and proximity operators. The software user needs to read through all retrieved segments, needs to interpret what she or he has read, and needs to begin to write up the interpretation. This takes time and some effort. Formulating the appropriate queries requires knowledge of the operators and how to use them.

On the other end of the spectrum, more and more packages add System 1 tools that allow the creation of quick overviews in various formats like tables, word clouds, word trees, graphical displays of code frequencies, and cluster analysis for document, attribute, or code similarity. This is complemented by automatic coding features whose purpose is to identify themes in the data or to analyze documents for sentiment. More sophisticated tools allow users to evaluate the automatically generated results in terms of whether it is a good or bad fit. Based on the decision made by the researcher, the algorithm learns, and the output improves over time. A modification of this is automatic pattern coding. This is a feature that codes data based on initial manual coding. How well this already works is discussed below in the section 'Critical appraisal'.

Looking at a longer time spam, what can also be observed is an increase in quantification of qualitative data. My assumption is that this is a user-driven demand and, as it is relatively easy to make a computer program to calculate something and to output numbers, software companies have responded to it. Examples are: word frequency counts, number of co-occurring codes, a coefficient indicating the strength of the relation between two codes, centrality measures for sociograms, document similarity measures, measures for intercoder agreement, tables showing number of codes (or number of words/characters of coded segments) by document or document groups. Often absolute and relative frequencies are offered, or an aggregated count per document/document group.

### Mapping tools
Over the years, more and more software packages have added tools to model or map data, if these were not already integrated from the very beginning - as in ATLAS.ti. In recent years, these tools have been enriched with analytic functionality, mostly with features and tools that appeal more to System 1. Instead of having to create a model manually and thinking through what should be linked to which item, why this should be done, and how items are related, predefined solutions are

offered like the one-case, one-code, or code-occurrence model, project maps, or comparison diagrams.

## Critical appraisal

### Approaching analysis

After adding data to a program, analysis can be approached in various ways. The classical way of qualitative data analysis is to start reading the data, beginning with coding in the second or third round of reading. This is feasible for instance with interview studies that contain about 20 to 40 transcripts. The larger a data set becomes, the less likely it is that we will be able to inspect the data in such a close manner. Consequently, a shift must occur from qualitative data analysis to analyzing qualitative data. This in turn requires different methodological approaches as well as different software tools.

We are currently observing the trend that the range of supported data file formats in CAQDAS and also the analytical toolbox are being extended. Some of these tools are geared towards 'big data' analysis and appeal more to System 1, at least at first sight, whereas other tools aim at better assisting close reading of the data, thereby addressing System 2. I start with discussing tools that appeal more to our intuition, tools that offer automated data processing and promise to provide quick insights; tools that thus are more appealing to System 1. The outcomes of such an analysis may spark the interest of System 2 to come in to either correct first intuitive thoughts or to go deeper and offer more detail and precision.

## Approaching analysis with System 1

Word clouds and text searches combined with auto coding, offered by most programs, are examples of System 1 tools. These are useful, for instance, when analyzing fairly structured data such as open-ended questions from a survey. If I have such a data set, I start the analysis by creating a word cloud to see which words occur and, based on what I find interesting, I auto code. Once this is done, System 2 needs to kick in as the auto-coded segments need to be checked, possibly adjusted, recoded, and resorted. New 'auto code by themes' features can take over the first step of this process by looking up words that occur in a text and making decisions which words to use for auto coding. I tested this new feature with a project that I had previously coded in the manner described above, starting with a look at word frequencies followed up by auto coding. When I left this first step up to the software, it sorted the data into thirteen main themes of which six were the same as I had previously selected when coding the data without the tool. This similarity diminished when comparing the automatically generated themes to the categories I developed after reviewing the auto-coded data. This is to be expected, as I consider the initial themes only as starting point and not as an end result. I probably would have developed a similar final coding schema based on the suggestions provided by the computer. However, I consider taking auto-coded data at face value for

decision-making purposes to be rather bold. There is a trade-off between generating quick results and accuracy. This resembles quite closely Kahneman's (2011) description of System 1 and System 2 modes of arriving at conclusions and making decisions. Intuitive hunches are often correct, at least more accurate than chance. And the more skilled a person is in an area, the better the intuition will be. Nevertheless, System 1 decisions are prone to errors as ambiguities are neglected and doubts are suppressed, and there is a bias to confirm beliefs, a focus on existing evidence, and an ignorance of absent evidence. What counts is the coherence of the story, regardless of the amount of data or the quality of the data it is based upon. Thus, if a person can make up a good story based on what the software delivers either by automatically coding the data, or in form of word clouds, word trees and charts, System 2 has no reason to take actions. In the end it may not even notice that System 1 has indeed provided an answer, but possibly an answer to the wrong question. Hence, when approaching data analysis via a text search followed up by auto coding, suggestions made by the computer can replace this first manual step. This however needs to be followed up by further close reading of the data, or we need more sophisticated tools than those currently available.

The social media hype is also leaving its marks on CAQDAS development. Even if one might become excited by playing around with the new possibilities, one soon reaches the limits of what is possible. CAQDAS has not been built to handle thousands of documents, and if you attempt to work with larger data sets, my experience, especially with NVivo, was that everything takes very long and crashes are frequent. This was in part also because I was attempting analysis in a System 2 way. When auto coding all Twitter hashtags in my sample project, I ended up with a list of over 2,300 codes. This list contained lots of similar code labels as the hashtags were not all spelled in the same way like Refugeesnotwelcome/Refugeeswelcome/ Refugeeswellcome, etc. To build a proper coding system – in a System 2 manner – would require lots of cleaning up, and one would want a bit more help from the software. A solution would be a somewhat smarter auto-coding feature that would recognize similar hashtags by stem. The performance is better if one runs a software-driven analysis applying different tools as suggested by the manual. The proposed approach is to sort, reorder and filter the data, display charts, compare numbers of followers and following, and create a cluster analysis or a computer-generated Twitter sociogram. These are all procedures more appealing to System 1. When testing the analysis of email messages, a similar conclusion can be drawn. The software was not able to handle larger amounts of data (1,000 email messages). A quick System 1 analysis – creating a network sociogram – worked well. When applying automated theme recognition, though, the software delivered useless results proposing themes like "a", "and", "is", "have" and the like.

Similarly disappointing was the experience with sentiment coding. The purpose of this feature is to discover positive and negative sentiments in a text. The results at the moment are not at all reliable. The following sentences for instance were classified as very negative:

- *I am not always happier than before I had children, but my life is richer for having had them.*

- *Minecraft is a great "puzzle" option (compared to the hunt and kill themes that are marketed towards kids).*

The following sentences were classified as very positive:

- *They look exhausted and frustrated most of the time, and are most willing to get away from their kids any chance they get.*
- *Minecraft is highly addictive for tweens and that is probably the biggest danger – tearing them away from this largely benign and creative game!*

Not all data is classified wrongly; overall, the program seems to be better in picking up positive than negative sentiment and the results are better for semi-structured than for narrative interviews. The algorithm for sentiment coding might improve over time and deliver more reliable results. Building in machine-learning capabilities may also help to improve results. The limiting factor might still remain the amount of data that can (or cannot) be handled. CAQDAS currently does not handle large amounts of data easily and as quickly as text-mining programs. Looking at the current state of the art, I lean more towards using specialized tools when wanting to analyze larger data sets instead of CAQDAS, combined with offering export and import options that facilitate data exchange and reusability of analyzes that have already been done elsewhere. Such an approach is also favored by Stulpe and Lemke (2016). They call it *blended reading*. As it is obvious that there are limits to what a human interpreter can handle and manage even with the support of software, it is also too optimistic to expect Digital Humanists to develop procedures that allow researchers to analyze data per mouse click in a flash producing publishable results. Data-mining tools can count, structure, and visualize. The more interesting question, however, is how these text statistics can be interpreted. And for this we need the recourse of the toolkit of the qualitative researcher. Distant reading is only one side of the coin and, given the available data material, a necessary one. As this is new, there is an understandable hype around it but, as Hitchcock formulates it, this does not mean that other analytical procedures are no longer relevant:

> "I end up ... feeling that in the rush to new tools and 'Big Data' Humanist scholars are forgetting what they spent much of the second half of the twentieth century discovering – the language and art, cultural constructions, human experience, and representations are hugely complex – but can be made to yield remarkable insights through close analysis. In other words, while the Humanists and 'Big Data' absolutely need to have a conversation, the subject of that conversation needs to change, and to encompass close reading and small data." (Hitchcock, 2014)

## Approaching analysis with System 2

The concept of blended reading can also be applied when looking at tools within CAQDAS. As shown above, most tools appealing to System 1 need to be followed up by System 2 with a more detailed analysis, if we choose to start the analysis in this way. Data analysis can, however, also take its point of departure with System 2. This is the case for all interpretive and inductive approaches. The available tools for this type of analysis also vary across software packages. A few packages like AT-LAS.ti, Dedoose, and Transana, offer a functionality that allows an analysis below the coding level. This means a highlighted data segment can become an entity of its own and does not need to be coded. Involving System 2, one can begin analysis by reading through the data, marking interesting data segments and creating quotations, excerpts, or quotes. The comment field can be used to formulate first interpretations (see Figure 1).

MAXQDA has just introduced 'tabled documents' that fulfil a similar function, as comments and interpretations can be added in columns and cells next to the data. We also see new features being developed that assist a System 2 analysis based on coded data. An example is the summary grid in MAXQDA (Figure 2). It assists users in writing up the analysis. All of these are features that need to be brought to life by the researcher. They support qualitative data analysis rather than offering to analyze the data. System 1 also has a role in this process, for instance, when writing down word associations, first intuitive thoughts about probable meanings and connections.

## Visualization

With regard to modelling tools, programs can be sorted into three types: those that offer tools that are designed to complete a particular task in a pre-defined, software-determined way. This makes them easy to use, but the drawback is that one needs to stay within the boundaries of what is predefined (Figure 3). Most of these tools serve a System 1 analysis.

Then there are programs that offer a number of predefined options, but allow expansion of the software generated maps, and the option of building models from scratch (Figure 4). Thus, a combination of System 1 and System 2 type of analysis is possible.

A third type consists of packages that had a built-in modelling function from the very beginning and thus offer the most integrated design (Figure 5). Visualization is not realized as an 'add-on', rather everything in a project can be visualized, and vice versa, everything one does in a model is reflected in the project. This integrated design has a number of advantages as there are no boundaries between the various mapping applications. This is coupled with tools that are more appealing to System 1, such as an automated import of co-occurring codes or neighboring objects facilitating a fluid change-over between System 1 intuition and more detailed follow up examination by System 2.
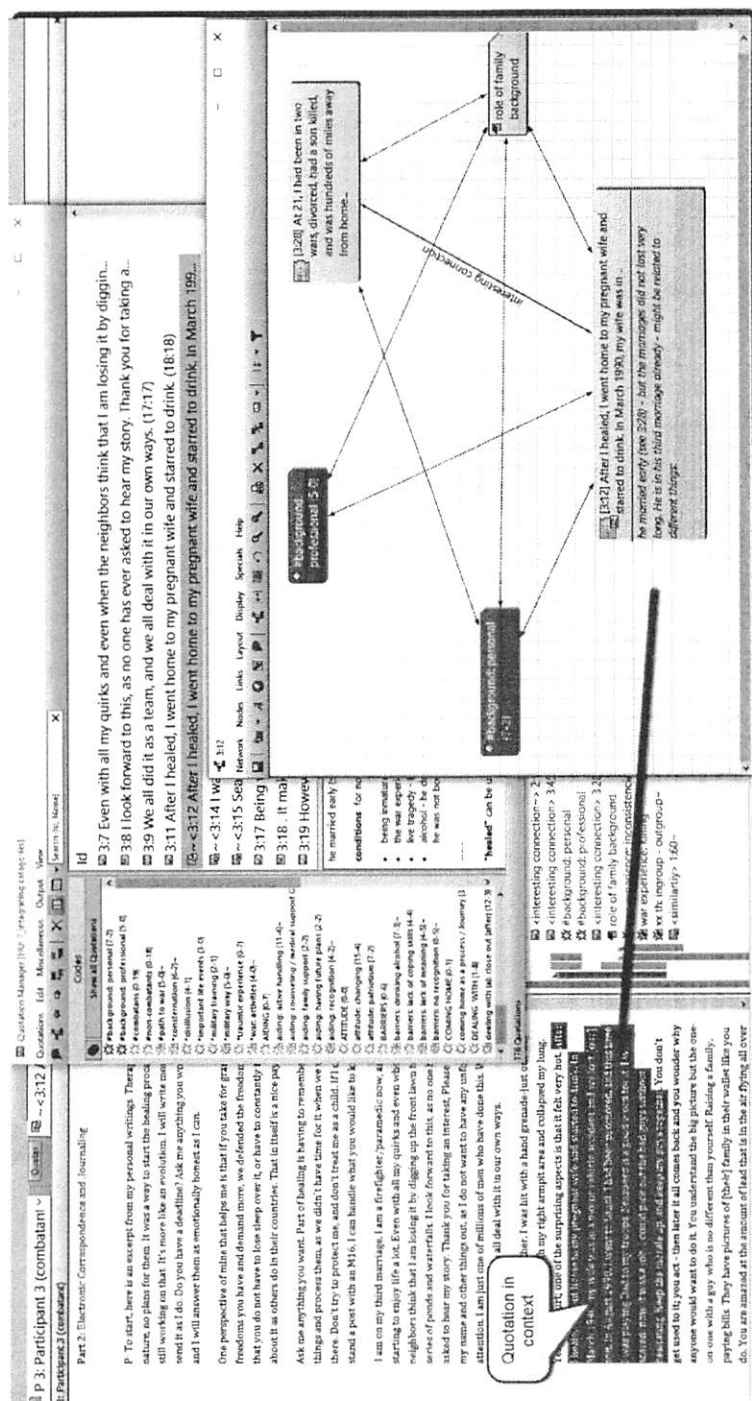
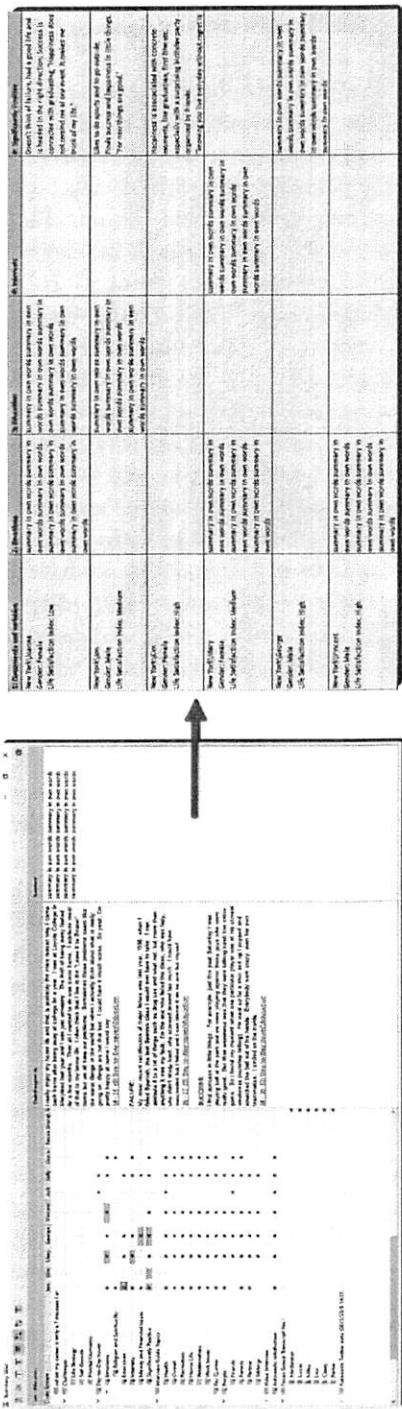*Figure 1    Working below the code level in ATLAS.ti*

Susanne Friese



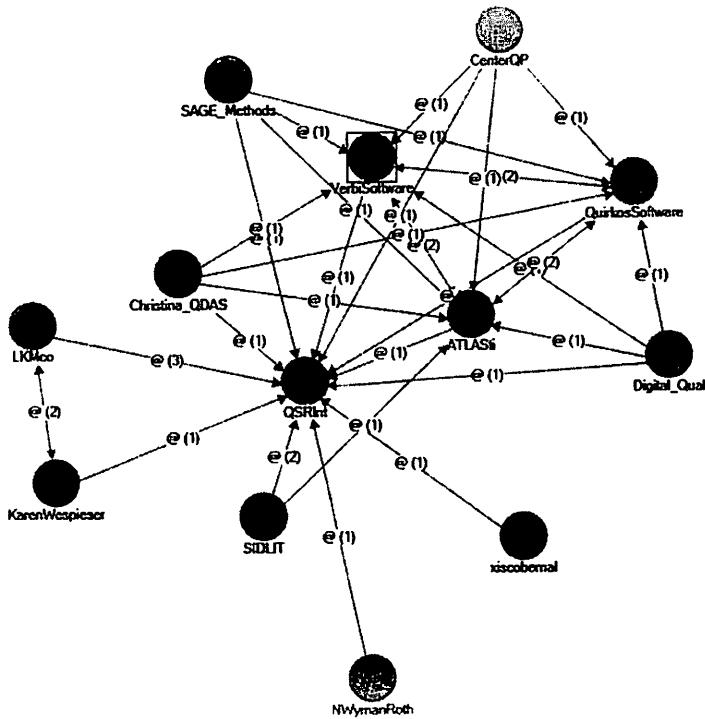*Figure 2     Summary Grid in MAXQDA*

*Figure 3*    *Predefined modeling tools in NVivo*
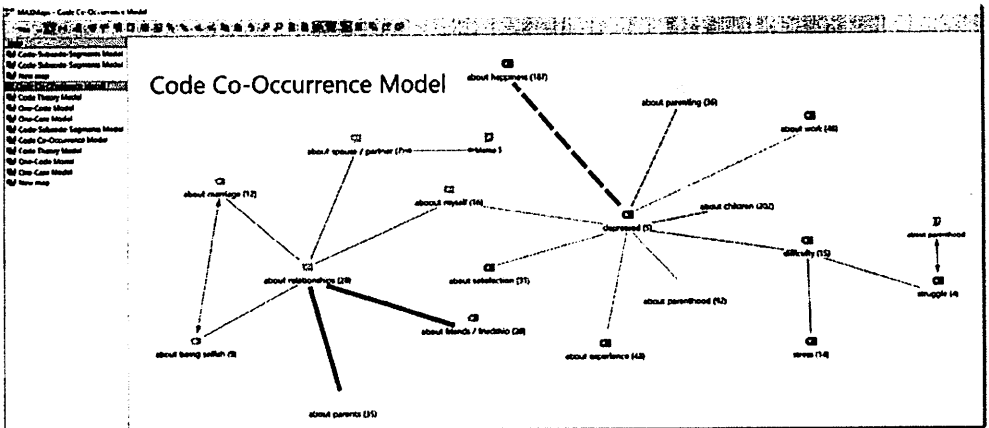


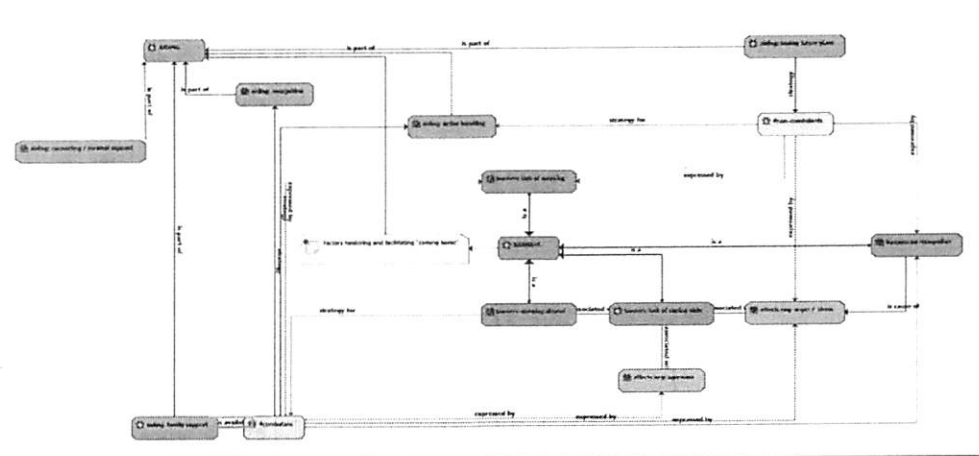*Figure 4*    *MAXQDA predefined map plus some manually added elements*

**Figure 5** *Integrated and flexible visualization in ATLAS.ti 8 for Windows*

## Summary

Over the past few months I have been 'playing around' with various packages and worked with proven old features and new ones. The proven old features have become much more sophisticated over time and new functionality has been added to existing tools. There is a trend to integrate social media data, with a strong emphasis on automation and tool-driven analysis, thus tackling the Big Data scene. Even if this sounds exciting at first sight, upon closer look it turns out that CAQDAS is not yet living up to expectations. The implemented algorithms appear to be quite simple, and most programs do not handle larger amounts of data well. The promise of being able to analyze social media data, like Twitter, Facebook, and emails with the combination of automated tools, raises expectations that software can also easily handle a few thousand posts or emails. The current state of the art is that more new data formats are supported and more automated analysis tools are added, but the development is still in its infancy.

It is plausible that developers of CAQDAS packages also want to jump on the bandwagon and join the Big Data crowd. The question that needs to be raised is whether it is sensible to incorporate specialized content analysis and text mining tools into CAQDAS, and if so, to what degree. Can the same type of sophistication be reached that is already available in text mining or even Big Data software? And if so, does it make sense to use these tools with smaller amounts of data? Or would good advice be: "cobbler stick to your last"? To end this paper, I would like to rephrase Hitchock (2014): CAQDAS and 'Big Data' absolutely need to have a conversation. The subject of that conversation should be how (or whether) to integrate *close reading and small data*, and *distant reading and large data*.

# References

Bosch, R. (this issue). Editorial: Qualitative research in the digital humanities. *KWALON 61*, 21(1).

Friese, S. (2016). Grounded Theory computergestützt und umgesetzt mit ATLAS.ti. In C. Equit & C. Hohage, *Handbuch Grounded Theory – Von der Methodologie zur Forschungspraxis* (pp. 483-507). Weinheim: Beltz Juventa.

Hitchcock, T. (2014). Big data, small data and meaning. *Historyonics*, http://historyonics. blogspot.sg./2014/11/big-data-small-data-and-meaning_9.html, Consulted: 2015-12-22.

Kahneman, D. (2011). *Thinking, fast and slow*. London: Penguin Books.

Stulpe, A. & Lemke, M. (2016). Blended Reading. Theoretische und praktische Dimensionen der Analyse von Text und sozialer Wirklichkeit im Zeitalter der Digitalisierung. In M. Lemke & G. Wiedemann (Eds.), *Text Mining in den Sozialwissenschaften: Grundlagen und Anwendungen zwischen qualitativer und quantitativer Diskursanalyse* (pp. 17-62). Wiesbaden: Springer VS.

Van Dijck, J. (this issue). Big data, grand challenges: On digitization and humanities research. *KWALON 61*, 21(1).

Wiedemann, G. & Lemke, M. (2016). Text Mining für die Analyse qualitativer Daten. Auf dem Weg zu einer Best Practice? In M. Lemke & G. Wiedemann (Eds.), *Text Mining in den Sozialwissenschaften: Grundlagen und Anwendungen zwischen qualitativer und quantitativer Diskursanalyse* (pp. 397-420). Wiesbaden: Springer VS.