



- Parallel Workshop 2-

The DIPEX database and NLP approaches to analysis (English)

Workshop Leaders: Gerold Schneider, Giovanni Spitale, Tilia Ellendorff

gschneid@ifi.uzh.ch

giovanni.spitale@ibme.uzh.ch

tilia.ellendorff@cl.uzh.ch



Content WS 2

1. Introduction Input
2. Leading questions for discussion
3. Output of discussion



1. Input

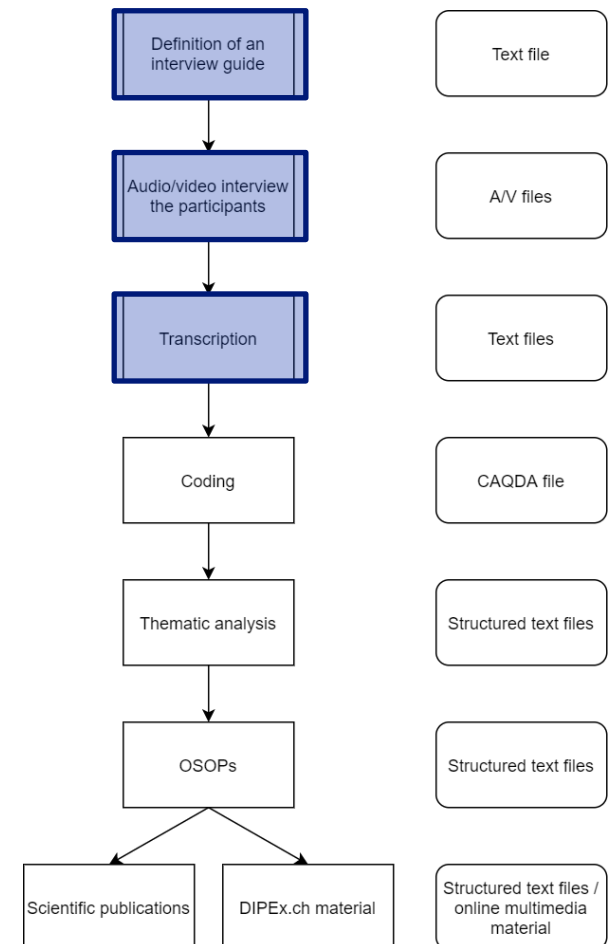
1.1 The Dipex Database

Interview guide

Semi-structured document listing the questions to ask and the prompts to give to the interviewee. Starts with an open section then follows specific topics of interest.

Interview files

The interview is audio and/or video recorded – according to the preferences of the interviewee – and then transcribed as text.



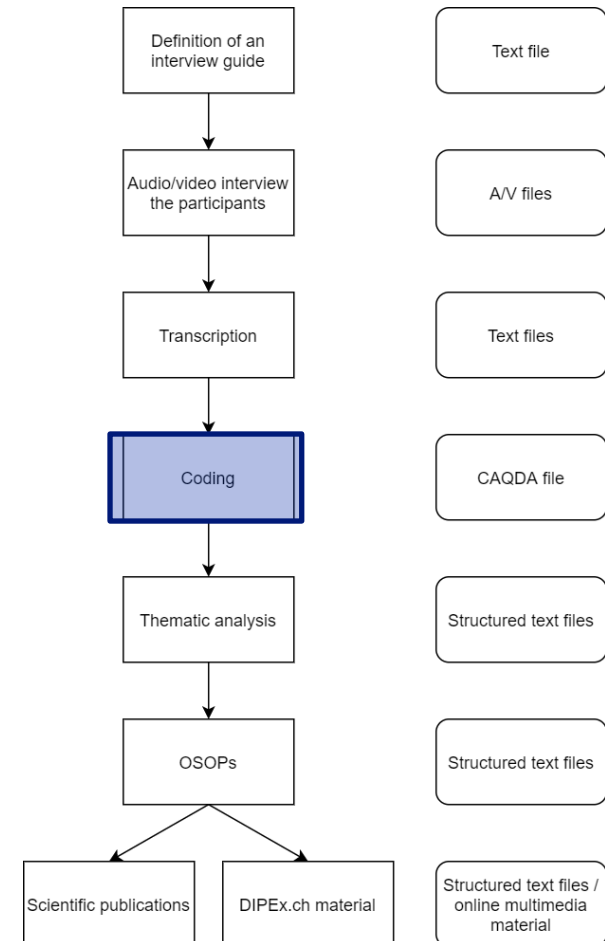


1. Input

1.1 The Dipex Database

Coding

The interviews are loaded in a software for computer assisted qualitative data analysis. We define a specific coding tree and manually code the text (= assign one or more labels to a meaningful passage)



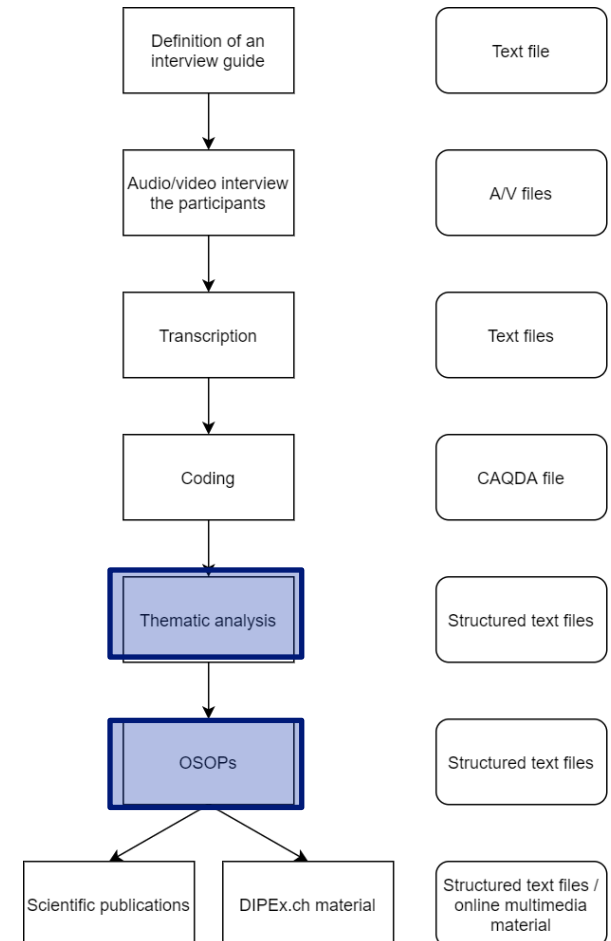


1. Input

1.1 The Dipex Database

Thematic analysis and OSOPs

We select specific topics that allow us to tell the ‘collective history’ of a given experience weaving individual voices together. We attribute codes to topics, retrieve the quotes, and put the story together.



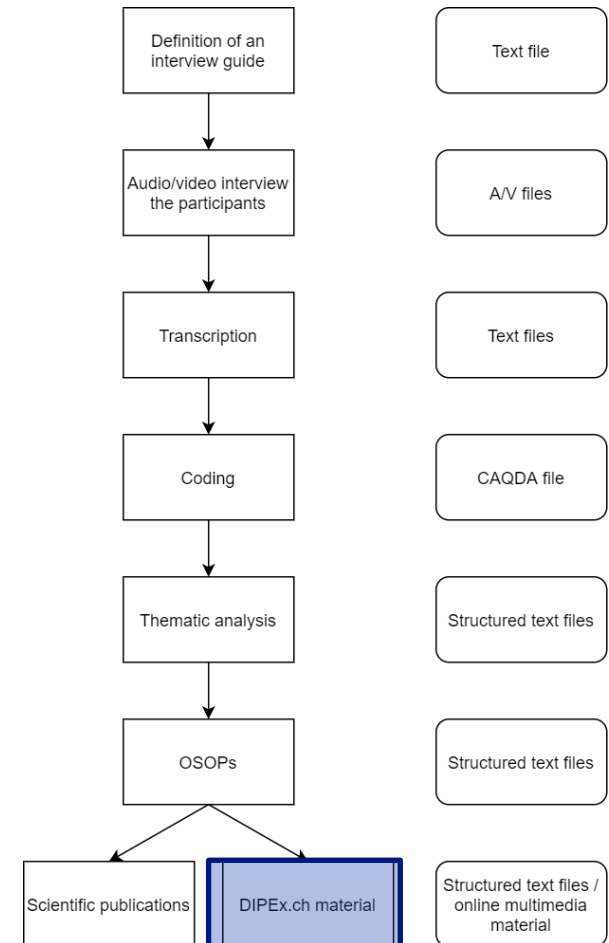


1. Input

1.1 The Dipex Database

Output

Based on the OSOPs, we prepare selected material to be put online and serve as an important resource for patients, relatives, caregivers, healthcare professionals, and students.





1. Input

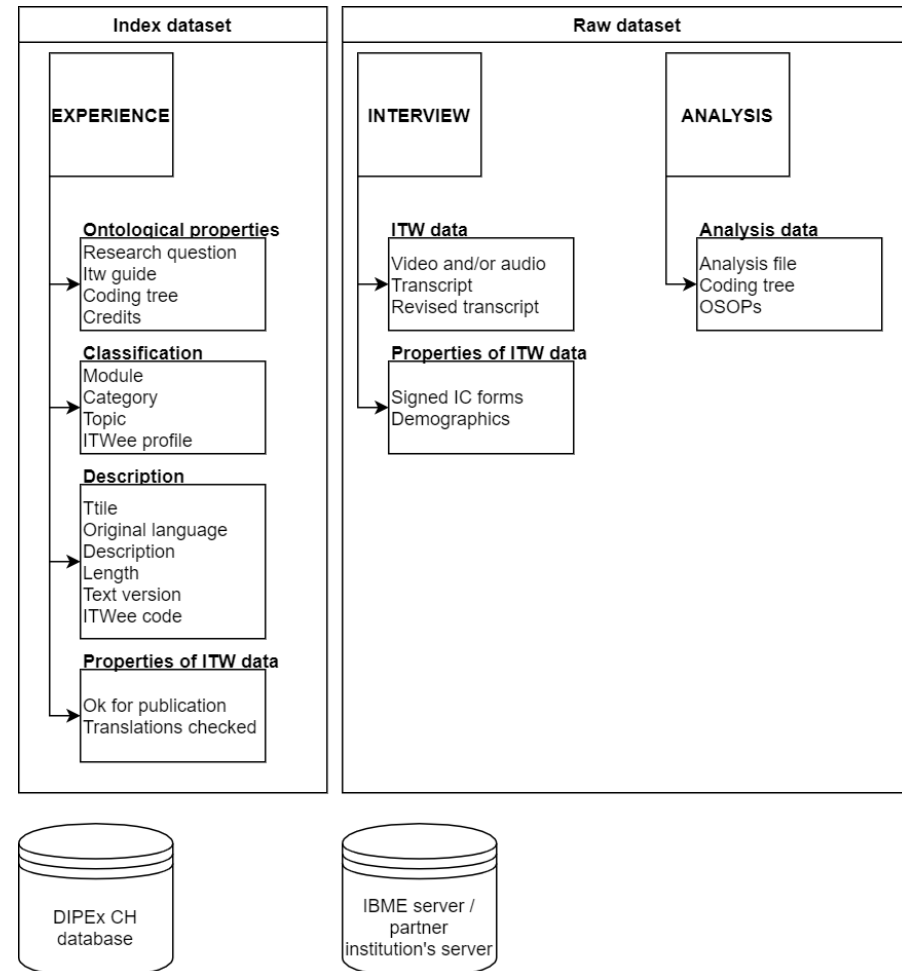
1.1 The Dipex Database

Raw dataset

The DNA of our research data; located on IBME's servers

Index dataset

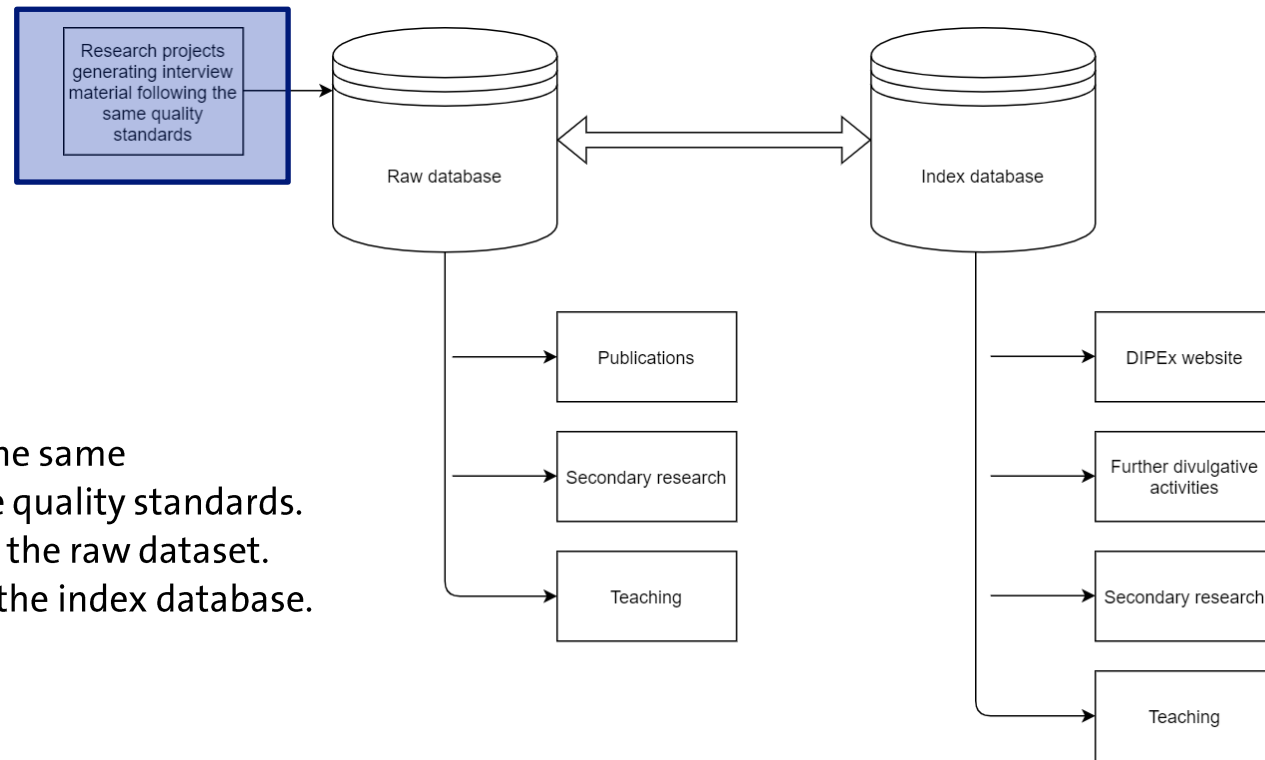
The index and mRNA of our research data (extended metadata); located on UZH's MariaDB instance





1. Input

1.1 The Dipex Database



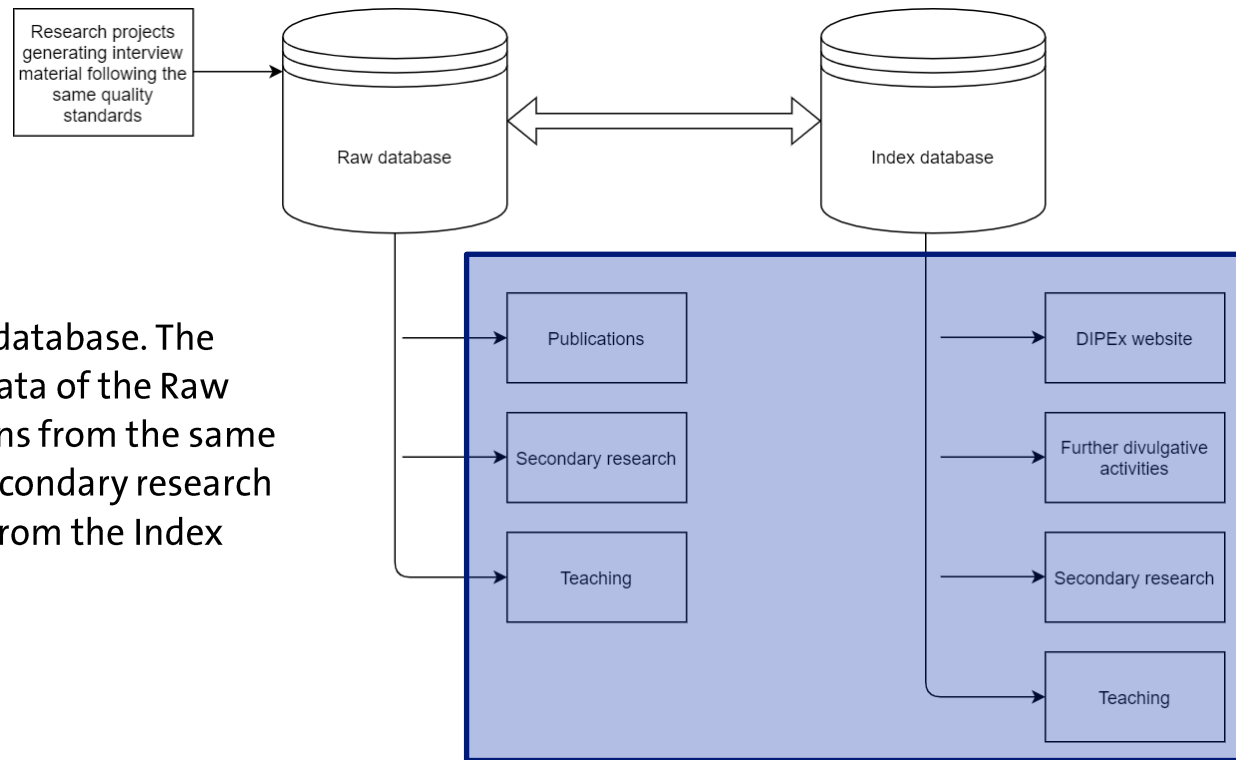
Data in:

Interview data generated with the same methodology and with the same quality standards. Raw interview data are stored in the raw dataset. Their content is then indexed in the index database.



1. Input

1.1 The Dipex Database



Data out:

Data can be explored via the Index database. The Index database contains the metadata of the Raw database, but also organized sections from the same data. Depending on the purpose, secondary research or teaching activities can use data from the Index dataset or from the Raw dataset.



1. Input

1.1 The Dipex Database

Filter experiences by topic

This creates a slice of the dataframe containing only the experiences belonging to a specific topic

```
[23]: experience_df_subset[experience_df_subset[EXP_topic] == 'Delusion/hallucination/ dreams/anxiety']
```

[23]:	ID_Experience	Belongs_to_module	ITW_code	ITWee_code	ITW_original_language	EXP_start_time	EXP_end_time	EXP_duration
21	31	CMI	CMI01_22102019_FR	CMI01	FRE	0 days 00:32:02	0 days 00:32:45	0 days 00:00:43
22	32	CMI	CMI03_05122019_EN	CMI03	ENG	0 days 00:05:35	0 days 00:07:20	0 days 00:01:45
23	33	CMI	CMI08_16122019_DE	CMI08	DEU	0 days 00:06:06	0 days 00:07:06	0 days 00:01:00
24	34	CMI	CMI05_05122019_FR	CMI05	FRE	0 days 00:27:48	0 days 00:29:16	0 days 00:01:28
25	35	CMI	CMI10_15012020_DE	CMI10	DEU	0 days 00:06:43	0 days 00:08:20	0 days 00:01:37
26	36	CMI	CMI20_29022020_DE	CMI20	DEU	NaT	NaT	NaT
27	37	CMI	CMI23_03072020_FR	CMI23	FRE	0 days 00:20:30	0 days 00:21:18	0 days 00:00:48

Getting all the experiences
belonging to a topic
...in 1 line of code



1.2 NLP Approaches to Analysis: Methods

How can we explore & analyse DipEx texts automatically?

Raw texts or annotated texts? → Supervised vs. Unsupervised

There are many useful methods. We just present three of them.

1. Document Classification (supervised)
2. Distributional Semantics (unsupervised)
3. Conceptual Maps (unsupervised or semi-supervised)

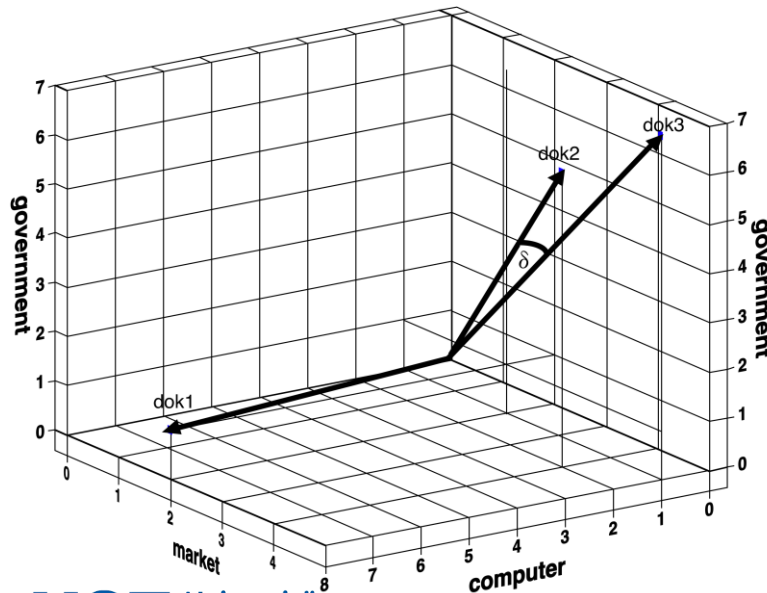
These approaches share that they need numerical text representations, typically vector spaces



1.2 NLP Approaches to Analysis: Vector Spaces

Document Classification uses
Document-Term Matrices

Frequency	<i>market</i>	<i>computer</i>	<i>government</i>
dok1	2	8	1
dok2	4	2	6
dok3	5	1	7



- Context Windows e.g. $[-10 \ w_0 \ +10]$ give us Term-Term Matrices

	<i>dog</i>	<i>hyena</i>	<i>cat</i>
runs	1	1	4
barks	5	2	0

TABLE 1 Distributional vectors representing the words *dog*, *hyena* and *cat*.

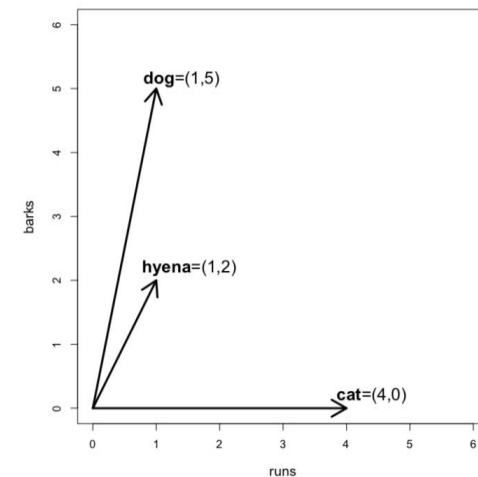


FIGURE 1 Geometric representation of the vectors in Table 1.



1.2 NLP Approaches to Analysis:

Teaser 2: Distributional Semantics with DIPEX MS

```
> closest_to(training_ms20,"leben")  
word similarity to "leben"
```

1	leben	1.0000000
2	gelassenheit	0.5548906
3	beziehungen	0.5310467
4	ausmacht	0.5218526
5	negative	0.5194965
6	krankheit	0.5168033
7	dingen	0.5158901
8	lebenswert	0.5157932
9	lebens	0.5121614
10	positiven	0.5043837

mein Kommentar: die Ebene an Reife ist beeindruckend: Gelassenheit, leben mit der Krankheit. Sich konzentrieren auf das was das Leben ausmacht: Beziehungen zu Menschen, die positiven Dinge als lebenswert annehmen.

```
> closest_to(training_ms20,"koerper")  
word similarity to "koerper"
```

1	koerper	1.0000000
2	physik	0.5999420
3	koerpers	0.5681156
4	wahr	0.5043965
5	sondern	0.4641749
6	bewusstsein	0.4570825
7	staerken	0.4353815
8	power	0.4327700
9	benoetigt	0.4313358
10	intensiver	0.4302129

mein Kommentar: gegen die Physik kommt niemand an. Stärker wahrnehmen, intensiver erleben, das Bewusstsein stärken.



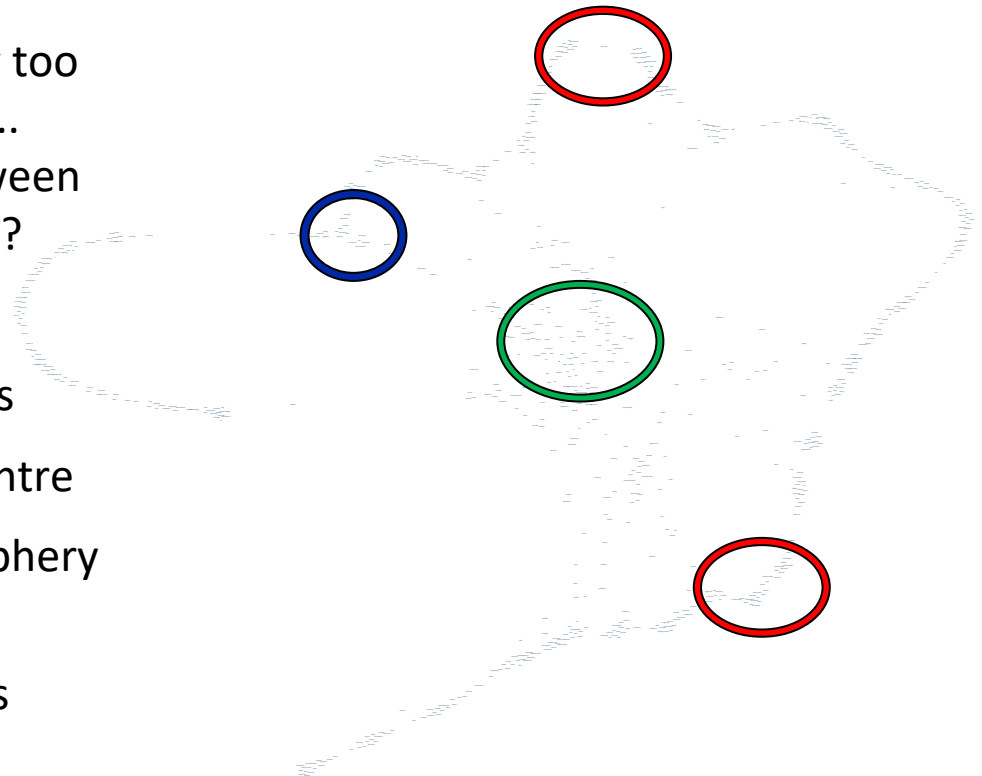
1.2 NLP Approaches to Analysis:

Teaser 3: Conceptual Maps with DIPEX CMI

DIPEX CMI is only 8000 words, way too small for distributional semantics ...
but can we explore the space between words and concepts, plot the plots?

The conceptual map (500 w) shows

- Common core (**green**) in the centre
- Shared issues (**blue**) in the periphery
- Individual adventures (coma, hallucinations, **red**) in the edges



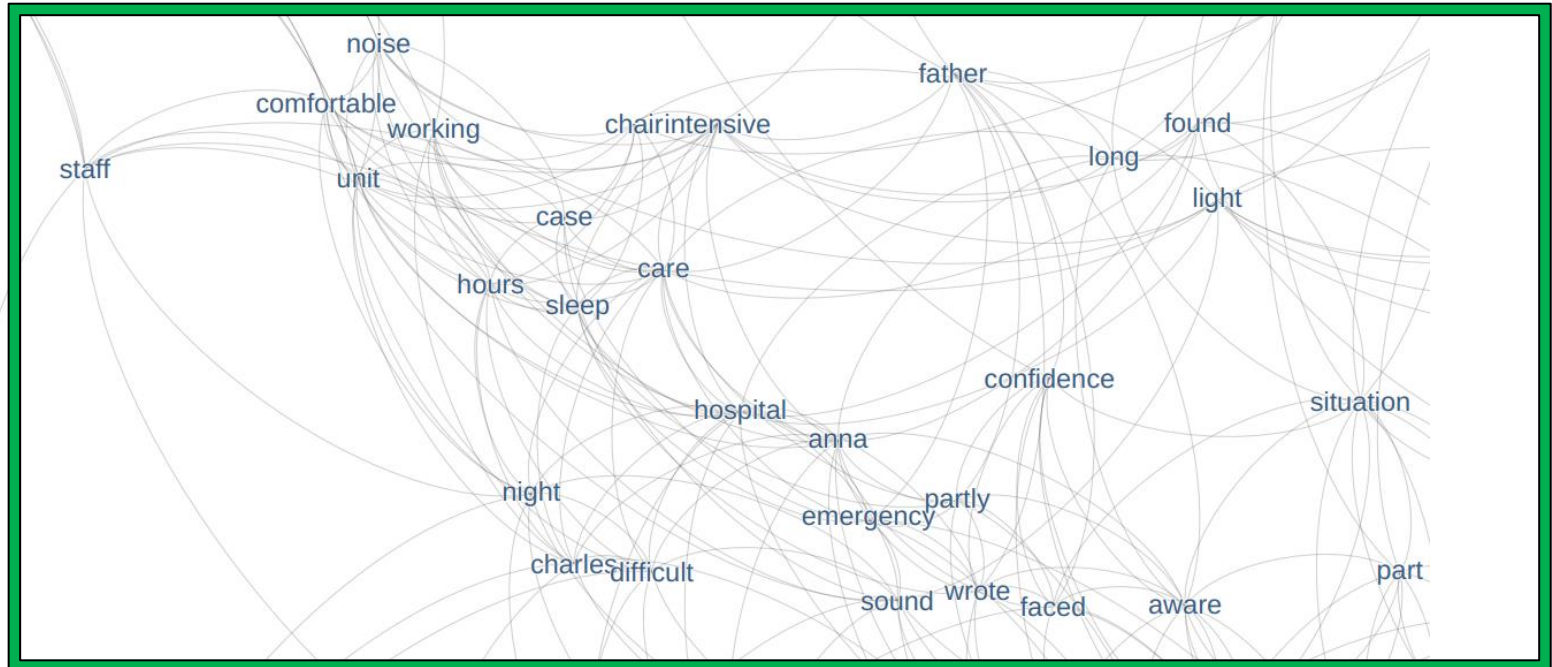


1.2 NLP Approaches to Analysis:

Teaser 3: Conceptual Maps with DIPEX CMI & MS

The conceptual map (500 w) shows

- Common semantic core (**green**): most patient (e.g. Charles or Anna) feel comfortable, but find it difficult to sleep, partly because of the noise



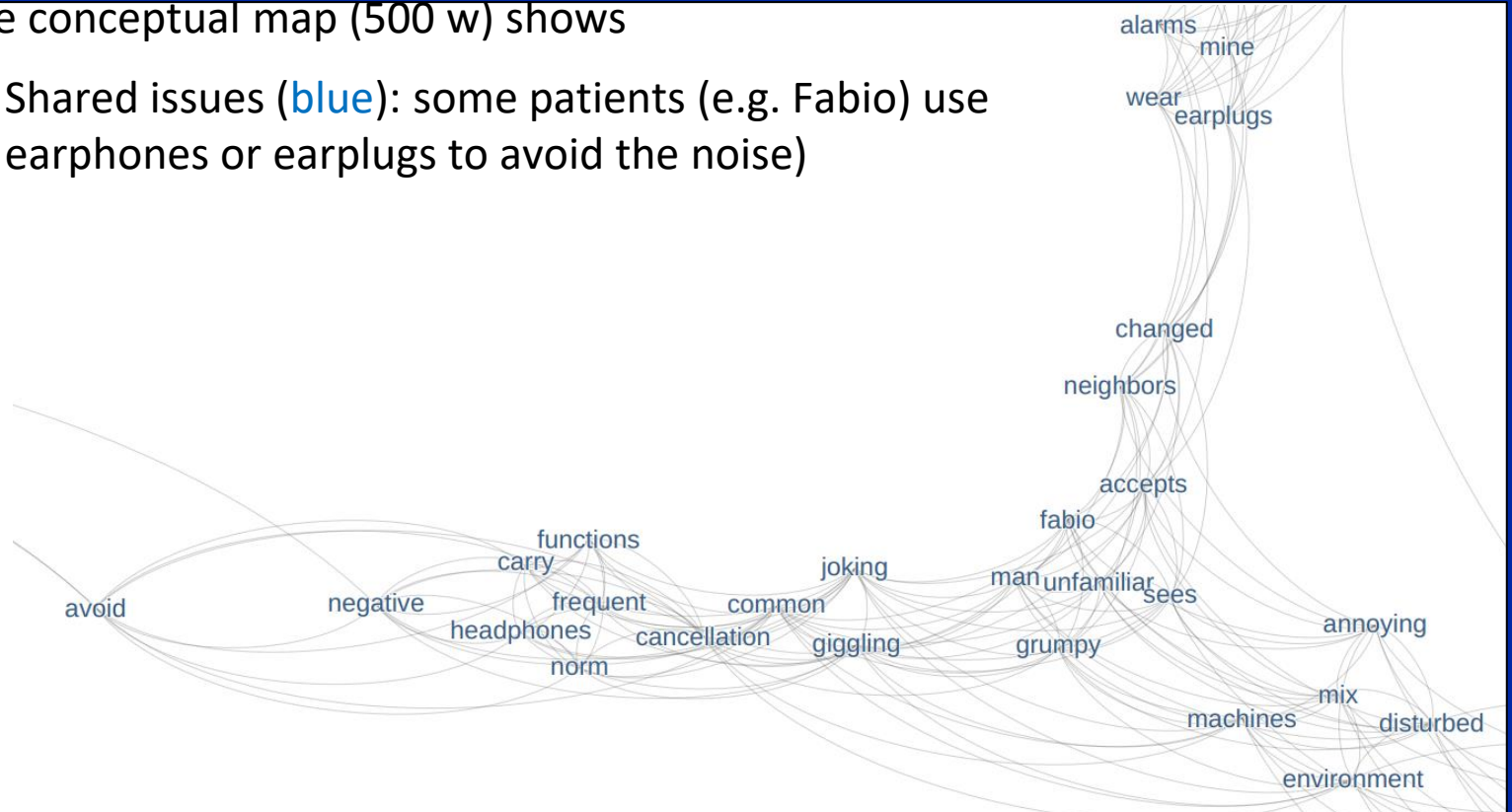


1.2 NLP Approaches to Analysis:

Teaser 3: Conceptual Maps with DIPEX CMI & MS

The conceptual map (500 w) shows

- Shared issues (**blue**): some patients (e.g. Fabio) use earphones or earplugs to avoid the noise)



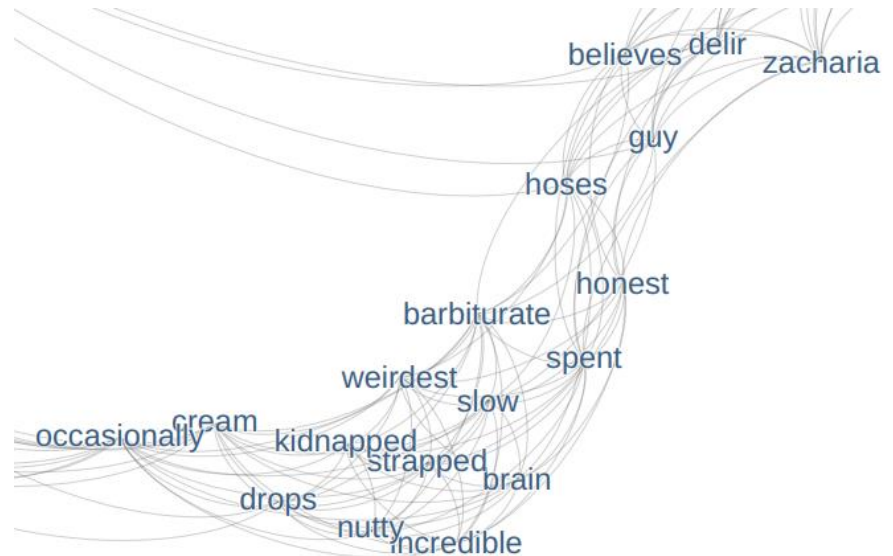


1.2 NLP Approaches to Analysis:

Teaser 3: Conceptual Maps with DIPEX CMI & MS

The conceptual map (500 w) shows

- Individual adventures (coma, hallucinations, **red**):
Zacharia, in delir, believes that he is kidnapped.
He finds that incredibly nutty, himself.





1.2 NLP Approaches to Analysis: Teaser 3: Conceptual Maps with DIPEX CMI & MS

From texts to maps ... and back



“One of them was very funny, he had a curly beard, curly hair and was very, very sensitive, I noticed that ... my friends who later came to visit me in the intensive care unit, they almost fell off their chairs when I first thanked them for being there and singing the song ... my wife and children said stay here, stay here. And then I just prayed to God that He would give me another chance, that I would be allowed to return again”



1.2 NLP Approaches to Analysis:

Teaser 1: Document Classification

Predict EXP_topic based on patient text and the description
(EXP_description_ENG, EXP_textversion_ENG) with logistic regression.

Model Evaluation Met		Model Confusion Matrix:		
Metric	Value	Act \ Pred	generalperecep...	stateofconsci...
Accuracy	0.6829	generalperece...	7	9
Kappa	0.294	stateofconsci...	4	21

Evaluation ^

Feature weights of
class "state of consciousness" →

Feature	Frequency	Feature Infl...
<input type="radio"/> tell	7	2.8189
<input type="radio"/> coma	9	2.7887
<input type="radio"/> dreams	4	2.7257
<input type="radio"/> remembers	7	2.573
<input type="radio"/> because	12	2.5301
<input type="radio"/> <QUESTIONMA...	13	2.356
<input type="radio"/> between	7	2.1189
<input type="radio"/> /	9	2.0057
<input type="radio"/> yeah	3	2.0006
<input type="radio"/> gave	7	1.9681
<input type="radio"/> later	7	1.9596
<input type="radio"/> myself	6	1.912
<input type="radio"/> back	7	1.8767
<input type="radio"/> my	14	1.775
<input type="radio"/> came	8	1.7002
<input type="radio"/> dream	4	1.6395
<input type="radio"/> ms.	2	1.6169
<input type="radio"/> reality	6	1.6115



Conclusion

- Demonstrated the DIPEX DB and three pilot studies
- New insights are possible, but interpretation is important
- The abstraction from individual stories to concepts is not always smooth
- We need more transcribed texts (with or w/o metadata)
- Outlook: DSI Health workshops on medical chatbots



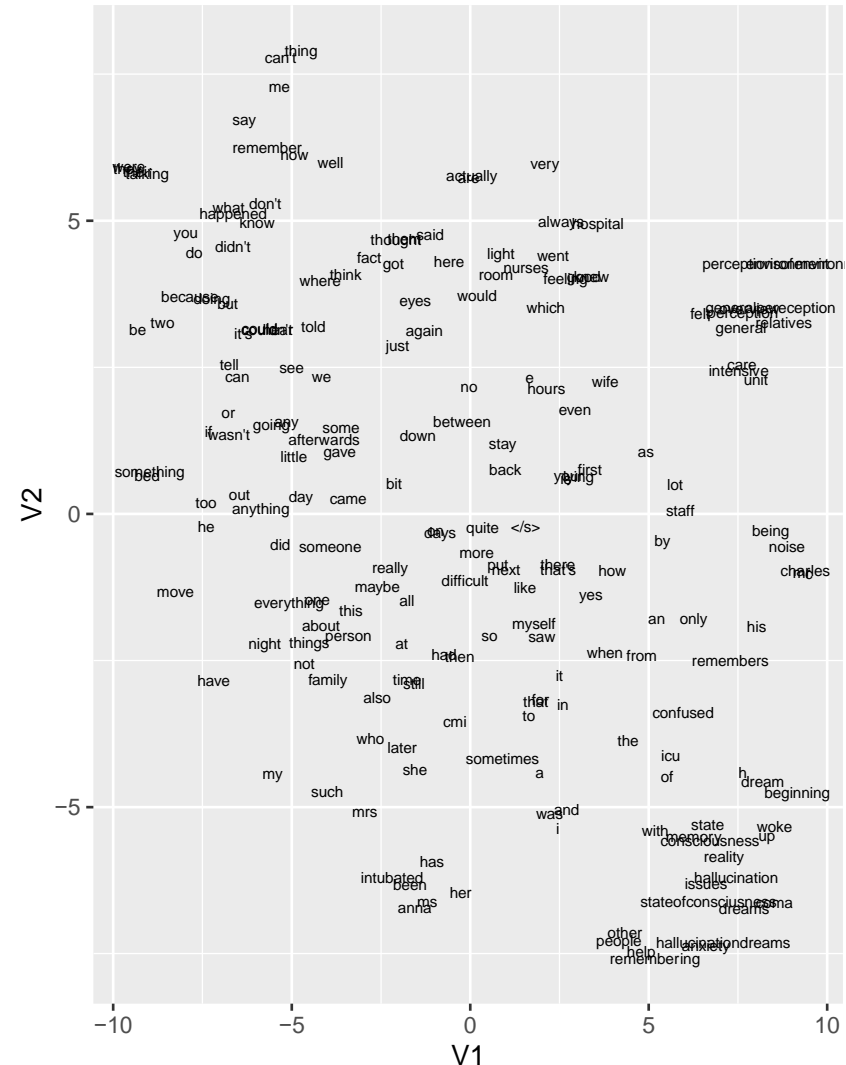
Leading questions for exchange

- Can we learn anything from these approaches?
- How can we evaluate and interpret the results?
- Due to sparseness, our approaches or only just about beginning to work. How can we get more texts?
- What further methods would you like to use?
- How can you imagine using these results in your work?



1.2 NLP Approaches to Analysis: Teaser 3: DIPEX CMI

DIPEX CMI is only 8000 words, too small
for distributional semantics ...





Universität
Zürich^{UZH}

Digital Society Initiative

Institute of Biomedical Ethics and History of Medicine



 **DIPEX.ch**
Erfahrungen mit Gesundheit, Krankheit und Medizin

 **DIPEX.ch**
Health Experiences

Thanks for your attention!