# Multiple Scale Visualization of Electronic Health Records to Support Finding Medical Narratives

S. van der Linden[1], J.J. van Wijk[1] and M. Funk[1]

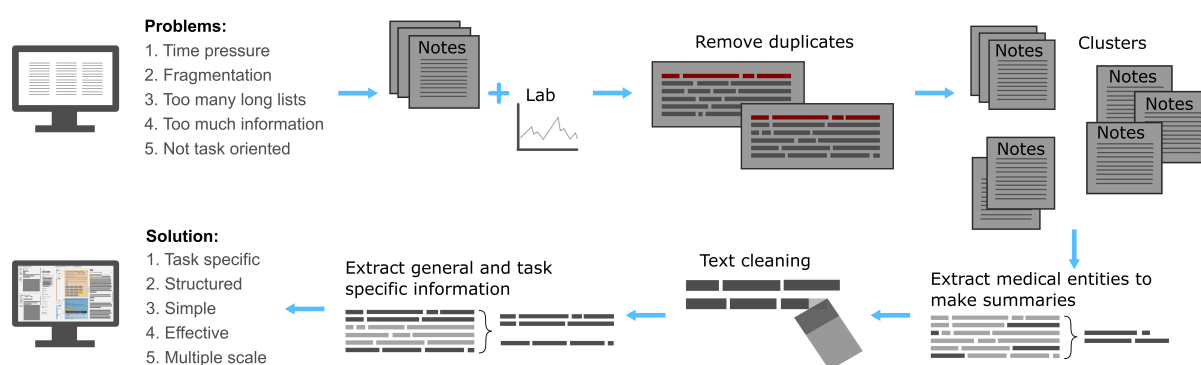[1]Eindhoven University of Technology (TU/e), Netherlands



**Figure 1:** *The system overview describes the current problems of EHR systems and how our system takes patients' notes and lab data, clusters the notes, extracts summaries, cleans text, extracts information and visualizes this to solve the current problems.*

## Abstract

*Electronic Health Records (EHRs) contain rich medical information about patients, possibly hundreds of notes, lab results, images and other information. Doctors can easily be overwhelmed by this wealth of information. For their daily work, they need to derive narratives from all this information to get insights into the main issues of their patients. Standard solutions show all the information in linear lists, often leading to cognitive overload; research solutions provide timelines and relations between the notes but provide too much fragmented information. We propose MEDeNAR, a system for enabling medical professionals to obtain insights from EHRs based on the different tasks in their workflow. The key aspects of our system are the introduction of an intermediate level that summarizes the information using clustering and NLP methods. The results are visualized along a timeline and provide easy access to the detailed descriptions in notes and lab results at the EHR level. We designed the system using an iterative design process in collaboration with 18 doctors, two nurses and 14 domain experts. During the final evaluation, the doctors ranked our system higher than a standard baseline solution and a variation for the used NLP methods.*

## CCS Concepts

• *Human-centered computing* → *Visualization toolkits;* *User interface toolkits;*

## 1. Introduction

An electronic health record (EHR) can consist of over thousand numeric data points, such as laboratory tests and several hundred clinical notes or images [Hal08]. For example, in a hospital in the US, a chronic kidney patient has on average 338 notes collected over 14 years [PE15]. Medical narratives are the textual information in EHRs; reports about the patient's condition and progress. Clinicians rely heavily on these narratives because they accurately identify diagnoses [JB15] and contain complex tem-

poral patterns for medical evidence [SBWC18]. People designed EHRs to, among other things, lower the amount of preventable medical errors [HTE*12]. However, the amount of data collected in the EHRs leads to new challenges [HTE*12], especially cognitive information overload [PE15]. This overload can have negative consequences regarding the patient's safety [PE15]. It exists due to several problems; see 'Problems' in Fig. 1. First, physicians and nurses need to balance the short amount of time they have per patient and taking sufficient time to go through that person's med-

ical history [HTE*12]. In general, doctors review less than 20% of the EHR notes [FRM*12]. Second, the structure of EHRs leads to fragmentation of the patient's data [SBWC18], which hinders the physician's ability to develop a coherent and complete patient narrative. Third, the current structure of EHR systems contains too many long lists of text. These fall short of organizing patient information. Fourth, there is too much information in these notes. To make it worse, due to the number of notes and lack of standardization, clinicians often copy and paste text to other notes [PE15; FRM*12], which leads to redundancy. Also, according to the interviewed doctors, EHRs treat text as static and unsearchable entities, which makes filtering parts of the text challenging. Fifth, the current visualizations do not cover all the tasks in the current EHR information retrieval workflow and do not differentiate between the different tasks of the doctors. In general, not many research solutions have been translated to clinical usage [SBWC18].

We aim to address these five problems while keeping the text accessible, not raising trust issues and increasing the translatability of this research to clinical use. Therefore, we developed an approach for browsing EHRs. Our contributions are:

- We provide a structured view with different levels of detail.
- We support the need of doctors to quickly retrieve information for their different tasks.
- We show how simple visualizations are sufficient to reduce the cognitive load.

We combined an iterative design process, user evaluations, NLP and visualization techniques to achieve our goal. The data analysis included removing duplicates, clustering notes around similar problems and making word summaries. Furthermore, the system extracted the task-specific information for the current task of the doctor. We designed a structured, multiple scales, simple and effective visualization called MEDeNAR (see Fig. 1). Doctors compared it to a baseline and alternative visualization in a user study.

## 2. Related Work

This section describes the related work for text analysis and visualization of EHRs from a single patient perspective.

**Text Analysis of EHRs for Single Patients.** Summarization techniques are a common approach to reduce the overwhelming amount of information. First, one can produce summaries by copying sentences or phrases from the input corpus [JB15; STB*17; SBWC18; HTL*15], known as *extractive summaries* [PE15]. However, they fall short of delivering readable results. Second, new text can be generated based on the original text [Hal08], known as *abstractive summaries* [PE15]. Currently, researchers often implement them by using neural networks [Lop19].

Researchers use techniques at several levels of depth to categorize the text or make extractive summaries. From low to high depth, first, NLP techniques, such as term frequencies, can extract and categorize phrases [HPM16; BHW*09], known as *word-level similarity*. Second, vector space comparison or ontology knowledge can identify similar semantic meaning between phrases [PE15], known as *concept-level similarity*. Third, patient personalized similarity identification can remove statement redundancy, known as

*statement-level similarity* [PE15]. We could only find two examples that partly use this; the V-model system [CP12] removes duplicated events, and another system highlights redundancy [FRM*12]. However, this does not include sentences similar in meaning.

**Visualizations of EHRs for Single Patient.** Worldwide, Epic is one of the most commonly commercially used systems. It has the biggest market share (29%) in the US [TW20]. According to interviewed doctors (D1-3), all commercial systems they used, including Epic, suffered from the five problems described in Fig. 1. In the research community, researchers focused on visualizing patient cohorts and numerical results, such as Wang et al. [WMH*21]. Our focus is on supporting doctors in treating individual patients. One of the first and most common systems to visualize single patients' histories is LifeLines. Its timeline with categories structure has become a baseline for EHR visualization [FBP13], and researchers use it abundantly to visualize narratives [HPM16; FLE06] for individual patients [MGHB04; HTE*12; Hal08; BHW*09; BAK07; HTL*15]. Researchers sometimes also replaced the text with graphical summaries, decreasing real-world adoption and trust because important text features are often omitted [SBWC18]. Sultanum et al. [SBWC18] explored how text visualization of individual patients could assist doctors. However, their timeline still caused an information overload, and it was not easy to locate the needed information. Recently, authors focus on cause and effect relations, such as V-model [CP12]. Moreover, the events on the timelines can be colored based on whether a finding improved or not [HTE*12]. Also, trees [STB*17] and networks [HTE*12] are used to show relations. Finally, Rajkomar et al. [ROC*18] used events for disease progression predictions.

Overall, it remains hard to quickly derive narratives and locate the needed information for each task, see problems 4 and 5 in Fig. 1. Also, there is a need to link instances on the timeline to related (possibly quantitative) material [CP12] to reduce fragmentation (problem 2 in Fig. 1). We aim to address these challenges.

## 3. User Research and Tasks

Different tasks and design requirements were identified based on interviews and workshops with doctors (D1-7) and nurses (N1,2) and a literature review. We only focused on reviewing the data in the EHR. Five main tasks were identified:

T1: Long-term check-up (the patient comes to the hospital for a check-up exam or progress talk with the doctor) where the doctor knows the patient.
T2: Long-term check-up where the patient is unknown.
T3: Diagnosing (newly) admitted patients.
T4: Check-up after an intervention while the patient is admitted.
T5: Diagnosing/admitting a patient in the emergency department.

See Table 1 for the components the narrative for the different tasks consists of and the user's actions to find this narrative. In all narratives, the doctors need the medical history to look up or explore past diseases/treatments or compare exam results. Only the location of this information in the EHR differs per task. The visualization should enable the user to obtain more of the patient's narrative in a certain period compared to the current hospital systems. To this end, we made the following design requirements:

| Tasks | Narrative | | | | | Interactions | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Patient's medical progress | Original diagnosis/ issues | Treatments/ issues in the past | Current issues/ physiology | | Look-up most recent notes | Summarize status of patient | Explore EHR for information | Compare findings | Browse through EHR | Locate specific information |
| T1 | X | | | X | | X | X | | X | | |
| T2 | X | X | X | X | | X | X | X | X | | |
| T3 | | | X | X | | | X | X | X | | |
| T4 | X | | | X | | | X | | X | X | X |
| T5 | | | X | X | | | X | X | X | | |

**Table 1:** *A description of the narrative per task and the actions conducted to find the narrative. The patient's medical progress column describes the reaction to the medication, changes in physiology, how the patient is feeling and the factors responsible for these changes. Look-up, explore, browse and locate are differentiated by permutations of the target and location of the data being known or not [BM13].*

DR1: All the notes need to be available alongside the abstracted visuals because the flexibility and expressiveness of medical narratives make it impossible to structure every aspect of it and make an all-encompassing vocabulary [STB*17].

DR2: The visualization should provide an aggregated form of the data to help the user summarize the data.

DR3: The visualization should display temporal patterns to understand the patient's illness [STB*17; MGHB04; CP12; Hal08] and recovery progress over time.

DR4: The visualization should provide task-specific information.

DR5: The visualization should help with the interpretation of the data. Interpretation of the data is difficult because EHRs fragment the narratives over multiple notes [STB*17; SBWC18]. Also, to help speed up information retrieval, balanced levels of granularity are needed to give insights on an overview and detailed level [STB*17; SBWC18; BHW*09].

## 4. Data Analysis

The data was aggregated (**DR2**) using clusters and word summaries. The notes and lab data were queried from the MIMIC-III data set [JPS*16] for each patient, see Fig. 1. MIMIC contains de-identified data of actual patients. For each patient, duplicated (parts of) sentences were removed from the notes using the algorithm from Rankin et al. [RBD20] if the amount of duplicated information was more than 15%. Below 15%, the duplicated information mainly consisted of false positives .

**Clustering.** A cluster, in our case, is a collection of notes of one patient about the same topic, see Fig. 1. Clusters should help to decrease the fragmentation of the narratives (**DR5**). D6 identified four different types of clusters; hospital admissions (1), check-ups (2), small interventions (3) and emergency room (ER) without hospitalization (4). The analyzed patients only had the first two types. Clustering each admission was easy because notes had a unique admission id. For check-ups, the data was pre-processed to stemmed tokens. Next, we compared three common clustering approaches (agglomerative, partitional and density-based clustering) to the results of clustering one patient by hand. The density-based clustering (HDBSCAN algorithm [MHA16]) was chosen because it gave the best results and detected outliers. We added two additional rules. First, two check-up notes with an admission occurring between them cannot be in the same cluster. This rule enforces a 'before' and 'after' admission, where the earlier typically relates to complaints leading to the admission. Second, check-up notes that

were written on the same day often belong together. Therefore, if multiple clusters from the HDBSCAN algorithm contained notes from the same day, these clusters were merged into one cluster.

**Entity Extraction and Salience Detection.** For each cluster, the five most important entities/words for diseases (including diagnoses), drugs and treatments were extracted from the MIMIC notes to make word summaries, see Fig. 1. D6,7 and N1 recommended limiting the summaries to these categories. Section 5 describes the alternatives we considered instead of word summaries. Five words were chosen because this was a good balance between providing aggregate information and keeping the information limited.

First, medical entities were extracted from the notes of one cluster by using two publically available tools; scispacy (bc5cdr model) [All20] and Medtagger [OHN13]. The scispacy model achieves performance within 3% of state-of-the-art models but cannot tag treatments specifically. Therefore, the common tool Medtagger, which tags all categories, was also used. Entities from the drugs and disease category needed to be extracted by both algorithms. For all entities, the semantic types from UMLS metathesaurus version 2019AB [NIH20] (if the entity could be mapped to a UMLS concept) were also checked to filter out misclassifications. Second, hierarchical structures or salient detection extracted the five most salient entities/words for each category. The paper from Moon et al. [MLC*19] was the only identified one which used salient detection in medical text. UMLS' vocabularies have a hierarchical structure which could help in detecting salient information. Therefore, we evaluated the results of these two approaches with D7. Based on this, the algorithm used ICD-10, international classification of diseases, from UMLS for the disease category and salient detection for the drugs and the treatment category.

For the hierarchical approach, a tree was temporarily made for each cluster with the entities from the notes as leaves. The entities were mapped to ICD-10 concepts. The number of ancestors and relations to the other entities of that leave were determined by the relations of the mapped ICD-10 concepts. The higher the level of depth, the more specific the entities were. Each node had a leaf count; the number of leaves the subtrees of that node led to. The algorithm took the following steps to get a salient entity:

1. It calculated a percentage for each child of the root by dividing their leaf count by the total number of leaves.
2. It selected the child with the highest percentage and leaf count $> 0$, e.g., child A.
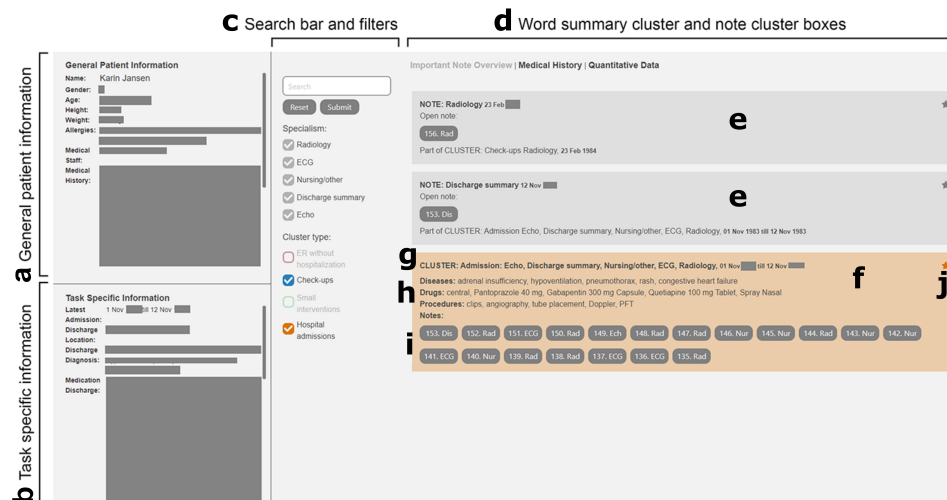
**Figure 2:** *Important note overview view. MEDeNAR always displays the general patient information (a) and task-specific information (b) on the side. It displays the search bar and filters in the middle (c). The default tab is the important note overview (d). The system automatically places the most important notes (e) and clusters of notes (f) in boxes based on the doctor's tasks. Each box consists of a title (g), words summarizing the content of the cluster based on diseases, drugs and treatments (h), links to the individual notes belonging to that cluster (i) and a star to mark notes as important (j). Currently, none of the notes in a cluster are opened. The MIMIC patient data is greyed out.*

3. It walked a path from child A to one of its descendants that is a leaf, where the node with the highest, nonzero leaf count was chosen at each level.
4. The entity of this leaf became one of the five chosen entities.
5. All leaf counts of the nodes in this path were decreased by one to ensure this leaf was not chosen again.
6. The percentage corresponding to child A was decreased by 20% because the algorithm needed to choose five salient entities.

For salience detection (inspired by Moon et al. [MLC*19]), the algorithm computed a dense matrix with tf-idf (term frequency-inverse document frequency; indicates the importance of a word in the document and corpus) scores per note category, such as ECG, for all words in all the categories. Next, it queried 100,000 notes at random of all patients in MIMIC-III for each note category. Afterward, the algorithm took the following steps to get a salient entity:

1. It calculated the occurrence of each note category in the set of all notes of one patient in percentages.
2. For each extracted entity, it retrieved the synonyms from UMLS. They were stemmed, tokenized and stop words were removed.
3. For each entity and synonyms, it calculated one tf-idf score for each of the note categories of that patient by summing over the tf-idf scores of the corresponding category of that entity and synonyms in the matrix and dividing it by the number of scores.
4. Similar to the hierarchical structures, the algorithm chose the entity with the highest tf-idf score from the note category with the highest occurrence percentage as a summary word.
5. This entity received a penalty so that it was not chosen again, also not for another category. The occurrence percentage of that note category was decreased by 20%.

These algorithms led to the word summaries for each cluster. The general and task-specific information (**DR4**) were queried using the

headers in the notes, see Fig. 1. We made these rules in discussion with doctors to determine which notes were placed in the important note overview, for example, the latest discharge note and the check-up clusters after discharge for T1.

## 5. Design Process

We developed MEDeNAR using an iterative design process; cycles of specifying the context and requirements, ideating solutions, developing a prototype and testing it with users. This section describes the insights per iteration, see the supplemental information for figures. Elements still present in MEDeNAR are marked italic.

**Iteration 1.** The first iteration was based on three co-creation sessions; designing solutions with the user in a workshop setting with domain experts P1-3, with P4,5 and with N1,2 and D4. It contained *general patient information and a timeline*. Everyone (except N1,2) mentioned these in their designs. N1,2 and D4 mentioned similar findings as in literature [STB*17], that including everything in the abstraction might be impossible. Therefore, the *links to the original information* should be present. This iteration displayed the data on two scales; the timeline and the notes. During these sessions, the participants did not use network or matrix visualization examples, only simple timelines and scatter plots. According to them, this was enough to get all the needed information quickly. Moreover, the participants mentioned that they would like to give the notes *some degree of importance*.

**Iteration 2.** The second iteration had one extra scale; *the timeline pointed to clusters containing summaries linked to the original notes and lab data*. This scale helped the doctor decide if he needed to read the notes of a cluster in detail. Moreover, this iteration had a compare view to compare different states of the patient to promote interpretation (**DR5**). Also, it was possible to *bookmark*

**Figure 3:** *Medical history view. It displays all patient data. MEDeNAR displays the general patient information (a) and task-specific information (b) on the side. It displays the search bar and filters in the middle (c). Notes are clustered in reversed chronologically ordered colored boxes (l, second scale) and displayed on a timeline (d, first scale). Per cluster box of patient notes, a colored bar (f) is displayed on the timeline (e) to indicate the duration, and it is connected to the corresponding cluster box (g). Dashed parts (h) of the timeline indicate time jumps when no patient data was recorded. These contain the start (i) and stop date (j) and the duration (k). The user can click on a word in the cluster word summary (m), then that word is highlighted in all notes (n). The user can open the notes of a cluster on the side (n, third scale); the system displays them in a list from new to old. Each note has a title (o), text (p), bookmark (q) and attached quantitative data button (r) to go to the lab data belonging to that note. The MIMIC patient data is greyed out.*

notes to prioritize them. D1-3 and P6-11 tested this iteration. They thought the timeline with the summary structure and lab data was time-efficient. Also, bookmarking notes was considered valuable. However, the space could be used more optimally, and there were not enough scales displaying different levels of granularity.

**Iteration 3.** The third iteration had one more scale (**DR5**). *Each cluster on the timeline had a summary* of a few lines of text, a more extended summary of a few paragraphs (new scale) and all the notes in that cluster. These bypassed the need to open the entire note. Also, the summaries contained snippets of the lab data. D2,6,7, N1,2 and P6-11 tested this iteration. Extractive summarization methods were used to make the summaries. We used cosine similarity with different embeddings (GloVe and tf-idf) to get the similarity between sentence vectors from the notes of one cluster. Cosine similarity measures the cosine angle between two vectors to see if they are similar (pointing in the same direction). Also, a skipgram Word2Vec model was used to get sentence vectors, and k-means clustering clustered similar sentences. Word2Vec is a neural network that learns which words are contextually relevant and represents them as vectors. Skipgram is the model architecture.

Doctors and nurses said that visualization's structure was beneficial to their workflow, would reduce scrolling and clicking and the clustering helped find relevant notes. However, the participants could not extract all the narratives partly because the summary quality was insufficient. Abstractive summarization techniques were explored to improve the summary quality (**DR4**). How-

ever, the current state-of-the-art still has drawbacks due to several challenges, such as ungrammatical [Lop19] sentences. Also, users were hesitant about the idea of newly generated text. Therefore, we decided to use *keyword-based abstractions* of medical entities to form word summaries. These should be well suited, considering the fragmented nature of the data [FLE06]. Only, the more extended keyword-based summaries were messy, according to users, so they wanted three scales again. Based on this MEDeNAR was designed.

## 6. Visualization

This section describes the system MEDeNAR and its interactions, see Fig. 1. It was implemented in Vue.js. See section 5 for some of the design choices and the supplemental information for a link to the video and code.

**General and Task-specific Information** The general patient information and task-specific information are always displayed. Based on interviews with D1-4,6,7 and N1,2 the general patient information consists of basic patient information, such as a brief medical history, see Fig. 2(a). Based on interviews with D6,7, the task-specific information displays relevant basic information based on the current task of the doctor (**DR4**). For example, the system displays the admission date, discharge location, discharge diagnosis and the medication after discharge for T1 after hospitalization, see Fig. 2(b). Currently, the system asks the doctor upfront which task he is doing, but this could partly be retrieved from the doctor's calendar.

**Figure 4:** *Buttons of the notes containing the search term or synonyms are marked yellow; by hovering over them, a snippet of the sentences containing the search term are displayed in a tooltip.*



**Figure 5:** *Illustrates how timeline handles consecutive clusters bars (a), clusters bars ending on the same date (b) and overlapping cluster bars (c).*

**Important Note Overview.** The system's main part consists of three tabs; the important note overview (default), the medical history and the quantitative lab data. It displays an automatically extracted, reversed chronologically ordered list of relevant notes or clusters of notes based on the doctor's task (**DR4**). For example, for T1 after hospitalization, the discharge note and all check-up clusters after discharge are displayed. The system displays the relevant notes in grey boxes (Fig. 2(e)) and the relevant clusters in colored boxes, Fig. 2(f). The color represents the type of cluster, see section 4. An orange box is an admission cluster, and a blue box is a check-up cluster (**DR5**). These colors are distinguishable, black text on top of it is readable and users preferred this color scheme. Especially, the grey note boxes are grey such that they are not mistaken for clusters. The cluster boxes ensure there is no spatial split attention effect; information sources about mutually referring information are not spatially separated [VS10]. This could reduce the cognitive load. Every box contains a title (see Fig. 2(g)), links to the notes as a reversed chronologically ordered list of dark grey buttons (see Fig. 2(i), **DR1**) and the word summary containing the most important words for diseases, drugs and treatments (see section 4) extracted from the notes in that cluster (see Fig. 2(h), **DR2**).

Alternatives of the box content (placement of elements, amount of information and which information) were discussed with users. For example, snippets of the lab results were not added to the boxes because, according to users, they cannot be compared to other values and, therefore, lose their value. The grey boxes, Fig. 2(e), have one dark grey button linking to one note because these boxes represent only one relevant note. The note category is displayed to indicate the type of information in the link location to enable navigation [DL07]. For example, 'Rad' for radiology. The user can remove boxes by deselecting the bookmark (star, see Fig. 2(j)). The important note overview gives a more isolated overview compared to displaying all the information. This relates to the simple-to-complex strategy [VS10], which could reduce the cognitive load.

**Filters and Search.** Search functions and filters are available, see Fig. 2(c). The search function uses the stemmed words and synonyms of the search term to search the notes for matches; these are displayed under the search box. The buttons of the notes in the cluster boxes that contain the search term or a synonym are marked yellow, see Fig. 4. The search term and synonyms are marked yellow in the notes. Two types of filters are available; the note category (such as radiology) and cluster box types (such as admission). The user can make combinations of these two filters.

**Medical History.** The medical history tab displays all the medical notes of the patient (including the ones from the important note overview) in colored cluster boxes on three scales; a timeline
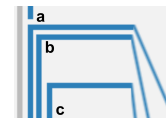
(Fig. 3(d), **DR3**), note clusters (Fig. 3(l)) and the notes (Fig. 3(n)). The cluster boxes have the same structure and order as in the important note overview. The timeline displays when clusters are occurring on a time axis (solid grey parts Fig. 3(e)) and skips over periods with no data (dashed grey parts Fig. 3(h)). The timeline uses as little space as possible while still displaying the information clearly. A vertical colored bar on the timeline displays the duration of each cluster box, Fig. 3(f). The vertical bars are placed such that:

- The start and end correspond with the oldest and newest date of the notes of the cluster.
- A line connects it to the corresponding cluster box, Fig. 3(g).
- Crossings are avoided in several ways; if multiple clusters (partly) overlap, the vertical colored bars are displayed next to each other (see Fig. 5c), and if two clusters have the same end time (top location of the vertical bars is the same) the bar placed on the right is slightly shortened and placed slightly lower (see Fig. 5b).
- There is always a small margin between two bars placed below each other, see Fig. 5a.
- There is a minimum and maximum length to ensure one bar does not use up all the space and bars do not disappear. The timeline has the same height as all the cluster boxes, and all the timeline elements are scaled in proportion to this.

The skips over periods with no data consist of an end date of the period (Fig. 3(i)), a start date (Fig. 3(j)) and the duration (Fig. 3(k)). The duration indicates the length of the interval compared to other intervals (horizontal grey bar), and the absolute duration is written under it, for example, one month and 17 days (Fig. 3(k)). These interval indications are beneficial for timelines with large chronological extents and for compressing timelines with non-uniformly distributed events [BLB*16], which is the case.

Furthermore, the user can open the notes of a cluster on the side in a reversed chronological list by clicking on a dark grey note button in a cluster, see Fig. 3(n). The background of that cluster is darkened. Each note has a title (Fig. 3(o)); a bookmarking star (Fig. 3(q)); the text (Fig. 3(p)); and a button to go the related lab data (Fig. 3(r), **DR5**). This content was determined together with users and experts. By clicking on the bookmark of a note or cluster, the user can place it in the important note overview. Moreover, by clicking on a term in the word summary in a cluster box, the user can highlight that word in the notes and the cluster boxes, see Fig. 3(m). These same actions of opening notes and highlighting words are also available in the important note overview view.

**Quantitative Data View.** The third tab is the quantitative data tab. The user can select it via the top menu. It displays all the lab data variables in a drop down menu. The user can select one vari-
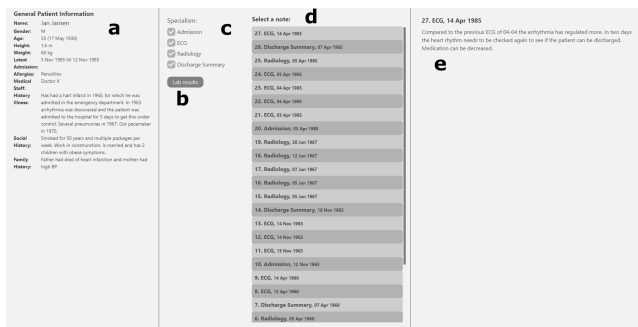
**Figure 6:** *Simple list visualization. The general patient information (a), filter options (c), a button to the lab results (b) and a reversed chronologically ordered list of all the notes (d) are displayed. When the user opens a note, it is displayed on the right (e).*



**Figure 7:** *A cluster box with the concept space word summaries. An accuracy percentage was added after each summary word.*

able and display it in a scatter plot. If the user clicks the 'Attached quantitative data' button of a note, this tab displays the related lab data with that note on the side to remove the split attention effect.

## 7. Evaluation Method

14 doctors (D1,5,7-18) from different specialisms and different regional hospitals participated in the final user test, with on average five years of working experience; five were male, and nine were female. Three different visualizations were compared; MEDeNAR, MEDeNAR with a different approach to make the word summaries, called 'concept space visualization' (Fig. 7) and a simple baseline list visualization (Fig. 6). The concept space visualization was made in collaboration with Fabian Viehmann, who based his approach on semantic concept spaces [EKC*19]. Semantic concept spaces is a new approach that considers semantic relations between keywords and documents in topic models [EKC*19]. It is interesting to see if users prefer the predictions from the concept spaces over the more conventional methods, e.g., the salience detection with tf-idf from section 4, to make word summaries. Also, the structure of the baseline visualization is similar to the structure of the notes in commercial EHR systems. All visualizations were displayed locally on our computer. The final user test contained data of three different patients with a similar amount of data and similar data types, e.g., echo, from the MIMIC-III [JPS*16] data set. On average, a patient had 101 notes from 3 admissions, 75 notes from check-ups and 2542 numeric lab data points. Furthermore, per visualization, the user needed to complete two tasks; first T4 and then T1, see section 3. For the latter, the check-up notes after the latest admission and the latest half of the admission notes were removed.

The user test protocol consisted of an introduction with an informed consent form and familiarization with the visuals (10 min), testing the three different visualizations (15 minutes per visualization) and a final semi-structured interview (5 minutes). Participants spoke out loud during the entire test, and the sessions were recorded, transcribed and anonymized afterward. All participants, except for D5, completed all the tasks for all visualizations. A double Latin square distribution determined the order of the visualizations and which patient data was displayed in which visualization
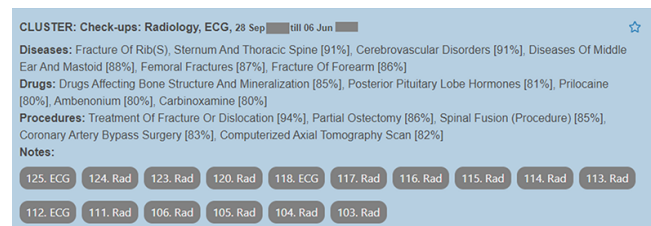
to reduce biases. Testing a visualization consisted of the following steps. After explaining T4 (1), we asked the participants which information they wanted to find (2). The participant had 5 minutes to try to find this information in the current visualization (3). During these 5 minutes, we asked the participants several times how much information in percentages they already gathered (4). These four steps were repeated for T1 with the same patient and visualization. After both tasks, the participant filled in the corresponding SUS [SBWC18], cognitive scales [Zij95; DM08] and scoring components questionnaires. All the remarks from the participants were analyzed using a thematic analysis. A theme needed the support of at least half of the participants.

The user test was kept open-ended to give a satisfying representation of insight capabilities [Nor06] and not restrict participants [SND05]. Based on the anonymous recordings, insights (distinct observations by the participant [SND05]) were identified. Based on recommendations from Saraiya et al. [SND05], the analysis consisted of the following aspects to measure medical insights:

1. A correctness score determined how much of the narrative (their answers in step 2) was found in the insights.
2. Errors determined incorrect or misinterpreted information.
3. A thoroughness score (1 to 5, where one is reading text and five making hypotheses) per insight determined the quality of the answers [SND05; Nor06].
4. Time to first insight and percentage of gathered information over time determined the time it took for the user to get immersed in the data [SND05] and the amount of information learned [SND05].
5. Each insight was expected or unexpected based on step 2.

## 8. Results User Evaluation

The quantitative and qualitative user test results are analyzed together because some participants gave the system lower scores due to the data set structure of the notes or disliking one small thing while being positive in general. This resulted in large standard deviations, while the feedback from the interview was positive. Furthermore, the quantitative results of D15 were omitted due to technical issues. Some measurements, see section 7, were similar for all tasks and visualizations; therefore, they are not discussed further.

**Retrieving Medical Narratives.** The first finding is about how (the components of) the three visualizations contribute to retrieving the medical narratives for the different tasks?'. Overall, MEDeNAR was ranked as the best 77% of the time, see Table 2. The concept

| Visual | 1st | 2nd | 3rd |
|--------|-----|-----|-----|
| MEDeNAR | 10 | 3 | 0 |
| Concept space | 3 | 5 | 5 |
| Baseline list | 0 | 5 | 7 |

**Table 2:** *How often the visuals were ranked first, second or third.*



**Figure 8:** *The results of deriving narratives.*

spaces visualization was ranked second and the list visualization third, which corresponds to the order of the 'quick overview' scores in Fig. 9a. All participants, except for D1, 5, 12 and 15, mentioned that the element contributing the most to the decision of this order was the cluster boxes, then the timeline and then the important note overview. However, the cluster boxes got the lowest usefulness scores (6.2 and 4.9, see Fig. 9b), probably due to the content of the word summaries. The concept space visualization was ranked second because the predictions were messy (containing more text), confusing and not trustworthy. This resulted in a higher complexity (3.95), see Fig. 9a. Only D13 and 18 preferred the concept spaces visualization because it could reveal missed things.

"In the list visualization, there was no real overview" - D11

"It [timeline with cluster boxes] gives a clear overview because it is bundled nicely." - D9

For T4, the participants extracted between 34% (list visualization) and 81% (concept space visualization) of the total narrative, see Fig. 8. The main problems were that the participants could not find the current treatment plan and extract an overview of the physiological changes from the lab data. The general and task-specific information (received one of the highest scores, see Fig. 9b), cluster boxes and important note overview (this component scored a 6.5, see Fig. 9b, but the participants used it 19 of the 28 times) helped the doctors extract the current narrative. Moreover, no particular reason was found why participants extracted more of the narrative in the concept space visualization.

For T1, the participants extracted between 72% (concept space visualization) and 79% (MEDeNAR) of the total narrative, see Fig. 8. For this task, the results are closer, and the standard deviation is smaller because the narrative was easier to find. The participants had problems extracting an overview of the physiological changes from the lab data and the complete progress since the discharge. The task-specific information, the discharge summary in the important note overview and the latest check-up clusters on the timeline helped the participants find the information.

"I think it is very convenient that the system pre-selects notes and puts them in the important note overview." - D10

**Cognitive Load.** The second finding is about the cognitive load in the three visualizations. Although the doctors were not used to the structure of MEDeNAR and the concept space visualization, the mean mental effort scores (43 and 47 respectively, see Fig. 9d) were lower than the known list structure. Moreover, the list visualization had the highest level of irritation (4.6) and difficulty (4.9), see Fig. 9c. Overall, the participants mentioned four explanations for the cognitive load scores. First, the participants disliked the structure of the text of the notes. Second, D5, 8-10, 13, 14, 17 and 18 mentioned that they needed to get used to the system because it
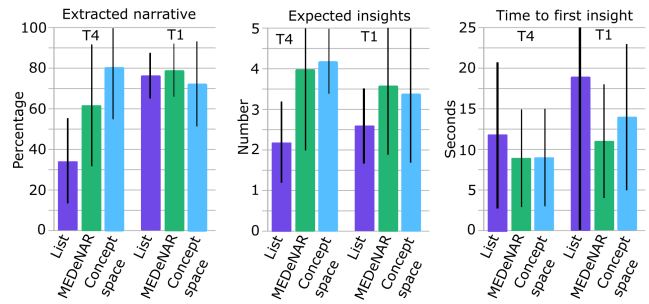
is new. Third, the predictions of the concept space visualization confused the participants (except D8 and 18). Fourth, the concept spaces and MEDeNAR provided a clearer overview compared to the list visualization. Furthermore, on average, participants needed slightly more time to find their first insight in the list visualization (12 and 19 seconds, see Fig. 8). It could be that the important note overview provided information quicker, but this should be investigated further.

"This [list visualization] was not hard because it looks like the current system. Everything is displayed below each other." - D10

## 9. Discussion

This paper describes the system MEDeNAR, which was designed to reduce the cognitive load of forming medical narratives for doctors based on the tasks in their workflow. This research focused on the information overload problem caused by time pressure and the structure of EHR systems, see problems in Fig. 1. The visualizations were kept simple. Also, all the textual notes are still available instead of using graphical summaries, which should increase real-world adaptation and not raise trust issues [SBWC18]. According to the final user evaluation, there were indeed no trust issues with MEDeNAR because the doctors could always go to the original notes and see where the summary words from the clusters occurred in the text. MEDeNAR was compared to a baseline list visualization and a visualization with different word summaries. Users preferred the proposed structure to visualize the notes over the current structure in EHR systems. Moreover, users thought the timeline with the clustering was beneficial to display the notes over time (**DR3**). Also, the word summaries helped with aggregating the data(**DR2**), and linking the notes and lab data helped reduce the fragmentation of the narrative (**DR5**). Participants used the important note overview and task-specific information abundantly and thought it was useful (**DR4**).

The content of the summaries and the display of the lab data could be improved. MEDeNAR could only display one numeric value in the quantitative data tab, which hindered the doctors from getting a quick overview of the most important numerical values, e.g., blood pressure and heart rate. The results have high standard deviations, and participants sometimes gave visualization components a low score, although they were positive in their feedback. The medical narratives could not be completely extracted by the participants, probably because of the data structure. Participants
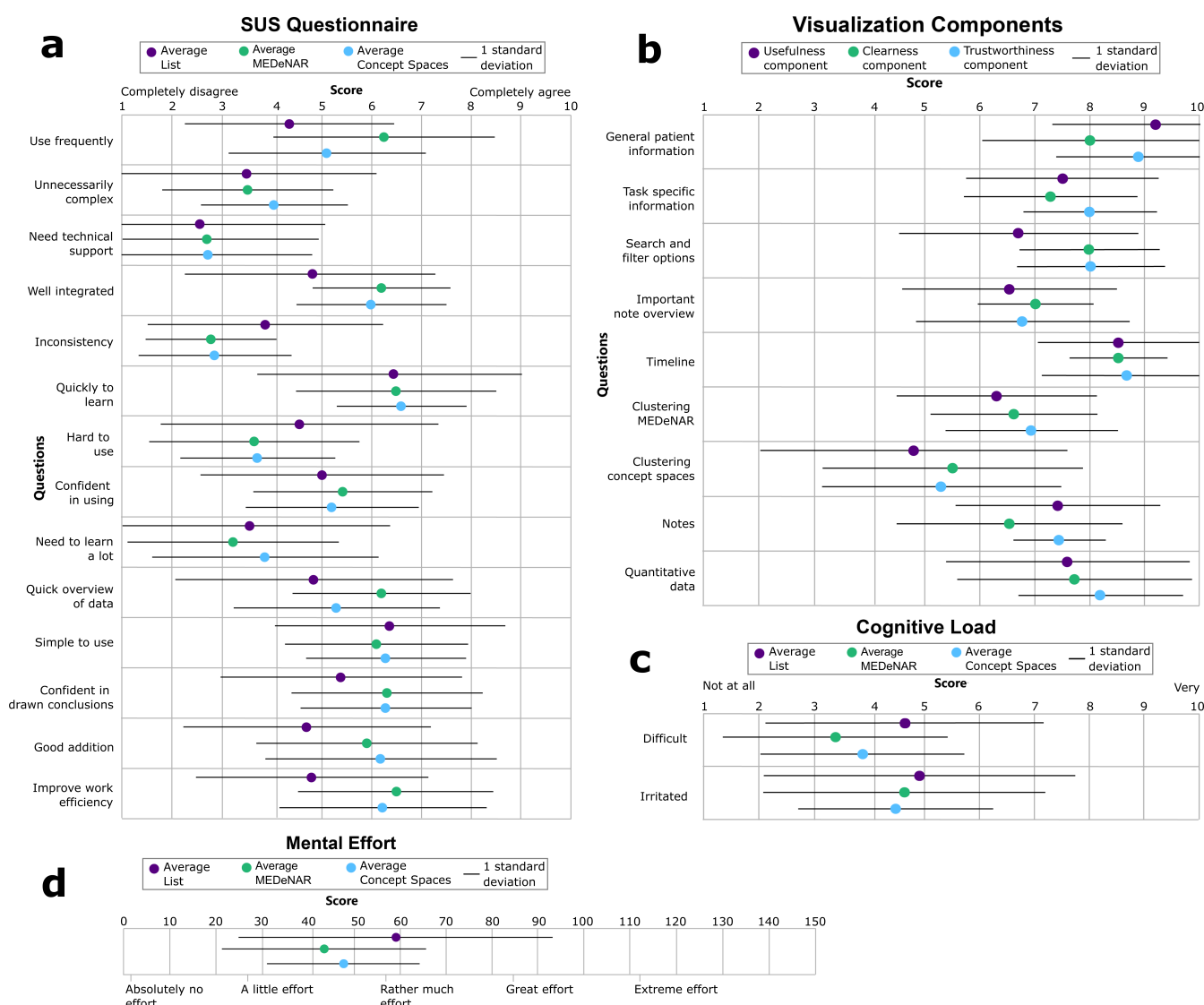
**Figure 9:** *The results of the SUS questionnaire (a), the results of scoring the different components of the visualization (b), the results for the cognitive load (c) and mental efforts (d) ratings.*

mentioned that notes were longer or differently structured compared to what they are used to.

Furthermore, there are ethical concerns that need to be taken into consideration. In current EHR systems, the doctors do not have the time to review all the information. MEDeNAR helps with this problem, but the data aggregations are imperfect, and doctors could still overlook important parts. Questions such as; 'Who would be responsible?' and 'Is it possible to include all the important information?' require further research as well.

## 10. Conclusion

In conclusion, the contribution is to propose a system that potentially reduces the cognitive load of forming medical narratives of

electronic health records. We proposed MEDeNAR. It offers structured and task-specific visualizations based on simple but effective visualizations. Also, MEDeNAR was compared to two other visualizations and tested with 14 doctors. The results suggest that MEDeNAR could reduce the cognitive load and was ranked highest by the participants for extracting the medical narratives.

Future work would include more user tests with a data set with a more familiar structure and testing all the different tasks. Also, it could include aggregations that are personalized for doctors and patients and link the data (**DR2 and DR5**). Currently, MEDeNAR provided doctors with a quicker overview of the data. Future research could look at how this quicker overview of the data helps identify gaps in the data and potential unanswered questions about the patient's disease.

## 11. Acknowledgement

I would like to thank Kees van het Maalpad for all his feedback.

## References

[All20] ALLEN INSTITUTE OF ARTIFICIAL INTELLIGENCE. *scispacy*. https://allenai.github.io/scispacy/. 2020 3.

[BAK07] BUI, A.A.T., ABERLE, D.R., and KANGARLOO, H. "Time-Line: visualizing integrated patient records". *IEEE Trans. Inform. Technol. Biomedicine* 11.4 (2007), 462–473. DOI: 10.1109/TITB.2006.884365 2.

[BHW*09] BASHYAM, V., HSU, W., WATT, E., et al. "Problem-centric organization and visualization of patient imaging and clinical data". *Radiographics* 29.2 (2009), 331–343. DOI: 10.1148/rg.292085098 2, 3.

[BLB*16] BREHMER, M., LEE, B., BACH, B., et al. "Timelines revisited: A design space and considerations for expressive storytelling". *IEEE Trans. Visualization Comput. Graph.* 23.9 (2016), 2151–2164. DOI: 10.1109/TVCG.2016.2614803 6.

[BM13] BREHMER, M. and MUNZNER, T. "A multi-level typology of abstract visualization tasks". *IEEE Trans. Visualization Comput. Graph.* 19.12 (2013), 2376–2385. DOI: 10.1109/TVCG.2013.124 3.

[CP12] CHOI, J. and PARK, H. "V-Model: A new innovative model to chronologically visualize narrative clinical texts". *Proc. 2012 CHI Conf. Human Factors Comput. Syst.* 2012, 453–462. DOI: 10.1145/2207676.2207739 2, 3.

[DL07] DESTEFANO, D. and LEFEVRE, J. "Cognitive load in hypertext reading: A review". *Comput. Human Behav.* 23.3 (2007), 1616–1641. DOI: 10.1016/j.chb.2005.08.012 6.

[DM08] DELEEUW, K.E. and MAYER, R.E. "A comparison of three measures of cognitive load: Evidence for separable measures of intrinsic, extraneous, and germane load". *J. Educational Psychol.* 100.1 (2008), 223. DOI: 10.1037/0022-0663.100.1.223 7.

[EKC*19] EL-ASSADY, M., KEHLBECK, R., COLLINS, C., et al. "Semantic concept spaces: Guided topic model refinement using word-embedding projections". *IEEE Trans. Visualization Comput. Graph.* 26.1 (2019), 1001–1011. DOI: 10.1109/TVCG.2019.2934654 7.

[FBP13] FAISAL, S., BLANDFORD, A., and POTTS, H.W.W. "Making sense of personal health information: challenges for information visualization". *Health Inform. J.* 19.3 (2013), 198–217. DOI: 10.1177/1460458212465213 2.

[FLE06] FANG, S., LWIN, M., and EBRIGHT, P. "Visualization of unstructured text sequences of nursing narratives". *Proc. 2006 ACM SAC*. 2006, 240–244. DOI: 10.1145/1141277.1141331 2, 5.

[FRM*12] FARRI, O., RAHMAN, A., MONSEN, K.A., et al. "Impact of a prototype visualization tool for new information in EHR clinical documents". *Appl. Clin. Inf.* 3.04 (2012), 404–418. DOI: 10.4338/ACI-2012-05-RA-0017 2.

[Hal08] HALLETT, C. "Multi-modal presentation of medical histories". *Proc. 13th Int. Conf. IUI.* 2008, 80–89. DOI: 10.1145/1378773.1378785 1–3.

[HPM16] HERNÁNDEZ, E., POMARES QUIMBAYA, A., and MUÑOZ, O. "HTL Model: A Model for Extracting and Visualizing Medical Events from Narrative Text in Electronic Health Records". *ICT4AgeingWell*. 2016, 107–114. DOI: 10.5220/0005863501070114 2.

[HTE*12] HSU, W., TAIRA, R., EL-SADEN, S., et al. "Context-Based Electronic Health Record: Toward Patient Specific Healthcare". *IEEE Trans. Inform. Technology Biomedicine* 16 (2012), 228–34. DOI: 10.1109/TITB.2012.2186149 1, 2.

[HTL*15] HIRSCH, J.S., TANENBAUM, J.S., LIPSKY GORMAN, S., et al. "HARVEST, a longitudinal patient record summarizer". *J. Amer. Med. Inform. Assoc.* 22.2 (2015), 263–274. DOI: 10.1136/amiajnl-2014-002945 2.

[JB15] JENSEN, L. and BOSSEN, C. "Factors affecting physicians' use of a dedicated overview interface in an electronic health record: The importance of standard information and standard documentation". *Int. J. Med. Inform.* 87 (2015). DOI: 10.1016/j.ijmedinf.2015.12.009 1, 2.

[JPS*16] JOHNSON, A.E.W., POLLARD, T.J., SHEN, L., et al. "MIMIC-III, a freely accessible critical care database". *Scientific data* 3.1 (2016), 1–9. DOI: 10.1038/sdata.2016.35 3, 7.

[Lop19] LOPEZ, J. "Automatic summarization of medical conversations, a review". *TALN-RECITAL 2019-PFIA*. ATALA. 2019, 487–498 2, 5.

[MGHB04] MAMYKINA, L., GOOSE, S., HEDQVIST, D., and BEARD, D. "CareView: Analyzing nursing narratives for temporal trends". *Proc. 2004 CHI Conf. Human Factors Comput. Syst.* 2004, 1147–1150. DOI: 10.1145/985921.986010 2, 3.

[MHA16] MCINNES, L., HEALY, J., and ASTELS, S. *Basic Usage of HDBSCAN\* for Clustering*. https://hdbscan.readthedocs.io/en/latest/basic_hdbscan.html. 2016 3.

[MLC*19] MOON, S., LIU, S., CHEN, D., et al. "Salience of Medical Concepts of Inside Clinical Texts and Outside Medical Records for Referred Cardiovascular Patients". *J. Healthcare Inform. Res.* 3.2 (2019), 200–219. DOI: 10.1007/s41666-019-00044-5 3, 4.

[NIH20] NIH. *Unified Medical Language System® (UMLS®)*. https://www.nlm.nih.gov/research/umls/knowledge_sources/metathesaurus/index.html. 2020 3.

[Nor06] NORTH, C. "Toward measuring visualization insight". *IEEE Comput. Graph. Appl.* 26.3 (2006), 6–9. DOI: 10.1109/MCG.2006.70 7.

[OHN13] OHNLP. *MedTagger Project Page*. http://ohnlp.org/index.php/MedTagger_Project_Page. 2013 3.

[PE15] PIVOVAROV, R. and ELHADAD, N. "Automated methods for the summarization of electronic health records". *J. Amer. Med. Inform. Assoc.* 22.5 (2015), 938–947. DOI: 10.1093/jamia/ocv032 1, 2.

[RBD20] RANKIN, S.K., BRIGHT, R., and DOWDY, K. "Bloatectomy: A method for the identification and removal of duplicate". (2020). DOI: 10.5281/zenodo.3909030 3.

[ROC*18] RAJKOMAR, A., OREN, E., CHEN, K., et al. "Scalable and accurate deep learning for electronic health records". *npj Digit. Medicine* 1 (2018). DOI: 10.1038/s41746-018-0029-1 2.

[SBWC18] SULTANUM, N., BRUDNO, M., WIGDOR, D., and CHEVALIER, F. "More Text Please! Understanding and Supporting the Use of Visualization for Clinical Text Overview". *Proc. 2018 CHI Conf. Human Factors Comput. Syst.* 2018, 1–13. DOI: 10.1145/3173574.3173996 1–3, 7, 8.

[SND05] SARAIYA, P., NORTH, C., and DUCA, K. "An insight-based methodology for evaluating bioinformatics visualizations". *IEEE Trans. Visualization Comput. Graph.* 11.4 (2005), 443–456. DOI: 10.1109/TVCG.2005.53 7.

[STB*17] SULTANUM, N., THAINE, P., BRUDNO, M., et al. "MedStory : unlocking the qualitative power of medical narratives". *Proc. 2017 Workshop VAHC*. 2017 2–4.

[TW20] TATE, C. and WARBURTON, P. *US Hospital EMR Market Share 2020 Shifting Perspectives Among Large and Small Hospitals*. https://klasresearch.com/report/us-hospital-emr-market-share-2020/1616. 2020 2.

[VS10] VAN MERRIËNBOER, J.J.G. and SWELLER, J. "Cognitive load theory in health professional education: design principles and strategies". *Med. Educ.* 44.1 (2010), 85–93. DOI: 10.1111/j.1365-2923.2009.03498.x 6.

[WMH*21] WANG, Q., MAZOR, T., HARBIG, T., et al. "ThreadStates: State-based Visual Analysis of Disease Progression". *IEEE Trans. Visualization Comput. Graph.* 2021 2.

[Zij95] ZIJLSTRA, F.R.H. "Efficiency in work behaviour: A design approach for modern tools". (1995) 7.