

BT-Nurse: computer generation of natural language shift summaries from complex heterogeneous medical data

James Hunter,¹ Yvonne Freer,² Albert Gatt,³ Ehud Reiter,¹ Somayajulu Sripada,¹ Cindy Sykes,² Dave Westwater⁴

¹Department of Computing Science, University of Aberdeen, King's College, Aberdeen, UK

²Simpson Centre for Reproductive Health, Royal Infirmary of Edinburgh, Edinburgh, UK

³Institute of Linguistics, University of Malta, Msida, Malta

⁴Data2Text, Aberdeen, UK

Correspondence to

James Hunter, Department of Computing Science, University of Aberdeen, King's College, Aberdeen AB24 3UE, UK; j.hunter@abdn.ac.uk

Received 17 February 2011

Accepted 27 May 2011

Published Online First

1 July 2011

ABSTRACT

The BT-Nurse system uses data-to-text technology to automatically generate a natural language nursing shift summary in a neonatal intensive care unit (NICU). The summary is solely based on data held in an electronic patient record system, no additional data-entry is required. BT-Nurse was tested for two months in the Royal Infirmary of Edinburgh NICU. Nurses were asked to rate the understandability, accuracy, and helpfulness of the computer-generated summaries; they were also asked for free-text comments about the summaries. The nurses found the majority of the summaries to be understandable, accurate, and helpful ($p < 0.001$ for all measures). However, nurses also pointed out many deficiencies, especially with regard to extra content they wanted to see in the computer-generated summaries. In conclusion, natural language NICU shift summaries can be automatically generated from an electronic patient record, but our proof-of-concept software needs considerable additional development work before it can be deployed.

INTRODUCTION

Medical professionals have access to increasing volumes of information about patients. This is particularly the case in the intensive care unit (ICU), where continuously monitored physiological data (eg, heart rate, oxygen saturation) and detailed records of observations and interventions are available.

Effective presentation and understanding of these data is important in real-time decision making, but is also very relevant during patient handover between clinicians where, despite an oral or written handover, the outgoing clinician may forget to mention important information, or the incoming clinician may not assimilate all the information presented during a brief exchange. In such cases the incoming clinician relies on the available data in order to fill in any gaps.

While large datasets are usually presented as graphs or tables, some studies have found that high-quality textual summaries can be more effective for decision-support in some circumstances.^{1,2} The summaries in these studies were carefully written by expert clinicians; this is only practical within a research context. However, as part of a larger project, BabyTalk,^{3,4} we have developed a Natural Language Generation system, BT-Nurse, which automatically generates English summaries

of the electronically recorded patient data over a 12 h nursing shift, for a baby in a neonatal intensive care unit (NICU).

CASE DESCRIPTION

The NICU at the Edinburgh Royal Infirmary uses Clevermed's *Badger* computer system to manage and display patient data. This system acquires and records several channels of continuous physiological data sampled once per second. A display is located beside each cot, and clinical staff routinely use this to enter additional information, including hourly physiological measurements, drugs and fluids administered, equipment settings, and care and treatment actions taken. Most of the data collected is preformatted but free-text entry is also available.

BT-Nurse analyses the patient data, decides which information is most important, and presents it as an English text; an extract is shown in figure 1A. Figure 1B shows a corresponding extract from a summary of the same shift data written by a research nurse. The summaries are structured according to physiological system (eg, *respiratory*); BT-Nurse only summarizes two of the 10 physiological systems, viz. *respiratory* and *cardiovascular*; it also reports on the patient's current status and problems.

METHODS OF IMPLEMENTATION

BT-Nurse is constructed around a standard data-to-text "pipeline" architecture (figure 2),^{5,6} where information is processed sequentially by modules which communicate via a domain ontology which includes mechanisms for modeling uncertainty and incomplete knowledge:

1. *Data Translation* transforms data from the format stored in the Badger system to that required by the ontology. A limited amount of information is extracted from free-text, using simple parsing and keyword extraction techniques. For example, the text bolus of sodium chloride infusing would be mapped to an instance of a *drug administration* event in the ontology, with properties indicating the drug (sodium chloride) and the method (infusion).
2. *Data Preprocessing* tries to fill in some of the omissions and gaps which are inevitably present in real-world patient data. For example, if the time of an intubation has not been entered, BT-Nurse attempts to infer approximately when it happened by examining the hourly ventilation observations.

(a)

Respiratory Support**Current Status**

Currently, the baby is on CMV in 27 % O₂. Vent RR is 55 breaths per minute. Pressures are 20/4 cms H₂O. Tidal volume is 1.5.

SaO₂ is variable within the acceptable range and there have been some desaturations.

The most recent blood gas was taken at around 07:45. Parameters are acceptable. pH is 7.3. CO₂ is 5.72 kPa. BE is -4.6 mmol/l. The last ET suction was done at about 05:15.

Events During the Shift

A blood gas was taken at around 19:45. Parameters were acceptable. pH was 7.18. CO₂ was 7.71 kPa. BE was -4.8 mmol/l.

Another ABG was taken at around 23:00. Blood gas parameters had deteriorated to respiratory acidosis by around 23:00. pH was 7.18. CO₂ had risen to 9.27 kPa by around 23:00. BE was -4.8 mmol/l.

The baby was intubated at 00:15 and was on CMV. Vent RR was 50 breaths per minute. Pressures were 20/4 cms H₂O. FiO₂ was 29 %. Tidal volume was 1.5. He was given morphine and suxamethonium. MAP was raised from 6 cms H₂O to 8 cms H₂O.

Between 00:30 and 03:15, SaO₂ increased from 88 % to 97 %.

Another ABG was taken at around 00:45. pH was 7.18. CO₂ dropped to 7.95 kPa. BE was -4.8 mmol/l.

...

(b)

Respiratory support needed due to immaturity**CURRENT MANAGEMENT / ASSESSMENT:**

CMV rate 55, pressures 20/4, in 27% oxygen, giving tidal volumes of 1.5 ml (3.3 ml/kg). Very recent ABG good: pH 7.31, CO₂ 5.72. He received morphine prior to intubation at 00:30; no spontaneous respiratory effort noted since being re-ventilated. Desaturates during cares and suction but recovers afterwards; otherwise SpO₂ has been fairly stable. Large ETT secretions, mucopurulent and blood stained.

EVENTS DURING THIS SHIFT:

While on BiPAP, oxygen requirement increased to 50% by 23:00

ABG at 23:10 showed CO₂ increased from 7.7 to 9.27 in 3 h

Electively re-intubated at 00:30 to CMV rate 50, pressures 18/4 in 30% oxygen

Difficult intubation; size 2 ETT

Was given morphine and suxamethonium

On ventilation, oxygen requirement reduced to 30% and ABG initially improved

...

Figure 1 Extracts from nursing shift summaries: (A) generated by BT-Nurse; (B) written by a research nurse.

3. *Signal Analysis* detects and removes artifacts from the physiological data and extracts a small number of events, both short-term (eg, bradycardias and desaturations) and long-term (eg, trends and abstractions such as “within normal range”). Figure 3 presents the data corresponding to the summaries in figure 1.
4. *Data Interpretation* uses medical knowledge to enhance the information recorded in the ontology: (i) by estimating the medical significance of events; (ii) by deriving higher-level abstractions (eg, the state of having respiratory acidosis); and (iii) by inferring causal and other relationships between events. Medical knowledge is expressed using forward chaining rules derived during extensive knowledge acquisition exercises with domain experts (a consultant neonatologist and senior neonatal nurse).

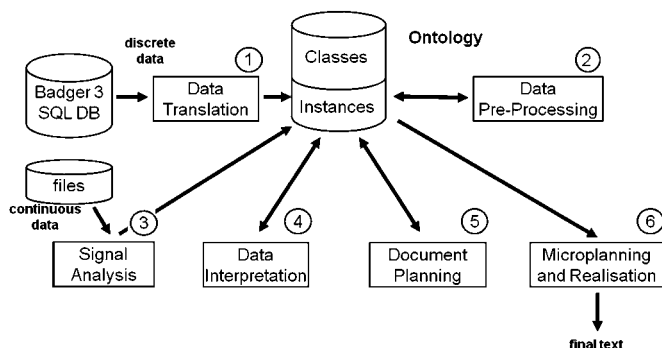


Figure 2 BT-Nurse architecture.

5. *Document Planning* decides on the content and structure of the generated text. Some sections of the text, such as *Current Status* (see figure 1A), essentially have fixed structures which are populated by relevant events. For other sections, such as *Events During the Shift*, the document planner identifies a small number of key events during the shift (based on medical significance) and generates a paragraph around each of these.
6. *Microplanning and Realization* maps the ontology instances selected by the Document Planner to English text by: (i) mapping each ontology instance to a semantic structure using rules which select linguistic predicates; (ii) aggregating the resulting representations into higher-order structures; and (iii) realizing the structures as English text using the SimpleNLG realization engine,⁷ which performs English syntactic and morphological generation.

EXAMPLE AND OBSERVATIONS

BT-Nurse was evaluated on-ward by nurses who read summaries about babies under their care. These summaries were constructed on demand toward the end of selected shifts using the live database and displayed to the nurse at the cot-side. They were generated in less than 1 min, with no noticeable impact on other users of the Badger system. No additional data-entry was required; all information was obtained from the Badger patient record.

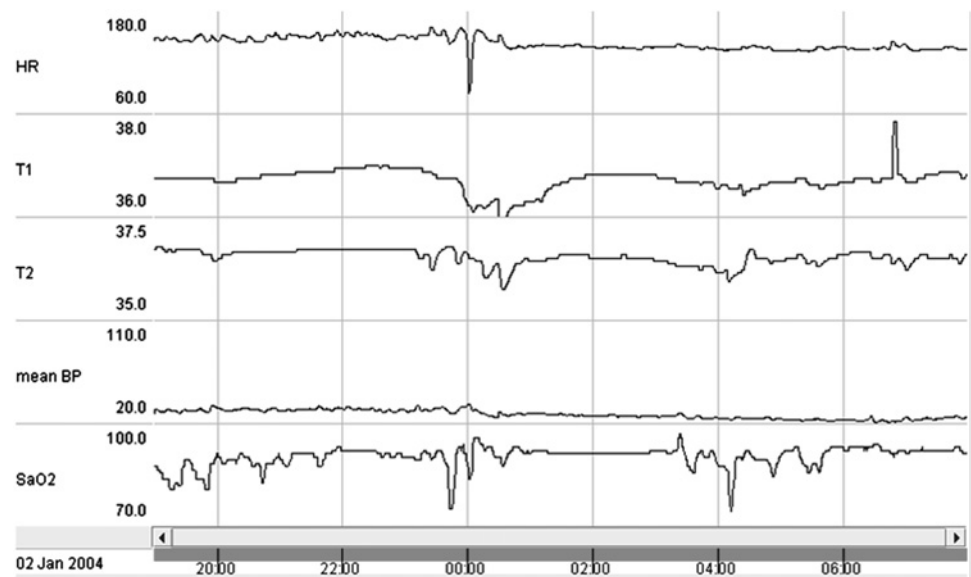
After reading a summary, the *outgoing* nurse in charge of that specific baby was asked to rate its *understandability*, *accuracy*, and *helpfulness in producing her own end-of-shift summary* (were she to produce one), by indicating *agreement*, *disagreement*, or *neutrality* with respect to these statements. She was also invited to enter additional comments on any aspect of the summary. The protocol for *incoming* nurses was the same, except that the final question referred to *helpfulness in care planning*. Incoming nurses had a verbal handover with the outgoing nurse and also had access to the computerized charts; hence they could judge the accuracy and helpfulness of a generated summary. We could not directly compare BT-Nurse texts against human-written texts, using measures such as recall and precision, because NICU nurses do not write detailed textual shift handover reports.

We conducted 165 trials (defined as an evaluation by one nurse of one shift summary): 73 with outgoing and 92 with incoming nurses. A total of 148 summaries were produced for 31 individual babies. In 131 cases, a summary was seen by only one nurse (outgoing or incoming); in the other 17 cases the summary was seen by both nurses. On average, each baby was seen by 2.3 nurses (maximum 6). Of a nursing staff complement of 93, 54 different nurses participated. On average each nurse saw 4.0 different babies; only four nurses saw more than 8 different babies.

We compared response frequencies for each of the categories in each of the three questions that nurses were asked, using a χ^2 test on response frequencies by items (ie, over trials). Overall, there were significant differences between the number of positive, negative, and neutral responses for *understandability* ($\chi^2=241.89$; $p<0.001$), *accuracy* ($\chi^2=110.22$; $p<0.001$), and *helpfulness* ($\chi^2=64.15$; $p<0.001$). As table 1 shows, the majority response was positive in all three cases. A multinomial logistic regression showed no significant differences between incoming and outgoing nurses for any of these questions (*understandability*: model $\chi^2=15.26$, $p=0.08$; *accuracy*: $\chi^2=19.99$; $p=0.1$; *helpfulness* $\chi^2=17.29$; $p>0.7$).

The comments were manually segmented, so that each segment addressed one specific aspect of the summary. This yielded 237 segments (125 for outgoing nurses and 112 for

Figure 3 Example physiological data sampled once per second.



incoming). The segments were annotated independently by three of the authors to indicate which aspect of a summary each was concerned with (*content*, *language*, and *overall*) and which of a predefined set of labels for each category applied. We used Cohen's κ statistic to calculate pairwise agreement between annotators on each dimension; using standard thresholds,⁸ we found *tentative* agreement in the *content* ($\kappa=0.73$) and *language* ($\kappa=0.66$) dimensions, and *good* agreement in the *overall* ($\kappa=0.83$) dimension.

Most segments concerned the *content* of the summary (185), most of these (109) noting *missing content*. Many of these referred to information which was not intended to be included (eg, content about nutrition, which BT-Nurse did not handle, as noted above). However, even disregarding such segments, requests for more content were much more common than requests for less, suggesting that BT-Nurse was under-reporting. As an example, one nurse wrote that the "baby ... is VERY small, and ... it should be pointed out that the ETT is size 2.0". BT-Nurse never reports ETT (endotracheal tube) size as usually this is not very important; however in some cases it is important and should be reported.

There were 46 segments concerning *incorrect content*; some were due to errors in the patient record data, but most were due to bugs in the software. For example, BT-Nurse sometimes listed *current problems* which in fact had been rectified; this was because of an error in reading the relevant database table.

There were only 11 segments about language, all of which were negative. Most of these criticisms reflected individual preferences. The small number of such segments raises the possibility that nurses only commented about language when they were unhappy with it; if this is true, then perhaps in most cases BT-Nurse's language, which is stylistically quite different from the somewhat telegraphic style that typifies free-text comments in normal shift summaries, was "good enough".

Among the 35 *overall* segments, eight concerned deficiencies from a high-level "narrative" perspective—for example, not describing causal links between observations and interventions, and not adequately describing the overall "big picture".

There were also some very encouraging comments about how BT-Nurse summaries were helpful, such as "BT picked up the change in HR trend that I had not noticed".

DISCUSSION

Given that we did not have time to develop a system which generated complete shift summaries, we think it is very encouraging that 58% of the nurses regarded BT-Nurse texts as helpful, and 90% regarded them as understandable.

Most of the criticisms do not concern the underlying technology and could be addressed by expanding BT-Nurse so that it generates complete summaries (including systems that were omitted, such as nutrition), and by doing more on-ward debugging. In terms of technology, the biggest challenges are dealing with incomplete input data and generating good narrative texts. Data entered manually will always have omissions and mistakes and dealing with these robustly is a major challenge for any medical data-to-text system. Generating good narratives which include causal links and make the big picture clear is also a key data-to-text challenge; indeed, one could argue that such narrative aspects are perhaps the primary benefits of textual summaries over tabular/visual presentations. BT-Nurse's medical knowledge base also needs to be expanded.

Data-to-text technology is very new, and systems have been developed in many areas, including weather forecasts,^{9–13} financial and statistical information,^{14–16} and engineering.¹⁷ Recent attempts to apply data-to-text in medical contexts^{18,19} have used input which is simple compared with BT-Nurse; these are akin to early automated interpretation and report generation systems for personality assessment based on questionnaire

Table 1 Nurses' views of the BT-Nurse summaries (% of trials)

	Understandability			Accuracy			Helpfulness		
	Agree	Neutral	Disagree	Agree	Neutral	Disagree	Agree	Neutral	Disagree
Incoming	92.4	7.6	0	73.9	23.9	2.2	56.5	35.9	7.6
Outgoing	87.7	8.2	4.1	65.8	24.7	9.6	61.6	30.1	8.2
Overall	90.3	7.9	1.8	70.3	24.2	5.5	58.8	33.3	7.9

responses.²⁰ We are not aware of any previous medical system which is as ambitious as BT-Nurse in the amount and diversity of data summarized.

BT-Nurse has shown that it is possible to use data-to-text technology to generate useful and helpful summaries of nursing shifts from a complex, state-of-the-art patient information system holding a large amount of heterogeneous data. Of course, BT-Nurse is a proof of concept, and would require considerable engineering effort before it could be realistically deployed, or indeed evaluated in a clinical trial which measured patient outcome instead of nurse's perceptions. In particular, report accuracy needs to be higher, and comparable to the overall accuracy of the information in the patient record system. However, our evaluation, which involved a deployment of the system within its target environment and running on live, previously unseen data, shows that data-to-text systems can generate shift summaries from clinical data extracted from an electronic patient record.

Funding The UK Engineering and Physical Sciences Research Council (EPSRC) funded the BabyTalk project with grants to the University of Aberdeen (EP/D049520/1) and the University of Edinburgh (EP/D05057X/1).

Competing interests None.

Provenance and peer review Not commissioned; externally peer reviewed.

REFERENCES

1. Law AS, Freer Y, Hunter J, *et al.* A comparison of graphical and textual presentations of time series data to support medical decision making in the neonatal intensive care unit. *J Clin Monit Comput* 2005;**19**:183–94.
2. van der Meulen M, Logie RH, Freer Y, *et al.* When a graph is poorer than 100 words: a comparison of computerised natural language generation, human generated descriptions and graphical displays in neonatal intensive care. *Appl Cognit Psychol* 2008;**24**:77–89.
3. Gatt A, Portet F, Reiter E, *et al.* From data to text in the neonatal intensive care unit: using NLG technology for decision support and information management. *AI Communications* 2009;**22**:153–86.
4. Portet F, Reiter E, Gatt A, *et al.* Automatic generation of textual summaries from neonatal intensive care data. *Artif Intell* 2009;**173**:789–816.
5. Reiter E, Dale R. *Building Natural Language Generation Systems*. Cambridge: Cambridge University Press, 2000.
6. Reiter E. An architecture for data-to-text systems. In: Buseman S, ed. *Proceedings of the 11th European Workshop on Natural Language Generation*. Schloss Dagstuhl, Germany: ACL, 2007:97–104.
7. Gatt A, Reiter E. SimpleNLG: a realisation engine for practical applications. In: Krahmer E, Theune M, eds. *Twelfth European Workshop on Natural Language Generation: Proceedings of the Workshop*. Athens, Greece: ACL, 2009:90–4.
8. Artstein R, Poesio M. Inter-coder agreement for computational linguistics. *Comput Ling* 2008;**34**:555–96.
9. Goldberg E, Driedger N, Kittredge RI. Using natural language processing to produce weather forecasts. *IEEE Expert* 1994;**9**:45–53.
10. Coch J. Multimeteo: multilingual generation of weather forecasts. *ELRA Newsletter* 1998;**3**:13–15.
11. Reiter E, Sripada S, Hunter J, *et al.* Choosing words in computer-generated weather forecasts. *Artif Intell* 2005;**167**:137–69.
12. Turner R, Sripada S, Reiter E, *et al.* Selecting the content of textual descriptions of geographically located events in spatio-temporal weather data. In: Ellis R, Allen T, Petridis M, eds. *Applications and Innovations in Intelligent Systems XV. Proceedings of the 27th SGAI International Conference on Innovative Techniques and Applications of Artificial Intelligence*. Berlin: Springer, 2007:75–88.
13. Belz A. Automatic generation of weather forecast texts using comprehensive probabilistic generation space models. *Nat Lang Eng* 2007;**14**:431–55.
14. Kukich K. Design of a knowledge-based report generator. In: Marcus M, ed. *Proceedings of the 21st Annual Meeting of the Association for Computational Linguistics*. Cambridge, MA: ACL, 1983:145–50.
15. Iordanskaja L, Kim M, Kittredge R, *et al.* Generation of extended bilingual statistical reports. In: *Proceedings of the 15th International Conference on Computational Linguistics*. Nantes: ICCL, 1992:1019–23.
16. Ferres L, Parush A, Roberts S, *et al.* Helping people with visual impairments gain access to graphical information through natural language: the iGraph system. In: Meisenberger K, Klaus J, Zagler W, *et al.*, eds. *Computers Helping People with Special Needs: 10th International Conference ICCHP 2006*. Berlin: Springer, 2006:1122–30.
17. Yu J, Reiter E, Hunter J, *et al.* Choosing the content of textual summaries of large time-series data sets. *Nat Lang Eng* 2007;**13**:25–49.
18. Dalal M, Feiner S, McKeown K, *et al.* MAGIC: an experimental system for generating multimedia briefings about post-bypass patient status. In: Cimino JJ, ed. *Proceedings of the AMIA Annual Fall Symposium*. Washington, DC: Henley and Belfus Inc, 1996:684–8.
19. Harris M. Building a large-scale commercial NLG system for an EMR. In: White M, Nakatsu C, McDonald D, eds. *INLG 2008: Fifth International Natural Language Generation Conference*. Salt Fork, OH: ACL, 2008:157–60.
20. Moreland KL. Validation of computer-based test interpretations: problems and prospects. *J Consult Clin Psychol* 1985;**53**:816–25.