# Objective Comparison Using Guideline-based Query of Conventional Radiological Reports and Structured Reports

MÁTÉ E. MAROS[1], RALF WENZ[2], ALEX FÖRSTER[1], MATTHIAS F. FROELICH[3], CHRISTOPH GRODEN[1], WIELAND H. SOMMER[3,4], STEFAN O. SCHÖNBERG[5], THOMAS HENZLER[5] and HOLGER WENZ[1]

[1]*Department of Neuroradiology, and* [5]*Institute of Clinical Radiology and Nuclear Medicine, University Medical Center Mannheim, Medical Faculty Mannheim, Heidelberg University, Mannheim, Germany;* [2]*Department of Life Sciences, Faculty of Natural Sciences, Imperial College London, London, U.K.;* [3]*Smart-Radiology, Smart Reporting GmbH, Munich, Germany;* [4]*Institute for Clinical Radiology, Ludwig Maximilian University Hospital, Munich, Germany*

**Abstract.** *Background: This feasibility study of text-mining-based scoring algorithm provides an objective comparison of structured reports (SR) and conventional free-text reports (cFTR) by means of guideline-based key terms. Furthermore, an open-source online version of this ranking algorithm was provided with multilingual text-retrieval pipeline, customizable query and real-time-scoring. Materials and Methods: Twenty-five patients with suspected stroke and magnetic resonance imaging were re-assessed by two independent/blinded readers [inexperienced: 3 years; experienced >6 years/Board-certified). SR and cFTR were compared with guideline-query using the cosine similarity score (CSS) and Wilcoxon signed-rank test. Results: All pathological findings (18/18) were identified by SR and cFTR. The impressions section of the SRs of the inexperienced reader had the highest median (0.145) and maximal (0.214) CSS and were rated significantly higher (p=2.21×10$^{-5}$ and p=1.4×10$^{-4}$, respectively) than cFTR (median=0.102). CSS was robust to variations of query. Conclusion: Objective guideline-based comparison of SRs and cFTRs using the CSS is feasible and provides a scalable quality measure that can facilitate the adoption of structured reports in all fields of radiology.*

This article is freely accessible online.

*Correspondence to:* Holger Wenz, MD, Department of Neuroradiology, Universitätsmedizin Mannheim, Theodor-Kutzer-Ufer 1-3, 68137 Mannheim, Germany. Tel: +49 6213832443, Fax: +49 6213832165, e-mail: Holger.Wenz@umm.de

Structured reporting (SR) represents a new direction in communicating radiological reports to clinicians in a way that is clear and uniform (1). However, follow-up studies on SR have shown mixed results regarding adaptation and adherence (2-4). In particular, neuroradiological residents most commonly complain that SRs are overly constraining and time-consuming (5). In spite of that, there is an increasing tendency to use online solutions to generate SR templates (6). For instance, "CT brain" and "MR brain" were the third and fifth most frequently viewed SR templates in the Radiological Society of North America (RSNA) online library (7). SR templates can serve as core frameworks, but because of heterogeneous representations of diseases and co-morbidities, there is a substantial need to allow for customization of each report – primarily in the form of additional free text. Consequently, the distinction between conventional free-text (narrative) reports (cFTR) and SR becomes even more blurred, which makes their objective comparison even harder. This raises two further questions: Are SRs indeed 'better' than cFTRs? How can cFTRs be compared with SRs in an objective, fast and scalable way if SRs are often a mixture of structured and free (narrative) text?

Firstly, there is a need to define how to objectively quantify 'better'. Commonly the method of choice for such quality assessments is either an expert opinion-based rating (5, 8) or a survey-based evaluation by physicians (4). Although both methods provide valuable insight, they are very laborious and time-consuming for experts and do not scale well for large datasets (4, 8). Therefore, we aimed to implement an approach using the adherence to imaging and clinical guidelines as a quality measure.

Secondly, in order to objectively compare cFTR with SR, we suggest an evaluation based on widely used text-mining technique using the 'bag of words' representation and cosine similarity (9, 10). The cosine similarity provides the core of many information retrieval systems (11, 12). Reports can be

ranked according to a pre-specified query of clinically relevant information based on imaging guidelines.

Addressing the questions raised above, we present a proof-of-concept study to provide a reproducible, objective, and scalable comparison of SR and cFTR by means of guideline-based key terms and the cosine similarity measure. Furthermore, we developed and provide a free online version (http://www.radreport-query.com) of this ranking algorithm with multilingual text retrieval pipeline, customizable query and real-time scoring.

## Materials and Methods

*Study cohort selection*. Data were retrieved for all patients (n=780) who underwent cranial computed tomography (CT) or magnetic resonance imaging (MRI) due to suspected cerebral ischemia between Jan 1st 2014 and Jan 1st 2015 from our institutional radiological database (Syngo, Siemens Healthcare Sector, Forchheim, Germany). To appropriately model a heterogeneous daily clinical case collective, a random sample (n=30) of that patient population was generated. This sample was screened for inclusion criteria of available cranial MRI with institutional standard stroke protocol including diffusion-weighted imaging in transversal and coronal planes, time-of-flight angiography, and transversal fluid-attenuated inversion recovery, T2*, T1-, and T2-weighted images. Additional MRI sequences were allowed.

This study was approved by the local Institutional Review Board (approval number: 2017-825-MA) and was therefore performed in accordance with the ethical standards laid down in the 1964 Declaration of Helsinki and its later amendments. Patient consent was waived for this analysis by the local Institutional Review Board due to its retrospective nature.

*Workflow*. We constructed a structured reporting template (in German) specifically designed for suspected stroke using an online reporting platform (www.smart-radiology.com; Smart Reporting GmbH, Munich, Germany) (13). The template included predefined text blocks and boxes for optional free text inputs. Using this online reporting tool, the study cohort was re-assessed by two independent readers [inexperienced: 3 years (M.E.M.); experienced: >6 years/Board-certified (H.W.)] blinded to all clinical (other than gender and age), previous imaging and laboratory data (Figure 1). To minimize differences from cFTRs, the assessment process was performed similarly to the regular working-hour processes, albeit in an isolated setting (without disturbances such as phone calls *etc*.). Next, all the corresponding cFTRs, signed-off by attending or senior radiologists, were retrieved from the local database (Syngo, Siemens Healthcare Sector, Forchheim, Germany).

A query (9) was then defined to cover the most important clinical aspects of stroke diagnostics and subsequent therapy prerequisites or its contraindications (Figure 2) according to the Joint Statement for Imaging Recommendations for Acute Stroke and Transient Ischemic Attack Patients by the American Society of Neuroradiology, the American College of Radiology, and the Society of NeuroInterventional Surgery (14), as well as the guidelines of the German Society for Neurology (DGN) for the Diagnostics of Acute Cerebrovascular Diseases AWMF-030/117 (15). It is paramount that key features are assessed, including the presence of i) hemorrhage, ii) ischemic core with its vessel territory,

iii) intracranial vessel status (*i.e.* stenosis, occlusion), iv) collateralization, v) tissue viability (perfusion, mismatch), vi) carotid vessel status or extracranial atherosclerosis, and vii) intracranial pressure signs. Different semantics were accounted for in order to address the fact that each radiologist has their own set of preferred terminologies.

*Statistical analyses and text mining*. All analyses were descriptive and performed using the R language and environment for statistical programming (R version 3.3.2, R Core Team 2016, Vienna Austria) (16) within the RStudio IDE (Version 1.0.136 (17)). Both SR and cFTR underwent pre-processing using the tm package (12) This included re-encoding of special characters (*e.g.* ä, ü, ß) using stringi package (18), lower case transformation, removal of punctuation, and (German) stop words (SnowballC package) (19).

Consequently, both reports and the query were viewed as a vector of their words, which enabled us to interpret the query as a short document (9). The cosine measure can then be utilized to quantify the similarity between the query and a given report (20). The cosine similarity score (CSS) was then used to objectively rank and compare cFTR and SR (9). CSS provides a continuous value of between 0 and 1 according to total disagreement and complete similarity, respectively. An example of CSS calculation is presented in Figure 3. A uniform weighting scheme for the query terms was used (11). In order to account for the fact that certain key terms can occur multiple times within a single report, a binary cut-off adjusted cosine similarity was used. This was done in a 'yes\no' fashion, meaning that a term was considered either present or not, irrespective of the number of occurrences (9, 10, 20). The findings and impression sections of both SR and cFTR were compared separately with the same guideline query. To assess the robustness of CSS distributions, the variations of the query were also tested as sensitivity analyses (data not shown). The Wilcoxon signed-rank test for paired data was used to compare the CSS distributions between users and report types (21). Associations were measured using Spearman's rank correlation ($rho_{Sp}$) and interrater agreement was quantified with Cohen's Kappa statistic. The Mann–Whitney–Wilcoxon test was used to compare distributions of the random sampling (21). Plots were generated using the ggplot2 package (22). *p*-Values of less than 0.05 were considered significant. In the case of multiple testing, the alpha-level was adjusted using the conservative Bonferroni correction.

## Results

*Study cohort and sampling*. To assess the random sampling, the age distributions of male (7/25; median=73 years, range=26-81 years) and female (18/25; median=73.5 years, range=41-93 years) cases were compared, and found to be similar ($p_{WMW}$=0.505). In half of the cases (13/25; 52%), additional sequences were performed (*e.g.* perfusion weighted imaging, MR neck or cerebral contrast-enhanced angiography, T2 constructive interference in steady state or contrast-enhanced T1). Both readers customized SRs with additional free-text in 80% (20/22) and 88%, (22/25), of cases respectively. All primary pathological findings (*e.g.* ischemia or intracerebral hemorrhage) and all incidental findings (*e.g.* meningioma, aneurysms) were identified
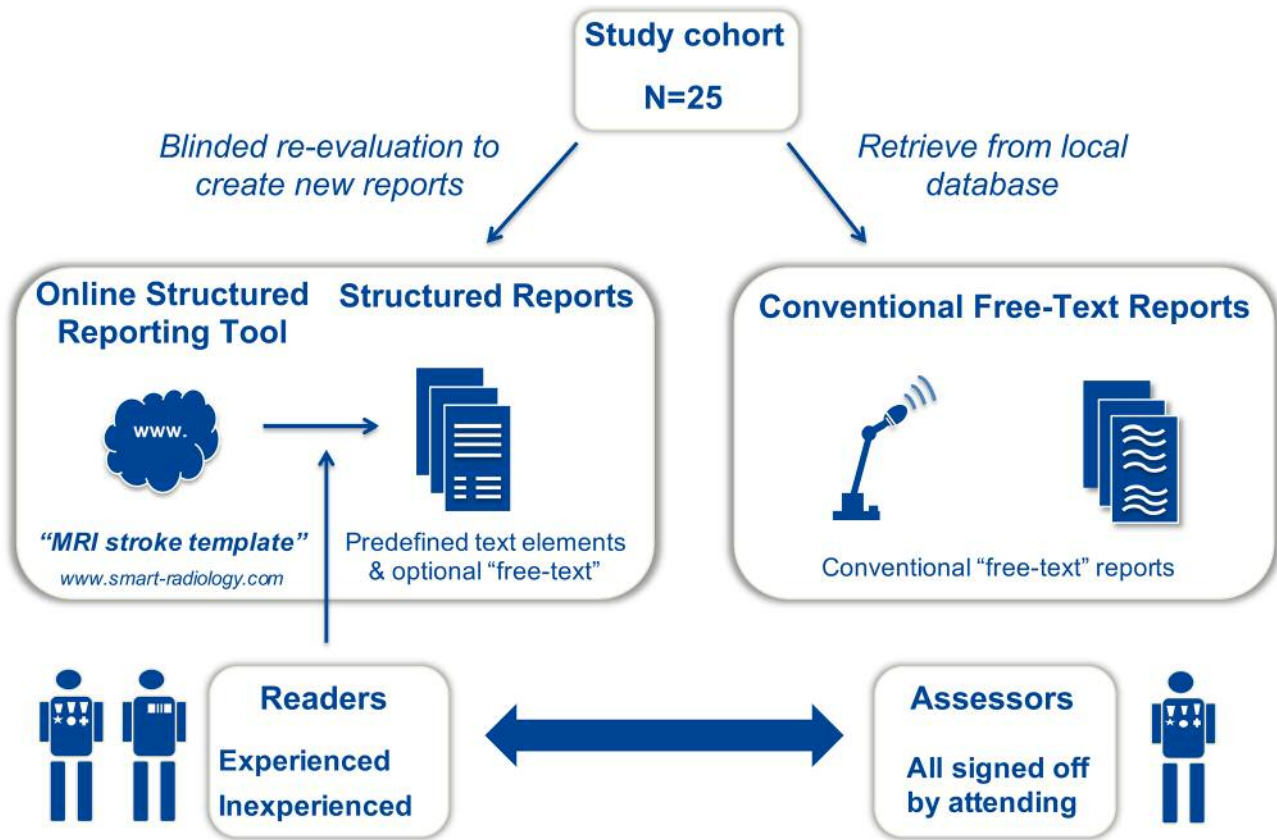
Figure 1. *Re-assessment of the study cohort. The retrospective random cohort of 25 patients with suspected stroke was re-assessed using an online reporting template (www.smart-radiology) by two independent blinded readers with 3 and >6 years of experience, respectively. These new structured reports were compared to the corresponding conventional free-text (narrative) reports signed-off by attending or senior radiologists.*

(18/25; 72%) and showed a perfect correlation ($rho_{Sp}$=1) and interrater agreement (Cohen's Kappa=1) between both readers of SRs and cFTRs. In 12 out of 25 cases (48%), the initial cFTR was written by residents in training with <5 years of experience. Eight different senior radiologists signed off cFTRs; just two being responsible for 44% (seven and four, respectively) of all reports, which had to be accounted for through inclusion of synonyms during the query definition.

*Query similarity.* The CSS for the key terms of the query and the impression sections of reports are presented in Figure 4A. The SR of the inexperienced reader had the highest median CSS and maximum values. The SR of the experienced resident had a comparable median score ($p_{Wilcox.SRT}$=0.240), with the smallest variance. cFTRs had significantly lower median CSS than the SR of inexperienced ($p_{Wilcox.SRT}$=2.21×10$^{-5}$) and experienced readers ($p_{Wilcox.SRT}$=1.40×10$^{-4}$) and they had the largest CSS variance.

The CSS distributions of the findings sections of reports (Figure 4B) revealed an even more significant positive association with increased key term content of SRs compared to cFTRs ($p_{Wilcox.SRT}$=5.96×10$^{-8}$). Corresponding to the results for impressions sections, the SR findings of the experienced reader had the most compact terminology with smallest CSS variance; the SR findings of the inexperienced reader had the highest maximum CSS, while cFTR had the lowest median and maximum scores. Although the CSS of the inexperienced reader showed higher variance, it had a similar distribution ($p_{Wilcox.SRT}$=0.0383) to the scores of the experienced reader when adjusted for multiple testing ($p_{alpha*-Bonferroni}$=0.0167). Details are presented in Figure 4B.

## Discussion

SR templates are thought to promote adherence to guidelines and increase critical information content (1, 3). However, the need to allow for free text customization of SRs blurs the margin between SRs and cFTRs (5).
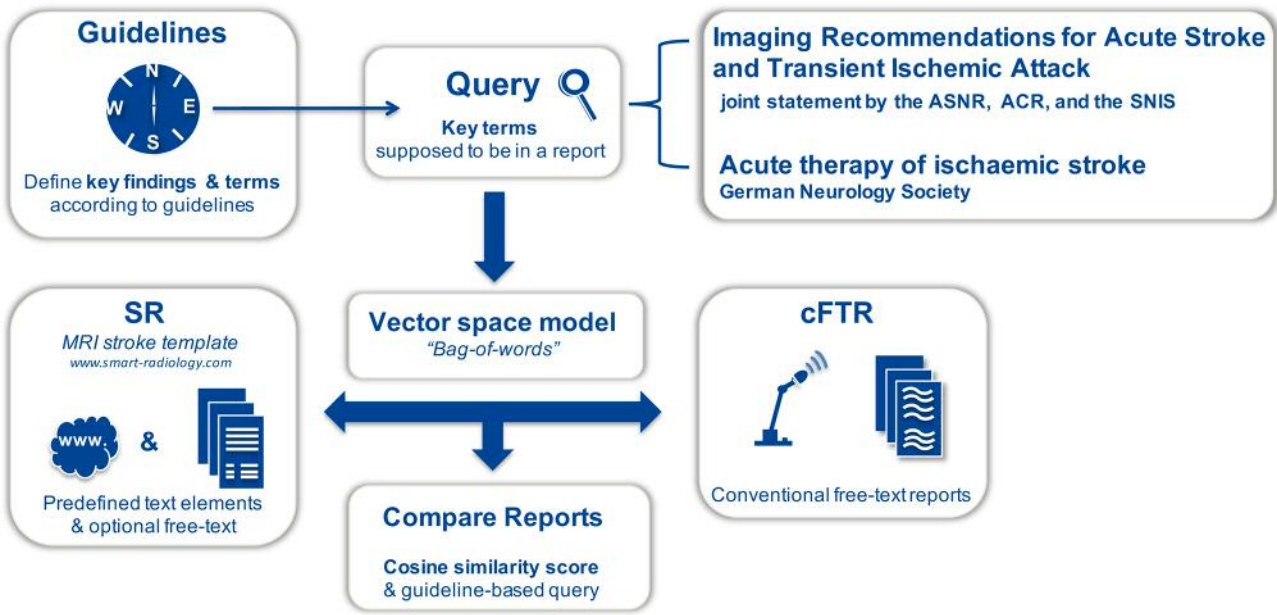
845

Figure 2. *Representation of workflow and query definition. We defined a query to cover the critical radiological aspects of stroke diagnostics according to the Joint Statement for Imaging Recommendations for Acute Stroke and Transient Ischemic Attack Patients by the American Society of Neuroradiology (ASNR), the American College of Radiology (ACR), and the Society of NeuroInterventional Surgery (SNIS) (14), plus guidelines of the German Society for Neurology for diagnostics of acute cerebrovascular diseases (15). The structured reports (SR) and the corresponding retrieved conventional free-text reports (cFTR) underwent a text-processing pipeline and both their finding and impression sections were compared to the guideline query using the cosine similarity score.*
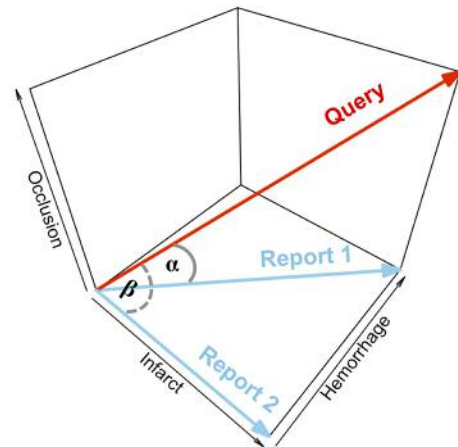


Figure 3. *Example of cosine similarity-based text ranking. The figure shows two hypothetical report impressions and a query. The query of three key terms represents a three-dimensional vector space (i.e. one axis for each term) in which the reports are scored (see column "vector representation"). The graph shows this 3D vector space and the corresponding cosine similarity scores of the reports and query measured as cos $\alpha$ and cos $\beta$.*

Firstly, we assessed how to objectively quantify the quality of both report types. For this, the quality of SRs and cFTRs was measured using guideline-based key information content as a query (3).

Secondly, we presented a task-specific application of a widely used information retrieval method (9, 10) to compare and rank reports using this query in a fast, scalable and reproducible fashion (12). Thereby, we showed that CSS is
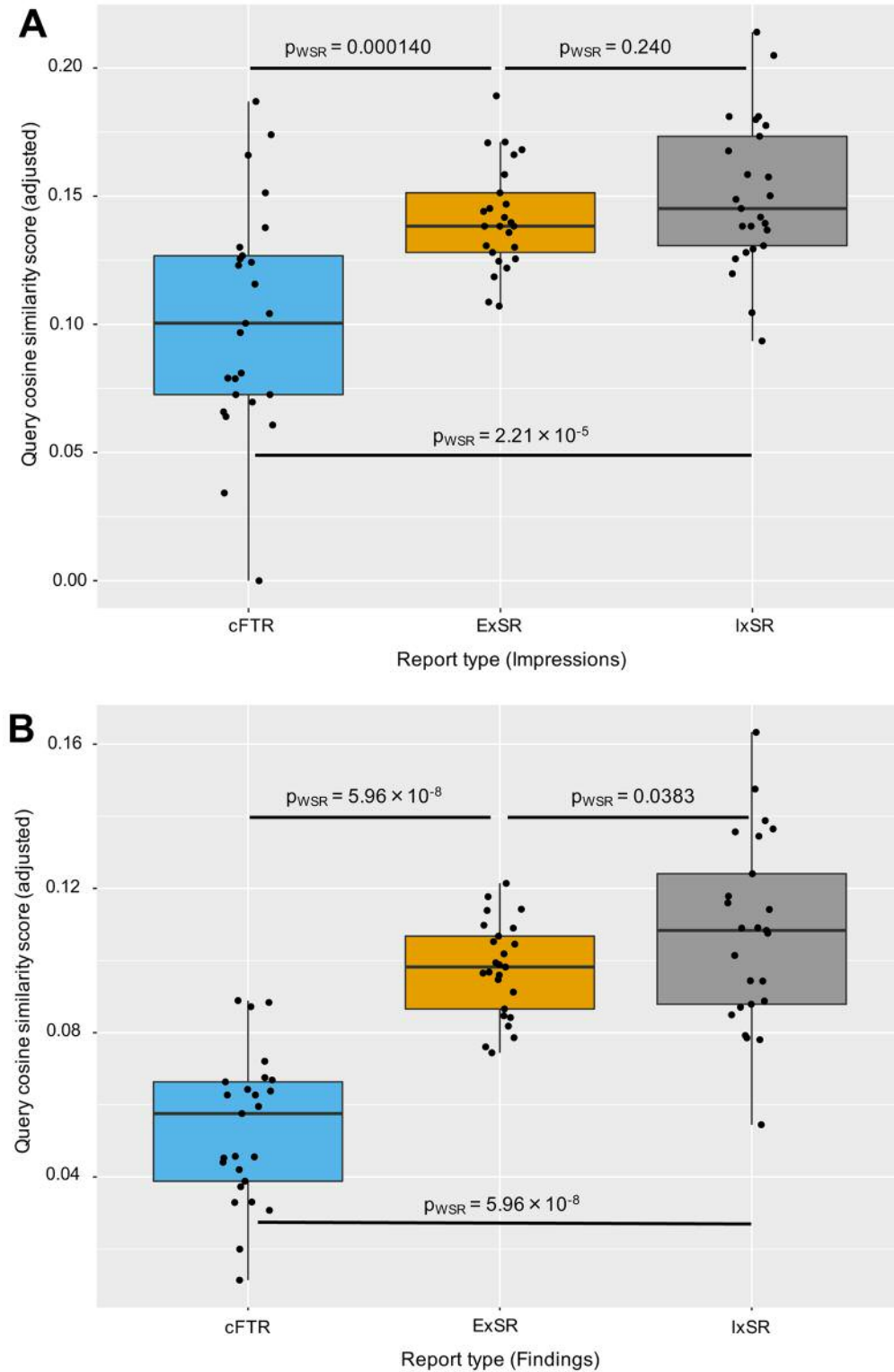
Figure 4. *Boxplot of query cosine similarity scores of report impressions and findings. The figure shows the guideline query-based cosine similarity score (CSS) distributions of conventional free text reports (cFTR) and structured reports (SR) of experienced (ExSR) and inexperienced readers (IxSR) for the impressions (A) and findings (B) sections. The impressions of both readers for SRs received significantly higher CSS than did cFTRs. Evaluating the findings section (B) showed an even more pronounced CSS difference in favor of SRs. There was no significant difference between the CSS of readers after adjusting for multiple testing using Bonferroni-correction ($p_{alpha*}$=0.0167).*

indeed a feasible and robust method to reliably compare reports, which we hope will aid other researchers in the future to assess and increase reproducibility of their work quickly and easily.

Information retrieval and text-mining methods have been extensively used to extract key information from unstructured radiological reports. Natural language processing methods can be used to find quality indicators of radiological practice (23), automatically analyze content (24), classify free-text reports (25) or assess the presence of sufficient clinical recommendations in radiological reports (26). Such methods successfully use more complex language models and retrieve contextual information, but they need to be tailored for specific conditions and require dedicated sets of examples for training and optimization (23).

In contrast, the query-based retrieval approach can be more flexible, with fast, expert-driven customization, while being more suitable for smaller datasets (23). Therefore, we used the latter approach that uses key terms taken from guidelines to build a query-based quality measure. As most clinicians focus primarily on the impression section of radiological reports (27), we solely relied on this section when constructing the query. In order to avoid introducing positive bias for SRs, through imbalanced query terms, after inspecting the document term frequencies, we included synonymous terms present in both report types.

Using SR templates favored the inexperienced reader and resulted in the highest median guideline CSS. The experienced reader had the lowest variance, suggesting most consistent radiological terminology. The CSS of the experienced and inexperienced readers using SR was similar. Additionally, SRs of both readers scored significantly higher than cFTRs. In a similar setting of cranial MRI of suspected stroke, Johnson *et al*. found that using SR templates to create reports did not seem to improve report clarity based on expert evaluation (4, 5). However, they noted that intrinsic report quality has not been tested (5) and SR should include definitions of key terms (4). Schwartz *et al*. compared the content, clarity, and clinical usefulness of cFTR and SR of body CT scans using clinical and radiological expert-based evaluation. Similarly to our findings, they reported that SRs scored significantly higher for content and greater clarity than cFTRs (8).

A limitation of query-based retrieval is that the composition of the query can substantially influence the results of both radiological imaging retrieval (28) and rankings (29). Therefore, we had to allow for a variety of terminologies because multiple radiologists signed-off on cFTRs. This resulted in a long query vector (~120 key words). However, the advantage of the applied vector space model is that the query can be interpreted as a short document (9). Thus, the CSS can be used to quantify the similarity between the query and reports (20). CSS compensates for the effect of report and query length (9, 10) and is considered a standard means of quantifying document similarity (12, 20). As a result of a long query, the CSS is generally small (*i.e.* close alignments in high-dimensional spaces become increasingly rare) but it enables a fine-scale evaluation of reports (9, 10, 12). Nevertheless, this simplified 'bag of words' considers neither the ordering of the words (9, 10) nor other linguistic features (12). For improving the consistency and robustness of CSS, we used binary-adjusted scores *i.e.* a query term was either present or not (10), regardless of the number of its occurrences (11). This accounts for the fact that multiple appearances of the same key term can slightly modify the CSS (*i.e.* report's position in the vector space) (9). Additionally, query terms were weighted equally (10, 11), which is crucial for using the query as an objective quality measure.

These shortcomings and the role of proper query definition underline the importance of the RadLex® initiative of the RSNA (30) to unify reporting language of radiologists and define a lexicon of preferred terms. Using RadLex®-based query terms could further improve the reliability of report quality rankings (31-33).

In order to facilitate reproducible, objective, and scalable comparisons of SRs and cFTRs, we created a freely available open-source version (radreport-query.org) of this ranking algorithm with multilingual text retrieval pipeline, customizable query and real-time cosine similarity scoring, with flexible plotting of the results.

## Conclusion

Our approach of using SR templates with additional free-text input provides three major advantages: i) a flexible and reliable tool with higher built-in adherence to guidelines; ii) increased report quality for both experienced and junior readers, thereby reducing the burden on attending radiologist when signing off reports, iii) cosine similarity is a robust and objective way to compare SR with cFTR using guideline-based query, that can facilitate the adoption of structured reports in all fields of radiology.

## Conflicts of Interest

The Authors declare that they have no conflict of interest in regard to this study.

## References

1 Morgan TA, Helibrun ME and Kahn CE Jr.: Reporting initiative of the Radiological Society of North America: progress and new directions. Radiology *273(3)*: 642-645, 2014.
2 Weiss DL and Langlotz CP: Structured reporting: Patient care enhancement or productivity nightmare? Radiology *249(3)*: 739-747, 2008.

3 Kahn CE Jr., Heilbrun ME and Applegate KE: From guidelines to practice: How reporting templates promote the use of radiology practice guidelines. J Am Coll Radiol *10(4)*: 268-273, 2013.

4 Johnson AJ, Chen MY, Zapadka ME, Lyders EM and Littenberg B: Radiology report clarity: A cohort study of structured reporting compared with conventional dictation. J Am Coll Radiol *7(7)*: 501-506, 2010.

5 Johnson AJ, Chen MY, Swan JS, Applegate KE and Littenberg B: Cohort study of structured reporting compared with conventional dictation. Radiology *253(1)*: 74-80, 2009.

6 Faggioni L, Coppola F, Ferrari R, Neri E and Regge D: Usage of structured reporting in radiological practice: Results from an Italian online survey. Eur Radiol *27(5)*: 1934-1943, 2016.

7 RSNA TRSoNA: Radiology Reporting Initiative; Report Template Library, (http://www.Radreport.Org/metrics.Php - accessed 03/16/2017).

8 Schwartz LH, Panicek DM, Berk AR, Li Y and Hricak H: Improving communication of diagnostic radiology findings through structured reporting. Radiology *260(1)*: 174-181, 2011.

9 Manning CD, Raghavan P and Schütze H: Introduction to Information Retrieval. Cambridge University Press, Cambridge, 2008.

10 Feldman R and Sanger J: The Text-Mining Handbook: Advanced Approaches in Analyzing Unstructured Data. Cambridge University Press, Cambridge, UK, 2007.

11 Salton G and Buckley C: Term-weighting approaches in automatic text retrieval. Info Process Manag *24(5)*: 513-523, 1988.

12 Meyer D, Hornik K and Feinerer I: Text mining infrastructure in R. J Stat Software *25(5)*: 1-54, 2008.

13 Smart Reporting GmbH: Smart-radiology, http://www.Smart-radiology. Com. 2016.

14 Wintermark M, Sanelli PC, Albers GW, Bello JA, Derdeyn CP, Hetts SW, Johnson MH, Kidwell CS, Lev MH, Liebeskind DS, Rowley HA, Schaefer PW, Sunshine JL, Zaharchuk G, Meltzer CC, American Society of Neuroradiology, American College of Radiology and Society of Neurointerventional Surgery: Imaging recommendations for acute stroke and transient ischemic attack patients: A joint statement by the American Society of Neuroradiology, American College of Radiology and Society of Neurointerventional Surgery. J Am Coll Radiol *10(11)*: 828-832, 2013.

15 Deutsche Gesellschaft für Neurologie: Diagnostic of acute cerebrovascular diseases awmf-030/117 2016. http:// docplayer.org/25127262-Diagnostik-zerebrovaskulaerer-erkrankungen.html [Last accessed 23 April 2018]

16 R Core Team: R: A language and environment for statistical computing. R Foundation for Statistical Computing, 2016. http://www.r-project.org/ [Last accessed 23 April 2018]

17 RStudio Team: Rstudio: Integrated development envinronment for R. RStudio Inc, Boston, MA, 2015. http://www.Rstudio.com

18 Gagolewski M and Tartanus B: R Package Stringi: Character String Processing Facilities, 2015. https://cran.r-project.org/web/packages/stringi/stringi.pdf [Last accessed 23 April 2018]

19 Bouchet-Valat M: Snowballc: Snowball Stemmers Based on the C Libstemmer UTF-8 library. R package version 05 1, 2014. https://cran.r-project.org/web/packages/SnowballC/SnowballC.pdf [Last accessed 23 April 2018]

20 Salton G and McGill MJ: Introduction to Modern Information Retrieval, McGraw-Hill, Inc. New York, NY, USA, 1986.

21 Hollander M, Wolfe DA and Chicken E: Nonparametric Statistical Methods. John Wiley & Sons, 2013.

22 Wickham H: Ggplot2: Elegant Graphics for Data Analysis. Springer, 2016.

23 Pons E, Braun LMM, Hunink MGM and Kors JA: Natural language processing in radiology: A systematic review. Radiology *279(2)*: 329-343, 2016.

24 Hong Y and Kahn CE Jr.: Content analysis of reporting templates and free-text radiology reports. J Digit Imaging *26(5)*: 843-849, 2013.

25 Dreyer KJ, Kalra MK, Maher MM, Hurier AM, Asfaw BA, Schultz T, Halpern EF and Thrall JH: Application of recently developed computer algorithm for automatic classification of unstructured radiology reports: Validation study 1. Radiology *234(2)*: 323-329, 2005.

26 Dang PA, Kalra MK, Blake MA, Schultz TJ, Stout M, Lemay PR, Freshman DJ, Halpern EF and Dreyer KJ: Natural language processing using online analytic processing for assessing recommendations in radiology reports. J Am Coll Radiol *5(3)*: 197-204, 2008.

27 Gershanik EF, Lacson R and Khorasani R: Critical finding capture in the impression section of radiology reports. AMIA Annu Symp Proc *2011*: 465-469, 2011.

28 Gerstmair A, Daumke P, Simon K, Langer M and Kotter E: Intelligent image retrieval based on radiology reports. Eur Radiol *22(12)*: 2750-2758, 2012.

29 Bianchini M, Gori M and Scarselli F: Inside pagerank. ACM Trans Internet Technol *5(1)*: 92-128, 2005.

30 RSNA TRSoNA: Radlex, https://www.Rsna.Org/radlex.Aspx. RSNA Informatics, 2017.

31 Bozkurt S and Kahn CE Jr.: An open-standards grammar for outline-style radiology report templates. J Digit Imaging *25(3)*: 359-364, 2012.

32 Hong Y, Zhang J, Heilbrun ME and Kahn CE Jr.: Analysis of radlex coverage and term co-occurrence in radiology reporting templates. J Digit Imaging *25(1)*: 56-62, 2012.

33 Kahn CE Jr., Genereaux B and Langlotz CP: Conversion of radiology reporting templates to the mrrt standard. J Digit Imaging *28(5)*: 528-536, 2015.