Clinical Studies

# Development of a natural language processing algorithm for the detection of spinal metastasis based on magnetic resonance imaging reports

Evan Mostafa, MD [a],*, Aaron Hui, BS [b], Boudewijn Aasman, BS [b], Kamlesh Chowdary, BS [b], Kyle Mani, BS [b], Edward Mardakhaev, MD [b], Richard Zampolin, MD [b], Einat Blumfield, MD [b], Jesse Berman, MD [b], Rafael De La Garza Ramos, MD [c], Mitchell Fourman, MD [a], Reza Yassari, MD [c], Ananth Eleswarapu, MD [a], Parsa Mirhaji, PhD [b]

[a] Department of Orthopaedic Surgery, Montefiore Medical Center, 111 E 210th St, Bronx, NY, 10467, United States
[b] Albert Einstein College of Medicine, 1300 Morris Park Ave, Bronx, 10461, NY, United States
[c] Department of Neurological Surgery, Montefiore Medical Center, 111 E 210th St, Bronx, NY, 10467, United States

## ARTICLE INFO

## ABSTRACT

*Background:* Metastasis to the spinal column is a common complication of malignancy, potentially causing pain and neurologic injury. An automated system to identify and refer patients with spinal metastases can help overcome barriers to timely treatment. We describe the training, optimization and validation of a natural language processing algorithm to identify the presence of vertebral metastasis and metastatic epidural cord compression (MECC) from radiology reports of spinal MRIs.

*Methods:* Reports from patients with spine MRI studies performed between January 1, 2008 and April 14, 2019 were reviewed by a team of radiologists to assess for the presence of cancer and generate a labeled dataset for model training. Using regular expression, impression sections were extracted from the reports and converted to all lower-case letters with all nonalphabetic characters removed. The reports were then tokenized and vectorized using the doc2vec algorithm. These were then used to train a neural network to predict the likelihood of spinal tumor or MECC. For each report, the model provided a number from 0 to 1 corresponding to its impression. We then obtained 111 MRI reports from outside the test set, 92 manually labeled negative and 19 with MECC to test the model's performance.

*Results:* About 37,579 radiology reports were reviewed. About 36,676 were labeled negative, and 903 with MECC. We chose a cutoff of 0.02 as a positive result to optimize for a low false negative rate. At this threshold we found a 100% sensitivity rate with a low false positive rate of 2.2%.

*Conclusions:* The NLP model described predicts the presence of spinal tumor and MECC in spine MRI reports with high accuracy. We plan to implement the algorithm into our EMR to allow for faster referral of these patients to appropriate specialists, allowing for reduced morbidity and increased survival.

## Introduction

Spinal column metastases are a common complication of malignancy, with the potential to cause pain, instability, and neurologic injury [1]. Timely diagnosis and treatment are associated with improved neurologic status, lower complication rates and improved survival [1–5]. Unfortunately, delays in the diagnosis, referral and treatment of spinal metastases are common for a variety of reasons [6]. Clinical care pathways designed to ensure the urgent referral of patients with spinal metastases to a multidisciplinary spinal oncology team have been shown to decrease time to treatment [7]. However, such pathways face numerous barriers to implementation, including a lack of clinician awareness, financial cost, and insufficient staff available to identify patients for inclusion in the pathway [8]. An automated system for

identifying and appropriately referring patients with spinal metastases can help to overcome these barriers and allow for the widespread implementation of spinal oncology care pathways even in resource-constrained environments. Natural language processing (NLP) involves computer analysis to extract and digitize information from narrative text. NLP has been used in a variety of medical applications, including the identification of complications from operative reports and the identification of pathology on radiology reports of X-ray, CT and MRI [9–13]. Here we describe the training, optimization and validation of an NLP algorithm to identify vertebral metastasis and metastatic epidural cord compression (MECC) from radiology reports of spine MRIs.

## Materials and methods

### Development of the training dataset

Ethical approval for this study was obtained from our institutional review board. We used radiology procedure codes to identify 37,579 radiology reports from 25,469 patients who had cervical, thoracic and/or lumbar MRI studies performed from January 1, 2008 to April 14, 2019. Radiology reports were labeled as "negative" if the patient did not have a diagnosis of cancer in the ICD-9 and ICD-10 codes entered in their electronic medical record. The remaining reports were then reviewed by a team of 5 annotators (3 board-certified neuroradiologists and 2 radiology residents) to assess for the presence of cancer. Image classification choices included "negative" (no pathology), or "positive" (presence of spinal tumor with or without spinal cord compression). We were thus able to generate a labeled dataset for model training.

### Development of the NLP model

Each radiology report was preprocessed using Gensim, a free open-source Python library program. Using regular expression (regex), the impression and interpretation sections were extracted from the reports. These were then further processed by conversion to all lower-case letters and removal of punctuation, digits and any other nonalphabetic characters (Fig. 1A and B). The reports were then tokenized and vectorized using the doc2vec algorithm [14]. Doc2vec is an open-source unsupervised learning algorithm which analyzes word sequences and creates a numeric vector to represent each document, allowing for tasks like clustering and classification. A distributed bag-of-words approach was used, in which the model attempts to predict words sampled randomly from the paragraph, based on the paragraph vector. Based on the error in prediction, the paragraph vectors are adjusted via backpropagation until the error is minimized. Once vectorized, the reports were then used to train a neural network to predict the likelihood of spinal tumor based on the MRI report. For each MRI report, the model provides a number from 0 to 1 corresponding to its impression of the likelihood of spinal cancer compression on the report.

### Testing and optimization of the model

We established a testing cohort consisting of 111 MRI reports from outside the training set, 92 of which were manually labeled as negative and 19 with spinal tumor. We used this to test the performance of the model. As the model outputs a number between 0 and 1 corresponding to the likelihood of tumor being present, we were able to adjust the cutoff for a "positive" result to optimize model performance for our clinical needs. Statistical analysis of model performance included calculation of the sensitivity, specificity, positive and negative predictive values, accuracy, and an F1 score. The F1 score is a harmonic mean of the sensitivity and positive predictive value (PPV) and is calculated as 2 x (sensitivity x PPV) / (sensitivity + PPV). These parameters were calculated at a variety of cutoffs for a "positive" result so that we could compare model performance at various thresholds. We also constructed a confusion ma-

**Table 1**

Calculated NLP model performance metrics at various cutoff values.

| Cutoff | Sensitivity | Specificity | PPV | NPV | F1 | Accuracy |
|---|---|---|---|---|---|---|
| 0.02 | 1 | 0.978 | 0.905 | 1 | 0.95 | 0.982 |
| 0.04 | 0.947 | 0.978 | 0.9 | 0.989 | 0.923 | 0.973 |
| 0.05 | 0.947 | 0.978 | 0.9 | 0.989 | 0.923 | 0.973 |
| 0.06 | 0.947 | 0.978 | 0.9 | 0.989 | 0.923 | 0.973 |
| 0.07 | 0.947 | 0.978 | 0.9 | 0.989 | 0.923 | 0.973 |
| 0.08 | 0.947 | 0.978 | 0.9 | 0.989 | 0.923 | 0.973 |

**Table 2**

Comparison of the NLP approach to untrained string search.

| Model | Sensitivity | Specificity | PPV | NPV | F1 | Accuracy |
|---|---|---|---|---|---|---|
| "Tumor" | 0.526 | 1 | 1 | 0.911 | 0.69 | 0.919 |
| "Mass" | 0.789 | 0.815 | 0.469 | 0.949 | 0.588 | 0.811 |
| "Lesion" | 0.737 | 0.924 | 0.667 | 0.944 | 0.7 | 0.892 |
| Any of 3 terms | 1 | 0.761 | 0.463 | 1 | 0.633 | 0.802 |
| NLP algorithm | 1 | 0.978 | 0.905 | 1 | 0.95 | 0.982 |

trix at each threshold to compare the number of true positives, false positives, true negatives and false negatives at each threshold.

### Comparison to untrained string search

We performed a string search for the terms "tumor," "mass," and "lesion" in each of the 111 MRI reports in our test set. The results of each of these untrained string searches were then compared to that obtained with our NLP algorithm. We also compared a search for any of the 3 terms ("tumor" OR "mass" OR "lesion") against our NLP algorithm.

## Results

The training set consisted of 37,579 radiology reports for 25,469 patients. Of those, 36,676 were labeled negative and 903 as positive for spinal tumor. The model performed favorably on the testing cohort of 111 patients across a variety of thresholds for a positive result, as shown in Table 1. As this model is intended to be used as an initial screening tool, we chose to optimize for a low false negative rate, even at the expense of a higher false positive rate. For that reason, we chose a cutoff of 0.02, which yielded a 100% sensitivity while still maintaining a relatively low false positive rate of 2.2%. The confusion matrix at the threshold of 0.02 and its corresponding receiver operating characteristic curve are shown in Fig. 2 and Fig. 3 respectively.

We compared the performance of our NLP algorithm against an untrained string search approach. The terms "tumor," "mass," and "lesion" were each used for string search. We also performed a search for any 1 of the 3 terms ("tumor" or "mass" or "lesion"). "Tumor "was the most specific term, with a positive predictive value of 100%, but a low sensitivity of 53%. "Mass" performed poorly on both sensitivity and specificity. "Lesion" had the best overall performance of each of the individual terms, but the performance was still significantly worse than our NLP algorithm. Using all 3 terms together yielded a sensitivity of 100% but at the cost of a high number of false positives, with a positive predictive rate of only 46% (vs. 100% for our NLP algorithm). The results of this analysis are summarized in Table 2.

## Discussion

The delayed diagnosis and treatment of spinal metastases can lead to substantial patient morbidity, requiring extensive surgery and long hospital stays, placing place a burden on both the patient's family and the healthcare system. Here we demonstrated that NLP analysis of MRI reports can identify patients with metastatic disease to the spine, allowing for urgent referral without undue resource allocation as part of a comprehensive spinal oncology clinical pathway.

Patient Name:

EXAMINATION: MRI CERVICAL SPINE WITH AND WITHOUT CONTRAST. MRI THORACIC SPINE WITH AND WITHOUT CONTRAST, MRI LUMBAR SPINE WITH AND WITHOUT CONTRAST

IMPRESSION:

Acute compression fracture deformity of the L5 vertebral body with marked associated osseous edema. No appreciable postcontrast enhancement or enhancing epidural soft tissue. This fracture is of indeterminate etiology, possibly pathologic. Recommend correction with tissue sampling.

Mild diffuse low marrow signal throughout the spine, which is nonspecific.

Otherwise no abnormal enhancement along the spinal canal

Bilateral pars interarticular defects at L5-S1

CLINICAL INDICATION: Stage IVFA cervical SCC w/ known metastatic hepatic lesions. :: Cancer

TECHNIQUE:
Multi-planar, multi sequence MR of the cervical, thoracic, and lumbar spine, was performed prior to and following administration of intravenous contrast.
Intravenous contrast: 6.89 mL of Gadavist

COMPARISION: CT dated 2/6/2024

INTERPRETATION:

LOCALIZER: no additional findings.

BONE AND BONE MARROW: There is an acute compression fracture deformity of the L5 vertebral body with marked associated osseous edema. There is no appreciable postcontrast enhancement or enhancing epidural soft tissue.

There are bilateral pars interarticular defects at L5-S1

There is mild diffuse low marrow signal throughout the spine, which is non specific.

ALIGNMENT: The craniocervical junction and atlantoaxial intervals are maintained. There is a grade 1 anterolisthesis of L5 on S1

INTERVERTEBRAL DISCS: There is degenerative loss of disc signal intensity

SPINAL CORD and EPIDURAL: The spinal cord is normal in size and signal intensity. There is no epidural mass or fluid collection.  There is no abnormal enhancement along the spinal canal.

EXTRAPSINAL: There is ill-defined heterogeneously lobe predominant hypoenhancing soft tissue within the pelvis, better evaluated on the recent CT of 2/6/2024.

CONTRAST: Please see above for discussion of area(s) of abnormal enhancement

FINDINGS BY LEVEL:

C2-C3: There is mild disc osteophyte complex asymmetric to the right with thickening of ligamentum flavum resulting in mild spinal canal stenosis. There is mild right neuroforaminal narrowing

C3-C4: There is uncovertebral spurring and facet hypertrophy resulting in mild bilateral neuroforaminal narrowing. There Is no significant spinal canal stenosis.

C4-C5 There is a disc osteophytes complex resulting in mild spinal canal stenosis. There is uncovertebral spurring and facet hypertrophy resulting in mild bilateral neuroforaminal narrowing.

C5-C6: There is a disc osteophyte complex with thickening of ligamentum flavum resulting in mild spinal canal stenosis. There is uncovertebral spurring and facet hypertrophy resulting in mild left greater than right neuroforaminal narrowing.

C6-C7: There is a disc osteophyte complex with thickening of ligamentum flavum resulting in mild to moderate spinal canal stenosis. There is uncovertebral spurring and facet hypertrophy resulting in mild severe right with moderate left neuroforaminal narrowing

C7-T1: There is no significant spinal canal or neural foraminal stenosis

Evaluation of the thoracic spine demonstrates no significant spinal canal stenosis

L1-L2: There is no significant spinal canal or neural foraminal stenosis

L2-L3: There is a mild disc bulge with facet arthropathy and thickening of ligamentum flavum without significant spinal canal stenosis. There Is mild right neuroforaminal narrowing.

L3-L4: There is facet arthropathy and thickening of the ligamentum flavum, resulting in mild spinal canal stenosis and mild bilateral neuroforaminal narrowing.

L4-L5: There is a disc bulge with facet arthropathy and thickening of the ligamentum flavum resulting in mild spinal canal stenosis and mild bilateral neuroforaminal narrowing.

L5-S1 There is a disc bulge with facet arthropathy and thickening of ligamentum flavum resulting in mild spinal canal stenosis and severe bilateral neuroforaminal narrowing.

ATTENDING RADIOLOGIST:

**Fig. 1.** (A) Preprocessed MRI. (B) Postprocessed MRI report.

Acute compression fracture deformity of the l vertebral body with marked associated osseous edema no appreciable postcontrast enhancement of enhancing epidural soft tissue this fracture is of indeterminate etiology possibly pathologic recommend correction with tissue sampling mild diffuse low marrow signal throughout the spine which is nonspecific otherwise no abnormal enhancement along the spinal canal bilateral pars interarticular defects at ls localizer no additional findings bone and bone marrow there is an acute compression fracture deformity of the l vertebral body with marked associated osseous edema there is no appreciable postcontrast enhancement or enhancing epidural soft tissue there are bilateral pars interarticular defects at ls there is mild diffuse low marrow signal throughout the spine which is non specific alignment the craniocervical junction and atlantoaxial intervals are maintained there is a grade anterolisthesis of l on s intervertebral discs there is degenerative loss of disc signal intensity spinal cord and epidural the spinal cord is normal in size and signal intensity there is no epidural mass or fluid collection there is no abnormal enhancement along the spinal canal extraspinal there is ill-defined heterogeneously lobe predominant hypoenhancing soft tissue within the pelvis better evaluated on the recent ct of contrast please see above for discussion of areas of abnormal enhancement findings by level cc there is mild disc osteophyte complex asymmetric to the right with thickening of ligamentum flavum resulting in mild spinal canal stenosis there is mild right neuroforaminal narrowing cc there is uncovertebral spurring and facet hypertrophy resulting in mild bilateral neuroforaminal narrowing there Is no significant spinal canal stenosis cc there is a disc osteophytes complex resulting in mild spinal canal stenosis there is uncovertebral spurring and facet hypertrophy resulting in mild bilateral neuroforaminal narrowing cc there is a disc osteophyte complex with thickening of ligamentum flavum resulting in mild spinal canal stenosis there is uncovertebral spurring and facet hypertrophy resulting in mild left greater than right neuroforaminal narrowing cc there is a disc osteophyte complex with thickening of ligamentum flavum resulting in mild to moderate spinal canal stenosis there is uncovertebral spurring and facet hypertrophy resulting in mild severe right with moderate left neuroforaminal narrowing ct there is no significant spinal canal or neural foraminal stenosis evaluation of the thoracic spine demonstrates no significant spinal canal stenosis ll there is no significant spinal canal or neural foraminal stenosis ll there is a mild disc bulge with facet arthropathy and thickening of ligamentum flavum without significant spinal canal stenosis. There Is mild right neuroforaminal narrowing ll there is facet arthropathy and thickening of the ligamentum flavum, resulting in mild spinal canal stenosis and mild bilateral neuroforaminal narrowing ll there is a disc bulge with facet arthropathy and thickening of the ligamentum flavum resulting in mild spinal canal stenosis and mild bilateral neuroforaminal narrowing ls there is a disc bulge with facet arthropathy and thickening of ligamentum flavum resulting in mild spinal canal stenosis and severe bilateral neuroforaminal narrowing.
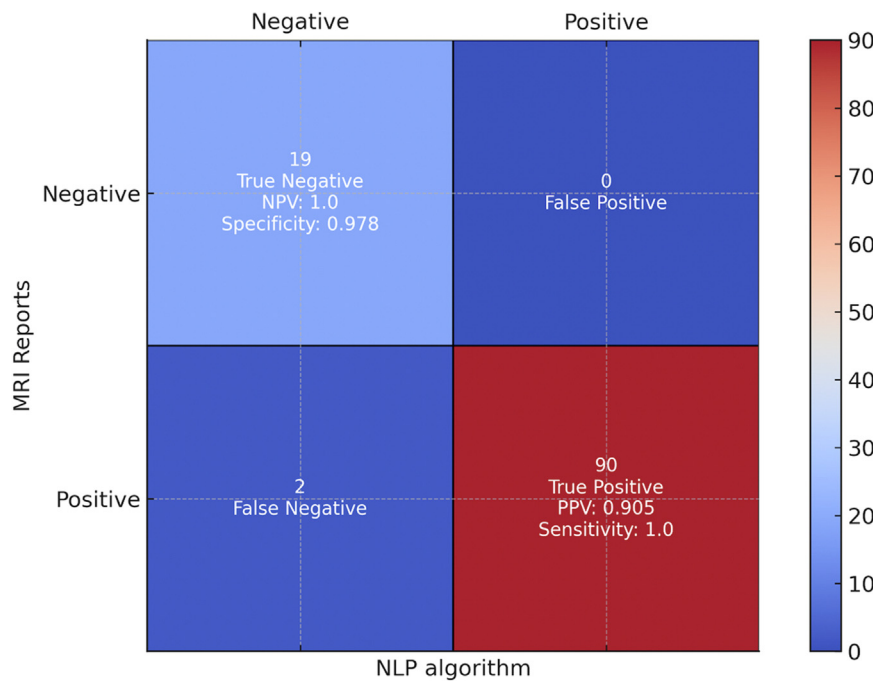
**Fig. 1.** Continued



**Fig. 2.** Confusion matrix at cutoff 0.02.

The rapid proliferation of electronic medical records (EMR) systems has created an opportunity to automate the identification and referral of patients with possible spinal metastasis to appropriate clinical care pathways. Harnessing the full potential of the EMR to gather information on a large number of patients requires scalable and inexpensive approaches such as NLP. The goal is to translate narrative text from the EMR into a structured format or discrete representation suitable for processing by computer algorithms.

However, there are several challenges to the use of NLP to extract data from EMRs, including the unstructured format of most clinical notes, the lack of large ground-truth training sets, the need for greater collaboration between data scientists and clinical domain-knowledge experts and hesitation in applying NLP algorithms to high-stakes clinical scenarios with low tolerance for error. While NLP has been used for a variety of clinical applications, both in spine surgery as well as in a variety of other fields such as neurology and cardiology, actual clinical implementation has been rare for these reasons.

A further barrier to the implementation of NLP algorithms is the "black box" nature of a technology that most healthcare practitioners do not fully understand. Fundamentally, many autonomous technologies such as NLP algorithms continuously refine their internal decision-making structure in such a way that even their human programmers cannot meaningfully understand, much less healthcare practitioners as the end users [15]. While it is possible for practitioners to compare and validate the outcomes of an NLP algorithm as we have done in this study, that does not imply that the user understands the inner workings of how the algorithm arrived at its conclusion. At scale, this black box nature of the algorithms could allow for the introduction of potential bias, accountability, and responsibility issues at large [16].

Despite these challenges, we believe NLP is a promising approach for extraction of useful information from unstructured clinical text. With proper design and training, an NLP approach can provide substantially better accuracy than that of a rules-based approach, as shown in our comparison of NLP to a string search for common metastasis-associated
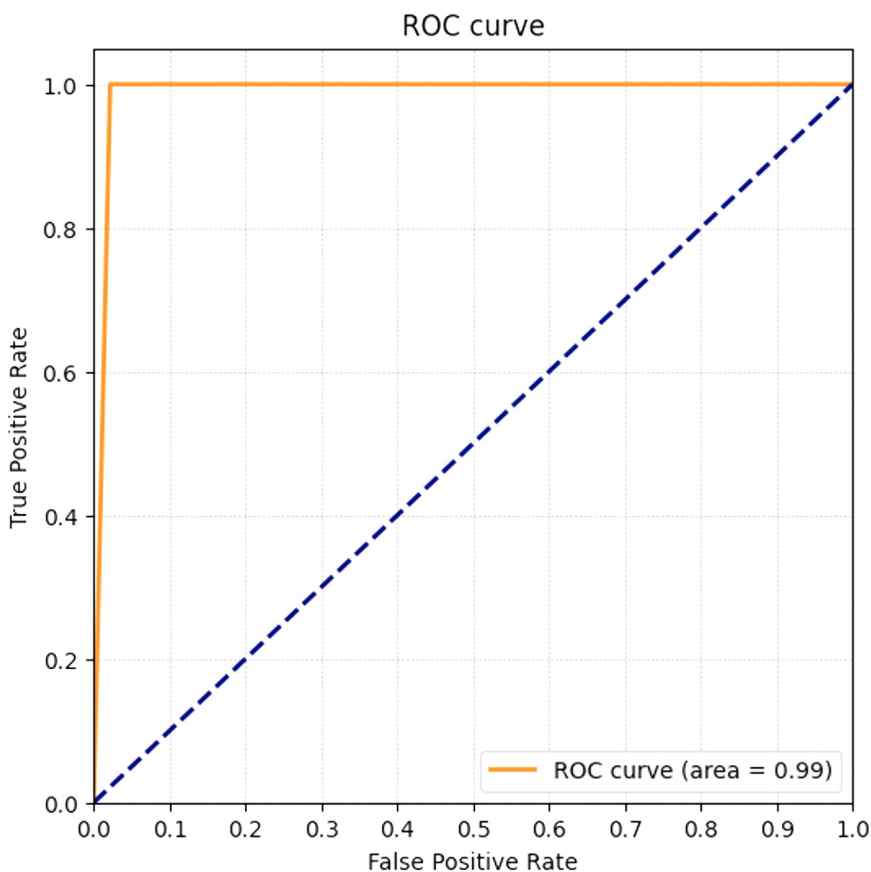
**Fig. 3.** Receiver operating characteristic curve for cutoff at 0.02.

terms in MRI reports. There are several reasons for this. One is that there is heterogeneity in the terminology used between different radiologists. A rules-based approach can also lack appropriate context, for example failing to distinguish an incidental finding of a lesion elsewhere in the body from a spinal lesion. There is also the problem of negation, as a statement such as "there is no evidence of a concerning lesion" could inappropriately be flagged as a positive result by a string search. The granularity provided by an NLP algorithm helps to overcome these obstacles.

There are several limitations to this work. The first is that the detection of possible spinal tumor by our algorithm is dependent upon an MRI report being generated by a radiologist, which may introduce some delay to the process. There have been several attempts in the literature at using computer vision to directly identify spinal tumors from MRI images, although to our knowledge none have been implemented into widespread clinical use [16,17]. The second limitation is that our model was trained using data only from our institution, and may have learned to identify semantic patterns common to our group of neuroradiologists. External validation will determine whether the model needs further training on a more diverse dataset of reports from multiple institutions.

Given the strong performance of the algorithm on our test cohort, we intend to implement our NLP algorithm into our EMR on a prospective observational and interventional basis to permit further internal validation and provide urgent referrals of patients with spinal metastases to our multidisciplinary spinal oncology team. While a finding of cancer on MRI should prompt the radiologist to contact the prescribing doctor, we have found cases at our institution where this does not lead to the patient receiving timely and appropriate care. This is particularly the case when an MRI is ordered for back pain by a primary care or emergency room doctor and the appropriate referral to spine oncology is not

made from there. We plan to use this algorithm to automate notification of our spine oncology team when an MRI shows evidence of spinal cancer. Our spine oncology team can then review these cases to ensure that the patients are receiving appropriate care. We also anticipate that this technology will be of use in the research setting for rapid identification of patient cohorts for retrospective studies. Future studies will include assessment of this algorithm's impact on time to treatment at our institution and external validation of the algorithm on MRI reports from other institutions.

**Conclusions**

The NLP model described in this manuscript predicts the presence of spinal tumor and metastatic cord compression in spine MRI reports with high accuracy. We anticipate this will allow for faster referral of these patients to appropriate specialists and allowing for reduced morbidity and increased survival.

**Declarations of competing interests**

One or more of the authors declare financial or professional relationships on ICMJE-NASSJ disclosure forms.

**Funding**

**References**

[1] Levack P, GJ, Collie D, et al. Don't wait for a sensory level—listen to the symptoms: a prospective audit of the delays in diagnosis of malignant cord compression. Clin Oncol 2002;14(6):472–80.

[2] Meyer HS, WA, Raufer A, et al. Surgery in acute metastatic spinal cord compression: timing and functional outcome. Cancers 2022;14(9):2249.

[3] Quraishi NA, RT, Manoharan SR, et al. Effect of timing of surgery on neurological outcome and survival in metastatic spinal cord compression. Eur Spine J 2013;22(6):1383–8.

[4] Van Tol FR, SK, Choi D, et al. The importance of timely treatment for quality of life and survival in patients with symptomatic spinal metastases. Eur Spine J 2020;29(12):3170–8.

[5] Van Tol FR, CD, Verkooijen HM, et al. Delayed presentation to a spine surgeon is the strongest predictor of poor postoperative outcome in patients surgically treated for symptomatic spinal metastases. Spine J 2019;19(9):1540–7.

[6] Guzik G. Analysis of factors delaying the surgical treatment of patients with neurological deficits in the course of spinal metastatic disease. BMC Palliat Care 2018;17(1):44.

[7] Allan L, BL, Dewar J, et al. Suspected malignant cord compression—improving time to diagnosis via a 'hotline': a prospective audit. Br J Cancer 2009;100:1867–72.

[8] Evans-Lacko S, JM, McCrone P, et al. Facilitators and barriers to implementing clinical care pathways. BMC Health Services Res 2010;10:182.

[9] Casey A, Davidson E, Poon M, et al. A systematic review of natural language processing applied to radiology reports. BMC Med Inform Decis Mak 2021;21(1):179.

[10] Kim C, ZV, Obeid J, et al. Natural language processing and machine learning algorithm to identify brain MRI reports with acute ischemic stroke. PLoS One 2019;28(14):e0212778.

[11] Groot OQ, Bongers MER, Karhade AV, et al. Natural language processing for automated quantification of bone metastases reported in free-text bone scintigraphy reports. Acta Oncol 2020;59(12):1455–60.

[12] Dewaswala N, CD, Bhopalwala H, et al. Natural language processing for identification of hypertrophic cardiomyopathy patients from cardiac magnetic resonance reports. BMC Med Informat Decision Making 2022;22(1):272.

[13] Huang BB, Huang J, Swong KN. Natural language processing in spine surgery: a systematic review of applications, bias, and reporting transparency. World Neurosurg 2022;167:156–164.e6.

[14] Le Q, Mikolov T. Distributed Representations of Sentences and Documents. *Proceedings of Machine Learning Research* 2014;32(2):1188–96.

[15] Hui AT, Ahn SS, Lye CT, Deng J. Ethical challenges of artificial intelligence in health care: a narrative review. Ethics Biol, Eng Med 2021;12(1).

[16] Durán JM, Jongsma KR. Who is afraid of black box algorithms? On the epistemological and ethical basis of trust in medical AI. J Med Ethics 2021;47:329–35.

[17] Ong W, Zhu L, Zhang W, et al. Application of artificial intelligence methods for imaging of spinal metastasis. Cancers (Basel) 2022;14(16):4025.