

# Development of a Scoring Tool for Chronic Nonbacterial Osteomyelitis Magnetic Resonance Imaging and Evaluation of its Interrater Reliability

Yongdong Zhao<sup>ID</sup>, T. Shawn Sato, Sabrina M. Nielsen, Meinrad Beer, Mingqian Huang, Ramesh S. Iyer, Michael McGuire, Anh-Vu Ngo, Jeffrey P. Otjen, Jyoti Panwar<sup>ID</sup>, Jennifer Stimec, Mahesh Thapa, Paolo Toma, Angela Taneja, Nancy E. Gove, and Polly J. Ferguson

**ABSTRACT. Objective.** Serial magnetic resonance imaging (MRI) examinations are often needed in chronic nonbacterial osteomyelitis (CNO) to determine the objective response to treatment. Our objectives in this study were (1) to develop a consensus-based MRI scoring tool for clinical and research use in CNO; and (2) to evaluate interrater reliability and agreement using whole-body (WB)-MRI from children with CNO.

**Methods.** Eleven pediatric radiologists discussed definitions and grading of signal intensity, size of signal abnormality within bone marrow, and associated features on MRI through monthly conference calls and a consensus meeting, using a nominal group technique in July 2017. WB-MRI scans from children with CNO were deidentified for training reading and an interrater reliability study. The reading by each radiologist was conducted in a randomized order. Interrater reliability for abnormal signal and severity were assessed using free-marginal  $\kappa$  statistics.

**Results.** Radiologists reached a consensus on grading CNO-specific MRI findings and on describing bone units based on anatomy. A total of 45 sets of WB-MRI scans, including 4 sets of non-CNO MRI examinations, were selected for the final reading. The mean  $\kappa$  of each category of bones was  $> 0.7$  with majority  $> 0.9$  demonstrating substantial/almost perfect interrater reliability of readings among radiologists. The agreement on signal intensity and the size of signal abnormality within the most commonly affected bones (femur and tibia) were lower than those of other bones.

**Conclusion.** The chronic nonbacterial osteomyelitis magnetic resonance imaging scoring (CROMRIS) tool is a comprehensive standardized scoring tool for MRI in children with CNO. Our interrater study demonstrated good interrater reliability and agreement of readings. (J Rheumatol First Release January 15 2020; doi:10.3899/jrheum.190186)

## Key Indexing Terms:

CHRONIC NONBACTERIAL OSTEOMYELITIS  
INTERRATER RELIABILITY

SCORING TOOL  
MAGNETIC RESONANCE IMAGING

From the Seattle Children's Hospital, Department of Pediatrics, and Department of Radiology, University of Washington; Center for Clinical and Translational Research, Seattle Children's Research Institute, Seattle, Washington; Department of Radiology, University of Iowa Carver College of Medicine, Iowa City, Iowa, USA; Musculoskeletal Statistics Unit, The Parker Institute, Bispebjerg and Frederiksberg Hospital, Copenhagen; The Rheumatology Research Unit, Department of Rheumatology, Odense University Hospital and University of Southern Denmark, Odense, Denmark; Department of Diagnostic and Interventional Radiology, University Hospital of Ulm, Ulm, Germany; Department of Radiology, Stony Brook University Hospital, Stony Brook, New York; Department of Radiology, Hackensack University Medical Center, Hackensack, New Jersey, USA; Department of Medical Imaging, Hospital for Sick Children, University of Toronto, Toronto, Ontario, Canada; Department of Radiology, Christian Medical College and Hospital, Vellore, India; Department of Imaging, Bambino Gesù Children's Hospital, Institute for Research and Health Care (IRCCS), Rome, Italy; Pediatric Rheumatology, Children's Healthcare of Atlanta, Emory University, Atlanta, Georgia; Department of Pediatrics, University of Iowa Carver College of Medicine, Iowa City, Iowa, USA.

This study was funded by a Childhood Arthritis and Rheumatology Research Alliance–Arthritis Foundation small grant. P.J. Ferguson is supported by R01AR059703 from the US National Institutes of

Health/National Institute of Arthritis and Musculoskeletal and Skin Diseases. Y. Zhao is supported by the Clinical Research Scholar Program from Seattle Children's Research Institute and Bristol-Myers Squibb. The Parker Institute, Bispebjerg and Frederiksberg Hospital (S.M. Nielsen), is supported by a core grant from the Oak Foundation (OCAY-13-309).

Y. Zhao, MD, PhD, Seattle Children's Hospital, Department of Pediatrics, University of Washington, and Center for Clinical and Translational Research, Seattle Children's Research Institute; T.S. Sato, MD, Department of Radiology, University of Iowa Carver College of Medicine; S.M. Nielsen, MSc, Musculoskeletal Statistics Unit, The Parker Institute, Bispebjerg and Frederiksberg Hospital, and The Rheumatology Research Unit, Department of Rheumatology, Odense University Hospital and University of Southern Denmark; M. Beer, MD, Department of Diagnostic and Interventional Radiology, University Hospital of Ulm; M. Huang, MD, Department of Radiology, Stony Brook University Hospital; R.S. Iyer, MD, MBA, Department of Radiology, Seattle Children's Hospital, University of Washington; M. McGuire, MD, Department of Radiology, Hackensack University Medical Center; A.V. Ngo, MD, Department of Radiology, Seattle Children's Hospital, University of Washington; J.P. Otjen, MD, Department of Radiology, Seattle Children's Hospital, University of Washington; J. Panwar, MD, FRCR, Department of Medical Imaging, Hospital for Sick Children, University of Toronto, and Department of Radiology, Christian Medical College and Hospital;

Chronic nonbacterial osteomyelitis (CNO) is a pediatric autoinflammatory bone disease challenging to physicians because of its occult nature and the difficulty of assessing disease activity. It is also known as chronic recurrent multifocal osteomyelitis and synovitis, acne, pustulosis, hyperostosis, and osteitis (SAPHO) syndrome. Physical examination and traditional inflammatory markers are not sensitive metrics to monitor disease progression because of occasionally minimal or absent findings on physical examination, normal laboratory values, and lack of correlation between them<sup>1</sup>. Radiographs are only 13–16% sensitive in detecting skeletal lesions in CNO<sup>2</sup> and bone scintigraphy was shown to be only 70% sensitive compared to magnetic resonance imaging (MRI)<sup>3</sup>. The current gold standard imaging modality is whole-body (WB)-MRI<sup>2,4,5</sup>, especially at the initial evaluation. However, the imaging findings of CNO can be nonspecific and bone biopsy may be necessary.

CNO can affect virtually any bone, and there is no uniform approach to assess all bones identically. Previously, CNO lesions on MRI were reported by the number of active lesions<sup>5,6,7,8,9,10</sup> and their anatomical locations. Detailed scoring systems have been reported<sup>11,12</sup>. WB-MRI has the dual advantages of greater sensitivity and lack of ionizing radiation when compared to skeletal scintigraphy<sup>3</sup>, and is more commonly used in pediatric rheumatology across the world<sup>2,4,13,14</sup>. Standardized reporting of each imaging characteristic across all bones of patients with CNO is critical in establishing imaging outcome measurements in CNO for future studies. Our objective is to develop a practical and consensus-based MRI scoring tool for clinical and research use in CNO. Further, interrater agreement and reliability will be evaluated using WB-MRI from children with CNO.

## MATERIALS AND METHODS

The development of the chronic nonbacterial osteomyelitis MRI scoring (CROMRIS) tool consisted of 3 steps: (1) a literature review of previously reported MRI scoring tools of CNO, (2) initial development of a standardized MRI scoring tool for CNO, and (3) a consensus meeting. Subsequently, the interrater agreement and reliability were assessed.

We did a literature review on previously reported MRI scoring tools of CNO as preparation for the meetings. The results of the review were presented at the conference call meetings and consensus conference. Members of an international CNO musculoskeletal radiologist working group initiated the process to develop a standardized MRI scoring tool for CNO at the Society of Pediatric Radiology annual conference in Vancouver,

British Columbia, Canada, in 2015. Since the first meeting, 11 pediatric radiologists, each with at least 5 years of experience reading musculoskeletal and CNO MRI from 7 different pediatric hospitals in North America and Europe, were identified by soliciting pediatric radiologists within the CNO work group. Group members discussed definitions and grading of signal intensity, size of signal abnormality within bone marrow and surrounding tissue, physis damage, and vertebral compression on MRI through monthly conference calls. Representative MRI images [short-tau inversion recovery (STIR) sequence except that skull used T2 sequence from 1.5T or 3T scanner] of active bone inflammation were assembled by members using a separate set of images to establish an atlas to illustrate the proposed scoring system.

**Consensus meeting.** There were 7 radiologists and 2 pediatric rheumatologists (YZ, PJF) at the face-to-face conference (Seattle, July 2017). The facilitators (YZ and PJF) participated in the discussion but were not eligible to vote. Nominal group technique was used to achieve consensus (defined as  $\geq 70\%$  agreement within the group) on all questions considered during the meeting.

**Interrater agreement and reliability.** The interrater agreement and reliability study was approved by the institutional review board from Iowa Children's Hospital (# 201609778). Written informed consent was waived owing to the retrospective nature and use of anonymized images. A total of 82 sets of preexisting WB-MRI scans (STIR sequence with 3–4 mm thickness from 1.5T or 3T scanner) between January 2013 and August 2016 from children with CNO or other diseases at the University of Iowa Children's Hospital were used for training reading and for assessing interrater agreement and reliability. A video tutorial was produced for training and interrater calibration exercise. Nine sets of MRI examinations were used for the training reading to improve familiarity with the tool before a reliability study. Of the 82 sets of MRI, these were excluded: 4 from subjects older than 18 years, 9 sets for training, and 1 set from a patient with leukemia. To assess interrater agreement and reliability, each radiologist read in a randomized order, among the remaining 68 sets of MRI from 45 patients (19 patients had MRI at more than 1 timepoint), 45 sets of MRI examinations from 45 patients (limit 1 set per patient and the set at the beginning of the disease course if more than 1 set is available), including 4 sets of MRI studies from non-CNO patients. Controls were included in the analyses to ensure variability in the sample. Data were recorded with a detailed scoring form (Supplement 1, available with the online version of this article). There was no gold standard defined for comparisons.

**Statistical analysis.** For the interrater agreement and reliability study, descriptive analysis was performed to assess the prevalence of abnormalities at each site defined as agreement among  $> 70\%$  of the radiologists. Data were presented combining similar types of bones per patient. Absolute agreement for each site was defined as the proportion of patients for whom the ratings were the same for all 11 radiologists.

We assessed interrater reliability (i.e., how well the persons can be distinguished from each other despite measurement errors) using the free-marginal  $\kappa$  statistic described by Brennan and Prediger<sup>15</sup>. The free-marginal  $\kappa$  statistic is recommended when raters are not instructed about the number of observations that should be assigned to each category<sup>15</sup> and when the distribution of ratings is highly skewed<sup>16</sup>. The  $\kappa$  coefficients were interpreted according to Landis and Koch<sup>17</sup>. Mean  $\kappa$  (and range) was calculated by categories of bones: the spine, complex bone, flat bones, hand/foot, and long bones. The long bones were further divided into proximal epiphysis, proximal metaphysis, diaphysis, distal metaphysis, and distal epiphysis. All analyses were conducted using R version 3.5.1<sup>18</sup>.

## RESULTS

**Literature review.** A search was conducted in PubMed using the following MeSH terms: (SAPHO[All Fields] OR “chronic recurrent multifocal osteomyelitis”[All Fields] OR “chronic nonbacterial osteomyelitis”[All Fields] OR “non-bacterial osteitis”[All Fields]) AND (“magnetic resonance

J. Stimec, MD, Department of Medical Imaging, Hospital for Sick Children, University of Toronto; M. Thapa, MD, MEd, Department of Radiology, Seattle Children's Hospital, University of Washington; P. Toma, MD, Department of Imaging, Bambino Gesù Children's Hospital; A. Taneja, MD, Pediatric Rheumatology, Children's Healthcare of Atlanta, Emory University; N.E. Gove, PhD, Center for Clinical and Translational Research, Seattle Children's Research Institute; P.J. Ferguson, MD, Department of Pediatrics, University of Iowa Carver College of Medicine. Address correspondence to Dr. Y. Zhao, MA 7.110, 4800 Sand Point Way NE, Seattle, Washington 98105, USA. E-mail: yongdong.zhao@seattlechildrens.org Accepted for publication September 18, 2019.

imaging”[MeSH Terms] OR (“magnetic”[All Fields] AND “resonance”[All Fields] AND “imaging”[All Fields]) OR “magnetic resonance imaging”[All Fields] OR “mri”[All Fields]) AND (Score[All Fields] OR scoring[All Fields]). Five peer-reviewed publications were identified and one<sup>4</sup> was excluded because it did not mention a scoring system. A total of 3 separate tools were reported in the remaining 4 eligible articles. Two reported an MRI score system for the osteitis lesions ranged from zero to 2 points and the highest score among lesions was used to indicate disease severity in SAPHO<sup>19,20</sup>. Bone marrow edema, bone erosions, or synovitis (with or without joint effusion) were ascertained. The presence of only 1 finding was scored 1 point and 2 or more findings, 2 points. A second tool used a semiquantitative approach to evaluate the characteristics of CNO lesions from MRI in children<sup>11</sup>. A comprehensive grading system for the evaluation of the extent of bone edema and soft tissue inflammation was reported, as well as the presence or absence of periosteal reaction, hyperostosis, physeal damage, and vertebral compression<sup>11</sup>. A third tool, a radiologic index for WB-MRI in patients with nonbacterial osteitis (RINBO), defined the size of active lesions by the absolute measurements and clustered the number of active lesions into 3 categories as unifocal, paucifocal (2, 3, or 4 lesions), and multifocal (5 or more lesions)<sup>12</sup>. Soft tissue inflammation, periosteal reaction, and hyperostosis were classified as extramedullary findings and spinal involvement was distinguished between active with abnormal STIR signal and chronic with deformation. Surrounding soft tissue inflammation was not included. Points were assessed for each of 4 areas of interest [number of radiologic active lesions (RAL), maximum size of RAL, extramedullary affection, and spine involvement], with a maximum score of 10.

Typical WB-MRI protocols include coronal images of the entire body and sagittal images of the entire spine, acquired with fluid-sensitive sequence (STIR, turbo-inversion-recovery–magnitude, or fat saturation) without contrast. Axial sequences of the pelvis and knees, and sagittal images of the ankles and feet, were also included at one of the centers (Iowa) that enhanced lesion identification in commonly affected sites. This protocol was therefore adopted by the group with consensus. T1-weighted images have been used to confirm findings from fluid-sensitive sequence in CNO. It was considered optional because it adds scanning time. Diffusion-weighted imaging (DWI) was reported<sup>21</sup>; however, it was not routinely performed in participating institutions. One study did not show difference in sensitivity of differentiating CNO lesions between STIR sequence alone and combining T1-weighted, DWI, and STIR sequences<sup>22</sup>. Thus, T1-weighted and DWI sequences were not included in the scoring system but use of DWI should be reconsidered when more data are available on its use in CNO. Detailed discussion based on the reported scoring tools led to the newly developed tool.

*The consensus process of the final CROMRIS tool.* At the 2017 conference, consensus defined as  $\geq 70\%$  agreement within the group<sup>23</sup> was reached on all questions considered during the meeting. The complete atlas developed following the consensus meeting includes evaluation of 20 sites using 4 different variables (Supplement 2, available with the online version of this article).

*Inclusion and definition of various characteristics of MRI findings in CNO.* As presented in Table 1, hyperintensity of bone marrow was defined as increased STIR signal within bone marrow compared to the nearby normal marrow, as per the interpreting radiologist’s assessment. Terminology of *bone edema* was discussed and replaced by *bone marrow hyperintensity* with consensus for scientific clarity and the uncertainty of pathology. Linear metaphyseal lines caused by bisphosphonate were included in the atlas to avoid misinterpretation as bone marrow hyperintensity. Periosteal reaction was deemed difficult to confirm by MRI whereas soft tissue inflammation was readily detectable. Thus, “hyperintensity of surrounding tissue” was included with consensus to report the presumed inflammation within soft tissue and periosteum. *Hyperostosis* was a common term used in radiography though identifiable on MRI as *bony expansion*. Thus the latter term was adopted by the group. Vertebral compression and joint effusion were included. Growth plate irregularity was discussed and voted not suitable for assessment in MRI with consensus. Kyphosis and limb hypertrophy were assessable in WB-MRI and thus included in this tool. Leg length discrepancy cannot be assessed reliably in MRI and thus was voted not to be included as part this tool. None of the above measures was assigned as acute or chronic at this stage because a prospective longitudinal study is required to distinguish among them.

*Grading scale of variables and definition of bone units.* In general, signal intensity of bone marrow was graded with 3 levels: absent, less than fluid signal, and similar to fluid signal. Confidence level of identifying abnormal signal was also recorded as low, medium, or high. The size of signal intensity within each unit/segment was graded using relative measurement because of various body sizes and bone sizes in affected patients. Small was defined as  $< 25\%$  of estimated volume, medium as  $25\text{--}50\%$  of estimated volume, and large as  $> 50\%$  of estimated volume. When imaging was inadequate for a confident estimate of the size, “unable to estimate the size” was recorded. The following variables were graded as present or absent: signal hyperintensity of surrounding tissue (soft tissue/periosteum), bony expansion, continuity of signal abnormality between diaphysis and adjacent segment in long bones, hypertrophy of limbs, signal intensity of posterior and/or lateral elements in spine, and kyphosis of entire spine. Vertebral compression was graded as normal, presence of some height loss, or plana (defined as complete flattening of a vertebral body).

The division of bone units and segments was discussed,



Table 1. Definition of variables in CROMRIS tool.

Variable	Definition
Hyperintensity within bone marrow	Presence of abnormal signal intensity within bone marrow on a fluid-sensitive sequence* with confluence pattern, per radiologist's discretion
Hyperintensity within surrounding tissue	Presence of abnormal signal intensity other than a normal luminal structure (i.e., bladder, intestine, cerebrospinal fluid space, vasculature) within surrounding tissue on a fluid-sensitive sequence*
Bony expansion	Enlarged bone contour that is greater than expected
Joint effusion	More than physiological amount of joint fluid within a joint space
Vertebral compression	Decreased height compared to the adjacent vertebra
Limb hypertrophy	Abnormally increased size of the limb comparing to contralateral side
Kyphosis	An exaggerated outward curve of spine on sagittal view

\* Short-tau inversion recovery, fat saturation, or turbo-inversion-recovery-magnitude. CROMRIS: chronic nonbacterial osteomyelitis magnetic resonance imaging scoring.

and the consensus was to follow anatomical divisions in complex bones and group bones into 1 unit in less commonly affected sites (hands and fore-/midfoot) and less well visualized sites. Long bones were divided into the following 5 segments anatomically: proximal epiphysis, proximal metaphysis, diaphysis, distal metaphysis, and distal epiphysis. The spine was graded as individual vertebrae from cervical to lumbar region. However, in addition to the grading of anterior vertebral body, there were reports of abnormal signals within “lateral and posterior elements” including pedicles, lamina, and posterior processes. Based on existing literature, the prevalence of signal hyperintensity within metatarsal bones is less common than in the talus and calcaneus<sup>24</sup>. Therefore, the consensus was to grade any signal hyperintensity within metatarsal bones as abnormal, and only signal hyperintensity with confluence in talus or calcaneus as abnormal.

Total scores as reported by RINBO<sup>12</sup> were not recommended because our first step was to describe and grade lesions from each individual bone unit reliably. Future studies will be needed to determine the exact weight of each characteristic using a much larger representative cohort.

**Interrater agreement and reliability.** The 45 subjects were mainly females with a median age of 11 years [interquartile range (IQR) 9–15] and a median disease duration of about 3.3 years (Table 2). About 80% of WB-MRI were collected with additional axial images of pelvis and knees, and sagittal images of ankles/feet, in addition to the coronal plane images of the entire body and sagittal sequences of entire spine, as done in 20% of subjects. The 11 raters were mainly from the United States, with a median 7 years of experience (IQR 6–10; Supplement 3, available with the online version of this article).

Lower extremities were the bones most commonly affected by CNO, with abnormal bone marrow signal (Figure 1). Upper extremities, including humerus, radius, and hand,

Table 2. Patient characteristics.

Variables	Patients, n = 45
Age, yrs	11 (9–15)
Height, cm <sup>†</sup>	148 (134–167)
Weight, kg <sup>†</sup>	41 (31–65)
Female, n (%)	31 (69)
CNO diagnosis <sup>‡</sup> , n (%)	41 (91)
Duration of disease <sup>†</sup> , mos	40 (16–56)
Basic WB-MRI <sup>§</sup> , n (%)	10 (22)
Complete WB-MRI <sup>¶</sup> , n (%)	35 (78)

Data are presented as median (IQR) unless stated otherwise. <sup>†</sup> Data were available for only 44 patients on height and weight, and for 41 patients on duration of disease. <sup>‡</sup> Other 4 patients had recurrent fever (1), juvenile idiopathic arthritis (1), and unknown conditions (2) at the time of MRI. <sup>§</sup> Basic WB-MRI protocol includes STIR sequence of coronal plane of entire body (upper extremities excluded) in 4–5 stations, and sagittal plane of entire spine in 2 stations. <sup>¶</sup> Complete WB-MRI protocol added STIR sequence of axial plane of pelvis and knees, coronal plane of upper extremities, as well as sagittal plane of ankles to the basic WB-MRI protocol. CNO: chronic nonbacterial osteomyelitis; IQR: interquartile range; WB: whole body; MRI: magnetic resonance imaging; STIR: short-tau inversion recovery.

were reported at 2–9% presence among these patients. Along the spine, the thoracic spine was the most commonly affected site. Pelvic bones, clavicle, and mandible were well represented. Lesions were absent within this cohort in the cervical spine, manubrium/sternum, rib, scapula, skull, and ulna. Hyperintensity within surrounding tissue was detected adjacent to tibia, femur, fibula, foot, humerus, peri-acetabulum, clavicle, and mandible. Bony expansion was present only in the femur, humerus, clavicle, and mandible. Vertebral compression was mostly present in the thoracic spine. Detailed data from individual bone units (i.e., left femur, right mandible) are available in Supplement 4 (available with the online version of this article).

The signal intensity of bone marrow hyperintensity had

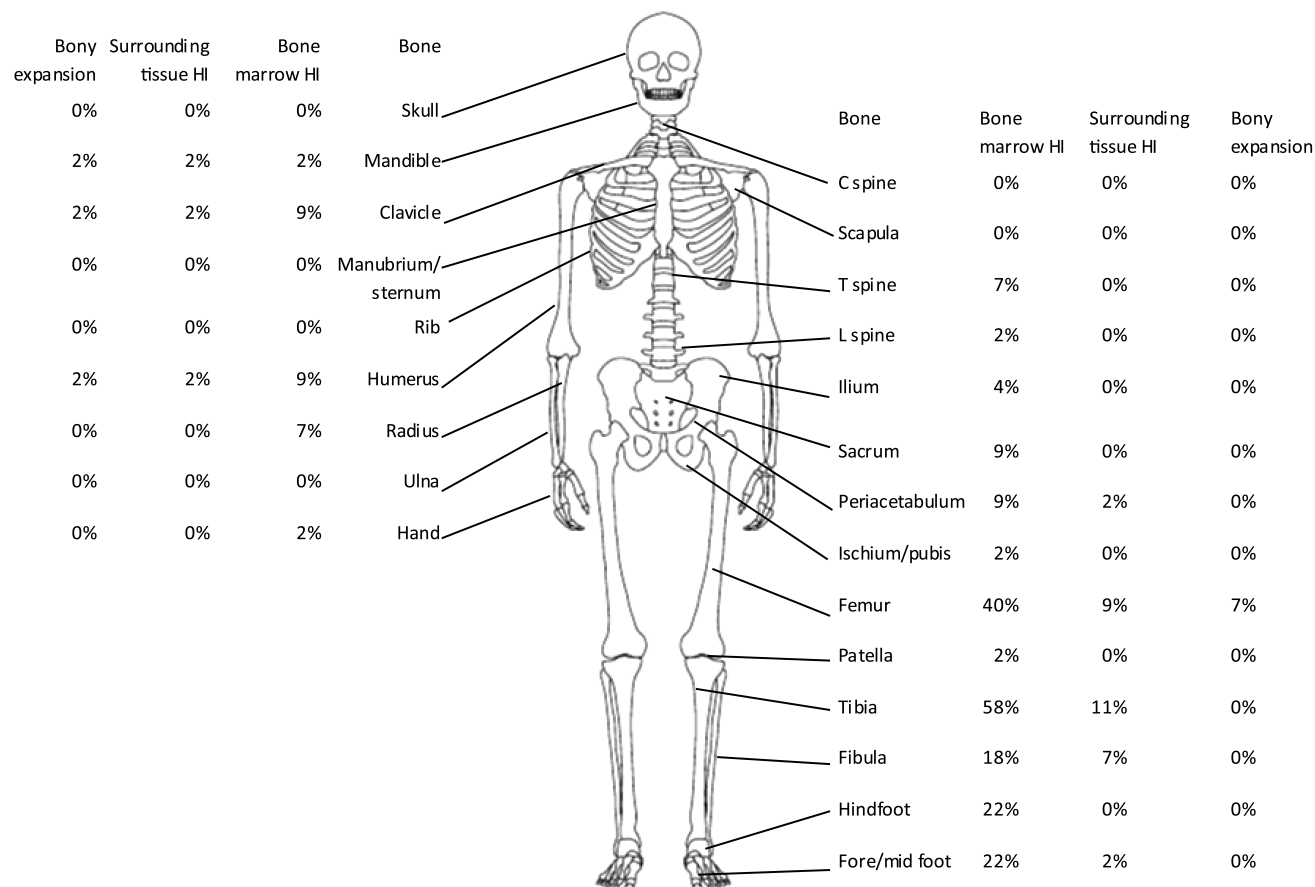


Figure 1. Mean prevalence of bone lesions based on > 70% agreement among 11 radiologists on the presence of bone marrow hyperintensity (HI), surrounding tissue hyperintensity, and bony expansion within the entire skeleton.

low absolute agreements (< 60%) in more commonly affected bones such as femur, tibia, fore-/midfoot, hindfoot, and clavicle (Figure 2A). The majority of less commonly affected bones, including the spine, pelvis, hands, scapula, patella, and radius, had near or greater than 80% of absolute agreements. The presence of hyperintensity within surrounding tissue and bony expansion agreed very well (> 80%) in all bones (Figure 2B, 2C). Detailed data from individual bone units are in Supplement 5 (available with the online version of this article). Most segments of femur and tibia had lower agreement for the size of bone marrow hyperintensity compared to other long bones (Figure 3A). Among other bones, all had good absolute agreement (> 80%) except for the clavicle, mandible, fore-/midfoot, and hindfoot (Figure 3B). The severity of vertebral compression assessed by radiologists has shown excellent absolute agreement in all patients (Figure 3C).

The mean  $\kappa$  of each category was > 0.7, with a majority > 0.9 demonstrating substantial/almost perfect reliability (Table 3). The lowest  $\kappa$  coefficient was observed in bone marrow hyperintensity for the tibia (right, 0.60, 95% CI 0.49–0.71) and the corresponding absolute agreement was only

29% (Supplement 6, available with the online version of this article). Spine, complex bones (pelvis), and flat bones had higher agreements in bone marrow hyperintensity than did hands/feet and long bones. The signal size of bone marrow hyperintensity within each category agreed perfectly, although hands/feet and proximal/distal metaphysis of long bones had the lowest  $\kappa$  scores. The reliability of presence of hyperintensity within surrounding tissue, presence of bony expansion, and vertebral compression were all almost perfect. Detailed data from individual bone units are available in Supplement 6. Joint effusion data showed excellent agreement (Supplement 6). Most low- and medium-confidence readings were from more commonly affected sites such as the femur, tibia, and foot (Supplement 7, available with the online version of this article).

## DISCUSSION

This is the first consensus-based MRI scoring tool for children with CNO and the first comprehensive assessment of interrater reliability of such a tool. Our tool includes the most commonly described characteristics seen in children with CNO from MRI and the grading system can be used as

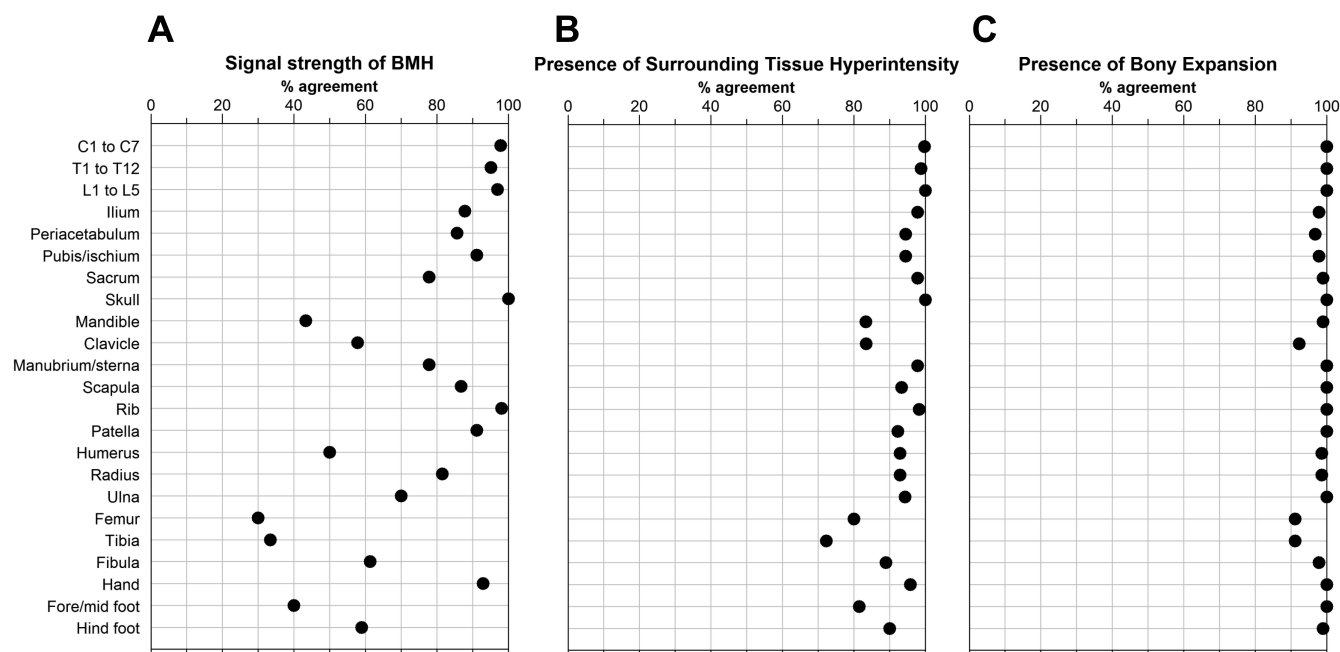


Figure 2. A. Mean absolute agreement of bone lesions among all 11 radiologists on the signal intensity of bone marrow hyperintensity (BMH). B. Presence of surrounding tissue inflammation. C. Presence of bony expansion. C: cervical; T: thoracic; L: lumbar.

a potential research tool after further development and validation. An atlas and training video were developed that may guide radiologists who are less familiar or less experienced in reporting MRI from these affected children.

We have further defined these variables and developed a semiquantitative scoring system as an assessment tool for longitudinal studies to measure the response to treatments. Comparing to the RINBO system<sup>12</sup>, our scoring tool included bone marrow hyperintensity (bone edema), size of bone lesion, vertebral compression, and bony expansion (hyperostosis). Several key differences between these 2 tools are (1) periosteal reaction was deemed not reliable by our group in consensus and so was not included in the current tool; (2) the size of lesion was reported in the current tool as relative to the bone, which is more appropriate for a pediatric population; and (3) a total score was not proposed because further studies are needed to determine the weight of each variable.

Defining the minimum abnormal signal is challenging because of individual scoring variations, as suggested by the low absolute agreement of signal hyperintensity in the commonly affected bones (tibia and femur). Therefore, we used a predefined 70% agreement as a threshold to determine whether a “true” abnormal signal existed in bone marrow. Based on this principle, we found a distribution of lesions among the entire skeleton similar to previous reports<sup>9,10,25,26</sup>. Abnormal signal within surrounding tissue and bony expansion were present at most long bones, but were uncommonly seen in the clavicle and mandible.

In addition, the absolute agreement of the intensity of signal abnormality was poor in commonly affected sites, suggesting that individual radiologists differ in their assessing of various levels either because of inadequate calibration/training or inherent challenge from defined classification. Most low- and medium-confidence readings were from commonly affected sites. These results suggest that adding mandated calibration exercise with a special focus on less conspicuous lesions might improve the interrater agreement. In contrast, the absolute agreements of abnormal signal in surrounding tissue and bony expansion were > 80% except for the tibia. Although the prevalence of these findings was less common than that of bone marrow hyperintensity, it was likely that these features were more distinguishable by radiologists and thus there was more agreement among radiologists.

The  $\kappa$  analysis showed moderate to substantial agreement on the MRI size readings of most commonly affected bones (tibia and femur). When grouped into large categories such as long bones or spine, the agreement significantly increased, which was likely due to the relatively fewer abnormal signals. Hands and feet were scored as regions by grouping multiple bones and the size of bone marrow hyperintensity may not be estimated well enough. It explained why the agreement of this variable is the lowest among all categories. Similarly, the signal size of bone marrow hyperintensity of proximal and distal metaphyses of long bones also had the least agreement because of the difficulty of clearly identifying the border/definition of this segment within long bones. These are very

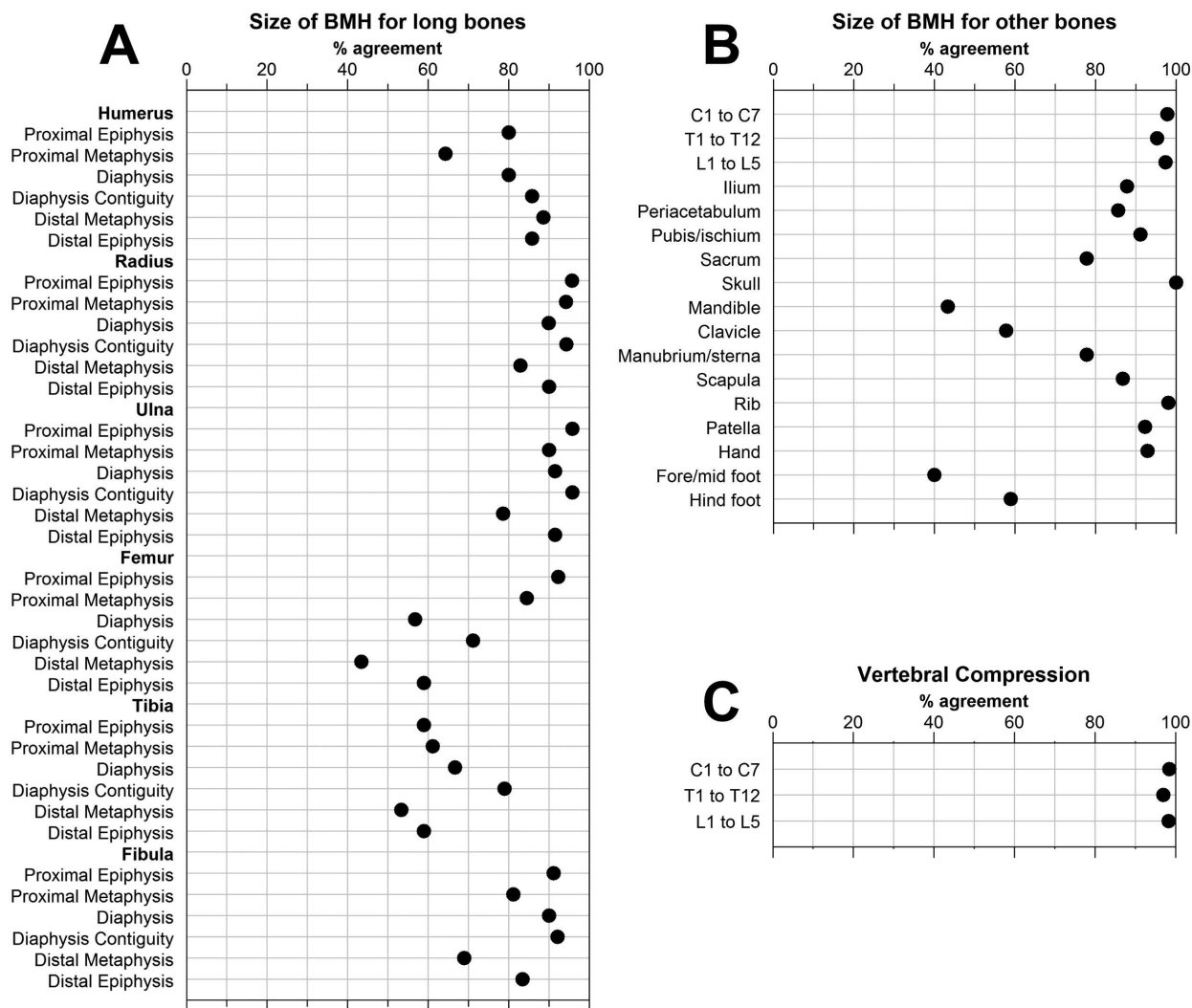


Figure 3. A. Mean absolute agreement of the signal size of bone marrow hyperintensity (BMH) in each segment of long bones. B. Mean absolute agreement of the signal size of BMH in other bones. C. Absolute agreement of the severity of vertebral compression. C: cervical; T: thoracic; L: lumbar.

helpful observations that will allow further improvements of our scoring system. Future studies will aim to answer the following questions: (1) is this scoring tool sensitive to change of clinical disease activities in CNO from a longitudinal study, and what is the intrarater reliability?; (2) what is the interrater reliability of this tool in a validation cohort?; and (3) how should the researcher integrate scores from each body site as a total score for disease activity on a whole-body level, and can this score differentiate patients with CNO from those without CNO?

There were limitations of our study. First, even with our large sample size of subjects, some bone sites were not well represented for interrater study. Future studies using a different subset of MRI with enriched prevalence of signal abnormalities in less commonly affected sites and inadequately scanned area (i.e., upper extremities) are needed to validate our findings. Second, joint effusion was not

adequately scored but owing to its complexity and less weight in managing these patients, we decided that this should be a separate effort. Thirdly, there was no gold standard of the abnormal signals identified by radiologists for our study. A more objective approach of identifying signal threshold is needed and may be accomplished through machine learning by creating a consensus reading result. Fourth, lower agreement and reliability may have been obtained as a result of unequal familiarity of the tool despite the training. Fifth, even with radiologists from 7 centers, this consensus may not be completely representative. Finally, the correlation of abnormal signals on MRI and the actual pathology from CNO was not confirmed. Therefore, a longitudinal study with detailed clinical characterization in children with CNO and healthy children may shed light on the clinical significance of these variables. Nevertheless, we developed a comprehensive MRI scoring tool for CNO with

Table 3. The  $\kappa$  results of MRI readings among radiologists for each category of bones.

Category	Signal Intensity of BMH	Signal Size of BMH	Presence of Surrounding Tissue Hyperintensity	Presence of Bony Expansion	Vertebral Compression
Spine	0.98 (0.96–0.99)	0.98 (0.96–0.99)	0.99 (0.98–0.99)	NA*	0.98 (0.97–0.99)
Complex bone	0.93 (0.89–0.98)	0.94 (0.91–0.96)	0.97 (0.96–0.98)	0.98 (0.97–0.99)	—
Flat bone	0.93 (0.72–0.99)	0.95 (0.87–0.99)	0.97 (0.91–0.99)	0.96 (0.93–0.99)	—
Hand/foot	0.80 (0.63–0.94)	0.83 (0.67–0.97)	0.94 (0.92–0.98)	0.99 (0.99–0.99)	—
Long bone	0.75 (0.60–0.90)	—	0.92 (0.76–0.99)	0.97 (0.93–0.99)	—
Proximal epiphysis	—	0.95 (0.87–0.99)	—	—	—
Proximal metaphysis	—	0.89 (0.78–0.96)	—	—	—
Diaphysis	—	0.94 (0.89–0.97)	—	—	—
Diaphysis contiguity	—	0.91 (0.80–0.98)	—	—	—
Distal metaphysis	—	0.81 (0.65–0.96)	—	—	—
Distal epiphysis	—	0.92 (0.85–0.96)	—	—	—

Data are mean (range). \*  $\kappa$  could not be calculated because of full agreement for all subsets within the spine. Spine includes cervical, thoracic, and lumbar vertebrae. Complex bone refers to pelvis that was divided into ilium, periacetabulum, pubis/ischium, and sacrum on each side. Flat and irregular bones included skull, mandible, clavicle, sterna/manubrium, ribs, patella, and scapula. Hands were graded as 1 unit including phalanges, metacarpal, and carpal bones on each side. Feet were divided into fore-/midfoot and hindfoot. Fore-/midfoot included phalanges, metatarsal, and tarsal bones. Hindfoot included talus and calcaneus. MRI: magnetic resonance imaging; BMH: bone marrow hyperintensity; NA: not applicable.

a consensus from experienced radiologists across 7 centers and 2 continents and showed excellent reliability and agreements in each category of bones and moderate to substantial reliability and agreements in readings from individual bones.

The CROMRIS tool was developed as a comprehensive standardized scoring tool for MRI in children with CNO. Our interrater study demonstrated good interrater reliability and agreement of readings from a group of radiologists. Because CNO is a rare disease and collaborative research is needed in this field, a consensus-based system, such as the CROMRIS tool, representing experienced radiologists from different centers and countries, will likely be adopted by future studies. This tool can be validated in a prospective study and may become a key element of disease activity assessment in CNO.

## ACKNOWLEDGMENT

The authors thank other physician participants who contributed in initial conference calls: Nancy Chauvin, Kirsten Ecklund, and Andrea Doria. We are grateful for the excellent advice from Dr. Maarten Boers to help improve the presentation of figures. Drs. Matthew Basiaga and Eric Allenspach critically reviewed the manuscript.

## ONLINE SUPPLEMENT

Supplementary material accompanies the online version of this article.

## REFERENCES

- Zhao Y, Ferguson PJ. Chronic nonbacterial osteomyelitis and chronic recurrent multifocal osteomyelitis in children. *Pediatr Clin North Am* 2018;65:783-800.
- Fritz J, Tzaribatchev N, Claussen CD, Carrino JA, Horger MS. Chronic recurrent multifocal osteomyelitis: comparison of whole-body MR imaging with radiography and correlation with clinical and laboratory data. *Radiology* 2009;252:842-51.
- Morbach H, Schneider P, Schwarz T, Hofmann C, Raab P, Neubauer H, et al. Comparison of magnetic resonance imaging and Technetium-labelled methylene diphosphonate bone scintigraphy in the initial assessment of chronic non-bacterial osteomyelitis of childhood and adolescents. *Clin Exp Rheumatol* 2012;30:578-82.
- Voit AM, Arnoldi AP, Douis H, Bleisteiner F, Jansson MK, Reiser MF, et al. Whole-body magnetic resonance imaging in chronic recurrent multifocal osteomyelitis: Clinical longterm assessment may underestimate activity. *J Rheumatol* 2015;42:1455-62.
- Roderick M, Shah R, Finn A, Ramanan AV. Efficacy of pamidronate therapy in children with chronic non-bacterial osteitis: disease activity assessment by whole body magnetic resonance imaging. *Rheumatology* 2014;53:1973-6.
- Beck C, Morbach H, Beer M, Stenzel M, Tappe D, Gattenlöhner S, et al. Chronic nonbacterial osteomyelitis in childhood: prospective follow-up during the first year of anti-inflammatory treatment. *Arthritis Res Ther* 2010;12:R74.
- Hospach T, Langendoerfer M, von Kalle T, Maier J, Dannecker GE. Spinal involvement in chronic recurrent multifocal osteomyelitis (CRMO) in childhood and effect of pamidronate. *Eur J Pediatr* 2010;169:1105-11.
- Miettinen PM, Wei X, Kaura D, Reslan WA, Aguirre AN, Kellner JD. Dramatic pain relief and resolution of bone inflammation following pamidronate in 9 pediatric patients with persistent chronic recurrent multifocal osteomyelitis (CRMO). *Pediatr Rheumatol Online J* 2009;7:2.
- Wipff J, Costantino F, Lemelle I, Pajot C, Duquesne A, Lorrot M, et al. A large national cohort of French patients with chronic recurrent multifocal osteitis. *Arthritis Rheumatol* 2015;67:1128-37.
- Jansson A, Renner ED, Ramser J, Mayer A, Haban M, Meindl A, et al. Classification of non-bacterial osteitis: retrospective study of clinical, immunological and genetic aspects in 89 patients. *Rheumatology* 2007;46:154-60.
- Zhao Y, Chauvin NA, Jaramillo D, Burnham JM. Aggressive therapy reduces disease activity without skeletal damage progression in chronic nonbacterial osteomyelitis. *J Rheumatol* 2015;42:1245-51.
- Arnoldi AP, Schlett CL, Douis H, Geyer LL, Voit AM, Bleisteiner F, et al. Whole-body MRI in patients with non-bacterial osteitis: radiological findings and correlation with clinical data. *Eur Radiol* 2017;27:2391-9.
- Guérin-Pfyffer S, Guillaume-Czitrom S, Tammam S, Koné-Paut I. Evaluation of chronic recurrent multifocal osteitis in children by whole-body magnetic resonance imaging. *Joint Bone Spine* 2012;79:616-20.



14. Darge K, Jaramillo D, Siegel MJ. Whole-body MRI in children: current status and future applications. *Eur J Radiol* 2008;68:289-98.
15. Brennan RL, Prediger DJ. Coefficient Kappa: some uses, misuses, and alternatives. *Educ Psychol Meas* 1981;41:687-99.
16. Quarfoot D, Levine RA. How robust are multirater interrater reliability indices to changes in frequency distribution? *Am Stat* 2016;70:373-84.
17. Landis JR, Koch GG. The measurement of observer agreement for categorical data. *Biometrics* 1977;33:159-74.
18. R Core Team. R: A language and environment for statistical computing. Vienna: R Foundation for Statistical Computing; 2013.
19. Assmann G, Kueck O, Kirchhoff T, Rosenthal H, Voswinkel J, Pfreundschuh M, et al. Efficacy of antibiotic therapy for SAPHO syndrome is lost after its discontinuation: an interventional study. *Arthritis Res Ther* 2009;11:R140.
20. Jung J, Molinger M, Kohn D, Schreiber M, Pfreundschuh M, Assmann G. Injection into sternocostoclavicular joints in patients with SAPHO Syndrome. *Semin Arthritis Rheum* Elsevier Inc.; 2012;42:266-70.
21. Leclair N, Thörmer G, Sorge I, Ritter L, Schuster V, Hirsch FW. Whole-body diffusion-weighted imaging in chronic recurrent multifocal osteomyelitis in children. *PLoS One* 2016;11:e0147523.
22. Merlini L, Carpentier M, Ferrey S, Anooshiravani M, Poletti P, Hanquinet S. Whole-body MRI in children: would a 3D STIR sequence alone be sufficient for investigating common paediatric conditions? A comparative study. *Eur J Radiol* 2017;88:155-62.
23. Orbai AM, de Wit M, Mease P, Shea JA, Gossec L, Leung YY, et al. International patient and physician consensus on a psoriatic arthritis core outcome set for clinical trials. *Ann Rheum Dis* 2017;76:673-80.
24. Shabshin N, Schweitzer ME, Morrison WB, Carrino JA, Keller MS, Grissom LE. High-signal T2 changes of the bone marrow of the foot and ankle in children: red marrow or traumatic changes? *Pediatr Radiol* 2006;36:670-6.
25. Borzutzky A, Stern S, Reiff A, Zurakowski D, Steinberg EA, Dedeoglu F, et al. Pediatric chronic nonbacterial osteomyelitis. *Pediatrics* 2012;130:e1190-7.
26. Roderick MR, Shah R, Rogers V, Finn A, Ramanan AV. Chronic recurrent multifocal osteomyelitis (CRMO) — advancing the diagnosis. *Pediatr Rheumatol Online J* 2016;14:47.