# Annotations, Ontologies, and Whole Slide Images – Development of an Annotated Ontology-Driven Whole Slide Image Library of Normal and Abnormal Human Tissue

**Karin Lindman[1], Jerómino F. Rose[2], Martin Lindvall[3], Claes Lundström[4], Darren Treanor[5,6,7]**

[1]Department of Clinical Pathology, Region Östergötland, Linköping, Sweden, [2]Center for Medical Image Science and Visualization, Linköping University, Linköping, Sweden, [3]Center for Medical Image Science and Visualization (CMIV), Linköping University, Linköping and Sectra AB, Sweden, [4]Center for Medical Image Science and Visualization, Linköping University, Linköping and Sectra AB, Linköping, Sweden, [5]Department of Clinical Pathology, Region Östergötland, Linköping, Sweden, [6]Department of Clinical Pathology, and Department of Clinical and Experimental Medicine (IKE), Linköping University, Linköping, Sweden, [7]Department of Cellular Pathology, St. James University Hospital, Leeds, UK

## Abstract

**Objective:** Digital pathology is today a widely used technology, and the digitalization of microscopic slides into whole slide images (WSIs) allows the use of machine learning algorithms as a tool in the diagnostic process. In recent years, "deep learning" algorithms for image analysis have been applied to digital pathology with great success. The training of these algorithms requires a large volume of high-quality images and image annotations. These large image collections are a potent source of information, and to use and share the information, standardization of the content through a consistent terminology is essential. The aim of this project was to develop a pilot dataset of exhaustive annotated WSI of normal and abnormal human tissue and link the annotations to appropriate ontological information. **Materials and Methods:** Several biomedical ontologies and controlled vocabularies were investigated with the aim of selecting the most suitable ontology for this project. The selection criteria required an ontology that covered anatomical locations, histological subcompartments, histopathologic diagnoses, histopathologic terms, and generic terms such as normal, abnormal, and artifact. WSIs of normal and abnormal tissue from 50 colon resections and 69 skin excisions, diagnosed 2015-2016 at the Department of Clinical Pathology in Linköping, were randomly collected. These images were manually and exhaustively annotated at the level of major subcompartments, including normal or abnormal findings and artifacts. **Results:** Systemized nomenclature of medicine clinical terms (SNOMED CT) was chosen, and the annotations were linked to its codes and terms. Two hundred WSI were collected and annotated, resulting in 17,497 annotations, covering a total area of 302.19 cm$^2$, equivalent to 107,7 gigapixels. Ninety-five unique SNOMED CT codes were used. The time taken to annotate a WSI varied from 45 s to over 360 min, a total time of approximately 360 h. **Conclusion:** This work resulted in a dataset of 200 exhaustive annotated WSIs of normal and abnormal tissue from the colon and skin, and it has informed plans to build a comprehensive library of annotated WSIs. SNOMED CT was found to be the best ontology for annotation labeling. This project also demonstrates the need for future development of annotation tools in order to make the annotation process more efficient.

**Keywords:** Annotation, digital pathology, image database, ontology, whole slide images

## INTRODUCTION

Digital pathology is today a widely used technology and includes whole slide imaging, "virtual microscopy," which involves the scanning of glass slides to create whole slide images (WSIs), high resolution images, which can be viewed on screen, annotated through computer-based annotation tools, and/or analyzed by computer-based image analysis tools.[1-3] WSIs, WSI markups, and WSI annotations can be integrated into databases and accessed through a local intranet or the internet for primary diagnosis, quality assurance, consultation, teaching, research, and image analysis.[2-4]

**Address for correspondence:** Dr. Karin Lindman, Department of Clinical Pathology, The University Hospital, 58185 Linköping, Sweden. E-mail: karin.lindman@regionostergotland.se

### Access this article online

**Quick Response Code:**

**Website:**
www.jpathinformatics.org

**DOI:**
10.4103/jpi.jpi_81_18

The digitalization of histopathology slides to WSIs allows the use of machine learning algorithms as a tool in the diagnostic process to make a more precise assessment of findings, for example, quantification of immunohistochemistry findings, nuclei detection, gland segmentation, or identification (ID) of other morphological features.[5-8]

In recent years, "deep learning" algorithms for image analysis have been applied to digital pathology with great success. However, the training of these algorithms requires a large volume of high-quality images and image annotations.[8,9] These large image collections are a very potent source of information, and to use, reuse, and share the information, standardization of the content through a consistent terminology is essential.[10-13]

In radiology, where the digitalization of images today is standard, large and annotated image datasets exist, like the lung image database consortium and the image database resource initiative.[12,14] In biomedicine, an example of a large and annotated image database is the human protein atlas.[15-17]

In the histopathology area, there are examples of large image datasets. The International Society of Urological Pathology has established a reference image database of representative images of several pathological entities in kidney, urinary bladder, and prostate.[18] Kostopoulos *et al.*'s group has built an image collection library covering the brain, breast, and laryngeal tumors.[19] The University of Leeds has developed an extensive and expanding database of pathology WSI.[20] However, most of these image databases do not cover the entire WSI, have no images of normal tissue, are rarely annotated, and if image annotations occur, they refer to the quality and content of the image, rather than the different tissue structures.

To the best of our knowledge, there are no existing large annotated image databases of different tissues and organs in histopathology today even though Royal Philips and LABPON have announced their plan to create a digital database of annotated pathology images.[1,21]

For the image annotations to be useful, data have to be coupled to said annotations to provide information related to them. One important challenge is to have a uniform system of nomenclature coupled to the annotations, creating homogeneous and easily reproducible information, allowing others to contribute to or continue the annotation process.

Ontologies are an example of systematic and consistent nomenclatures; they are structured vocabularies consisting of terms designed to represent the type of entities in the domain of reality that each ontology has been devised to capture. These terms are organized hierarchically, ordered by subtype relations.[22,23] In medicine, many different ontologies and controlled vocabularies exist and are evolving: the International Classification of Disease (ICD), Systemized nomenclature of medicine clinical terms (SNOMED CT), Generalized Architecture for Languages Encyclopedias and Nomenclature in Medicine (GALEN), medical subject headings (MeSH), Foundational Model of Anatomy Ontology (FMA), Unified Medical Language System (UMLS), the open biomedical ontologies (OBO), National Cancer Institute Thesaurus (NCIt), and so on.[24-27] An example of a more specific diagnostic ontology is the well-known radiology ontology RadLex.[28,29] To the best of our knowledge, there are no well-known and specific ontology in the histopathology area, although Quantitative Histopathology Image Ontology (QHIO) is under development. QHIO is an ontology covering terms representing the different types and subtypes of histopathological images, imaging processes and techniques, and computational algorithms.[22,30]

To date, machine learning development has focused on specific disease, abnormality, or simple quantification of immunohistochemical stains. The image data has consisted of limited, manually selected regions of WSI or tissue microarrays rather than exhaustive WSI annotations.[5,6,8,14,31] These limited regions do not provide complete information when compared to the information given by the pathologists while examining microscope slides or WSIs. In this context, exhaustive WSI annotations, where all pixels of the WSI will be included in the annotations, could be more useful.

As far as we know, no publicly available dataset of exhaustive annotated WSIs of normal human tissue exists, even though examination of normal tissue is a large and time-consuming part of the histopathological analytic process. To develop machine learning algorithms for diagnosing and classifying different types of tissue, a large database of WSIs of different types of normal and abnormal human tissue will be required.

The aim of this project is to develop a pilot dataset of exhaustively annotated WSIs of normal and abnormal human tissue and link the annotations produced by this process to appropriate ontological information.

## MATERIALS AND METHODS

### Ontology investigation

A systematic search of different biomedical ontologies and controlled vocabularies was made at the BioPortal webpage by a specialist in clinical pathology (KL): http://bioportal.bioontology. org/. This webpage is a comprehensive repository of biomedical ontologies and is provided by the National Center for Biomedical Ontology (NCBO). The goal of NCBO is to support biomedical researches by providing online tools and a web portal enabling them to access, review, and integrate ontological resources.

The ontologies and controlled vocabularies were investigated at the BioPortal webpage, and a PubMed search was also made, to examine their content and structure. The goal of this research was to find the most suitable ontology for the project's purpose, and the selection criteria required an ontology that covered anatomical locations, histological subcompartments, histopathologic diagnoses, histopathologic terms, and generic terms such as normal, abnormal, and artifact.

### Collection of cases

To decide which tissues or organs to be chosen for annotation, the specialist in clinical pathology (KL) made annotation

suggestions of different organs and tissue types: colon, bladder, bone, breast, bronchus, ductus deferens, lung, ileum, liver, lymph node, pancreas, prostate, salivary gland, seminal vesicles, skin, spleen, stomach, thyroid, and uterus. These suggestions and annotations were discussed with a consultant pathologist (DT), and decisions were made by consensus.

Colon and skin were chosen because of their well-defined and histological layered structures, making them very suitable for exhaustive and reproducible annotations. The colon cases were randomly collected from colon resections diagnosed at the Department of Clinical Pathology in Linköping in the year 2015. Small resections of adenomas in the colon were excluded. The skin cases were randomly collected from skin excisions, including pouches, diagnosed at the Department of Clinical Pathology in Linköping in the year 2016. Normal skin excisions and skin excisions diagnosed with neoplasia were included.

The number of 200 WSIs was decided to be enough for the study objective, related to the time and effort taken in the creation of manual and exhaustive annotations. One hundred and one WSIs from the colon and 99 WSIs from the skin were collected.

To make the collection random, colon and skin cases ending with 1, 5, or 8 in their clinical case ID number were chosen.

In cases with both normal and abnormal tissue, one WSI of each type was chosen. In cases with normal or abnormal tissue, one WSI from each case was chosen. The chosen WSI had the best quality, i.e., the least artifacts. The WSIs were manually selected by the specialist in clinical pathology (KL).

## Staining, scanning, image retrieval, and workstation

All of the slides were stained with hematoxylin and eosin stain and scanned by Scanscope AT (Aperio, US), NanoZoomer XR (Hamamatsu, Japan), or NanoZoomer XRL (Hamamatsu, Japan) at a resolution equivalent to 20 times magnification (approximately 0.5 microns per pixel). All of the cases were viewed and selected in Sectra workstation IDS7 Px (Sectra, Sweden), and the WSIs were retrieved from the digital image archive in the clinical pathology picture archiving and communication system (PACS) at the Department of Clinical Pathology in Linköping during 2016–2017.

All the annotations were made with the Sectra workstation IDS7 Px (Sectra, Sweden) and stored in the Sectra IDS PACS system. The computer screen used was an EIZO RadiForce RX850 monitor (EIZO, Japan), and the annotations were made on a Wacom Cintiq 27QHD Touch display (Wacom, Germany) with a Wacom Pro Pen (Wacom, Germany) [Figure 1]. Each WSI from colon covered one tissue level. In skin, each WSI covered 1–6 tissue levels, but only one level was annotated in each WSI.

## Annotator contribution

The colon cases were annotated by the specialist pathologist (KL). The skin cases were annotated by a research assistant (JR),
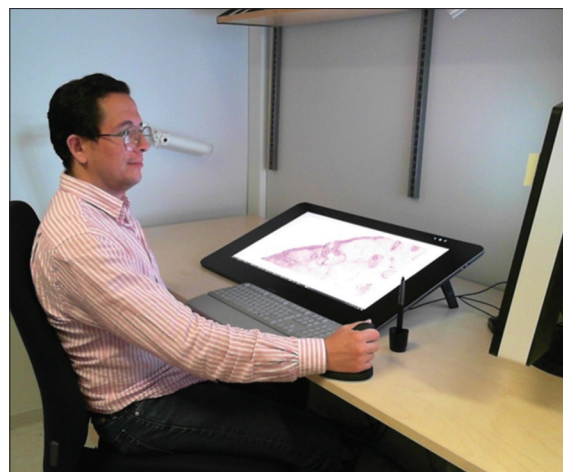


**Figure 1:** Workstation. The complexity of the annotation process and workflow required an adequate and ergonomic workstation, which was composed of several tools, including a high-resolution screen, high-resolution touch screen, precision pen, computer, ergonomic mouse, keyboard, chair, and table. The height of the table and the chair were able to be adjusted to the proportions of the person doing the annotations

after an initial training in the annotation of WSI (supervised by both KL and DT). The specialist pathologist (KL) had regular follow-ups with the consultant pathologist (DT) and research assistant (JR) regarding the annotation procedure and annotation rules. When annotation difficulties occurred (e.g., how to annotate different findings such as structures and abnormalities), the annotator made suggestions of how to annotate the area with difficulties. In the colon cases, the specialist pathologist (KL) and the consultant pathologist (DT) made decisions of annotation rules by consensus. In the skin cases, the specialist pathologist (KL), the consultant pathologist (DT), and the research assistant (JR) made decisions of annotation rules by consensus.

## Annotation rules

The primary aim was to identify the predominant tissue patterns, discerned by a human observer, in each WSI. The annotations were supposed to include 50% normal and 50% abnormal areas. The annotations also covered a range of appearance of the same tissue types (e.g., dark and light staining). Background pixels such as white areas (glass) were not annotated.

In the exhaustive annotation strategy, each pixel on the tissue image was allocated to a nonoverlapping class, the annotations delineated morphologically different subcompartments, and the entire tissue image was annotated. Annotations were made in as a low magnification as possible, but still enough to delineate major anatomic subcompartments in the tissue.

The major subcompartments annotated in the colon WSIs were mucosa, submucosa, muscularis propria, and fatty tissue. The major subcompartments annotated in the skin WSIs were epidermis, dermis, adnexal structures, and subcutaneous fatty tissue. In the skin, cartilage tissue in excisions from the ear also was annotated.

In both the colon and skin, the subcompartments were annotated as normal or abnormal. If abnormal, the abnormality was specified. Artifacts were also annotated. Normal tissue was defined as tissue including expected structures without neoplasia, fibrosis, edema, hemorrhage, or inflammation. Abnormal tissue was defined as morphological abnormal looking tissue (e.g., diverticular mucosa may be a disease, but it is morphological normal when it comes to many of the subcompartments of the lesion).

If an abnormality involved multiple subcompartments and the subcompartment borders could not be morphologically distinguished, the abnormality was annotated as a whole, but the annotation labels included all the subcompartments. If an abnormality involved multiple and easily distinguished subcompartments, each subcompartment was annotated and labeled. Continuous laying lesions and subcompartments were annotated as a whole, and discontinuous laying structures were annotated separately. If a tumor consisted of tumor stroma (e.g., basal cell carcinoma), the stroma was annotated as a part of the tumor. In diffuse lesions, where exact borders and tumor cells were hard to distinguish, the borders were defined as the region where the normal tissue started/ended. In lesions with abnormal architecture but normal cell morphology, the annotations were labeled as normal. Tissue folds, focal thick areas because of tissue preparation, and separately lying tissue parts were annotated as artifacts.

## Annotation strategy

Image annotations were made in a systematic way: the first step in the annotations process was to identify the parts of the tissue and structures that were to be annotated. Then, the artifacts were annotated, both those derived as a result of specimen preparation as well as those resulting from the scanning process. After that, the annotations of the tissue itself were done from the epithelial side to the innermost layers, and from left to right. For the colon WSI, the process started from the mucosa to the serosa/adventitia, doing individual annotations for each one of the layers, starting with the abnormalities found in the tissue, and then continuing with the normal colonic tissue. In the case of the skin WSI, the process started with the epidermis layer being annotated, followed by the tumor or abnormality present and then the adnexal structures were annotated. After that, if there was inflammatory tissue present, this was annotated; otherwise, the dermis area was annotated as a whole, excluding areas with abnormal tissue present and differentiating between areas with inked margins as separate annotations. Finally, the subcutaneous tissue was annotated, also doing the differentiation between areas with inked margins as separate annotations. During the annotation process, all the pixels of the WSI that contained tissue were included, and the overlay and crossing of different annotations was avoided (except for things that could be covered in the tissue layers, e.g., adnexal structures of the skin that can normally be found in the dermis), to avoid conflicting annotations [Figures 2 and 3].



**Figure 2:** Annotated whole slide images from the colon. The annotations were done focusing on different subcompartments found in the colon, as well as pointing out abnormalities if present

## Annotation labeling rules

Terms and codes for organ, anatomic location and subcompartment were taken from the SNOMED CT hierarchy "Body structure" class. The most specific SNOMED CT concept and code was used for organ, anatomic location, subcompartments, abnormalities and disease (e.g. descending colon instead of colon, submucosa of colon instead of submucosa etc.). The normal, abnormal and artifact concepts and codes were also included in the annotation labelling.

## Annotation labeling information

The annotations were stored in the Sectra IDS PACS system. Individual ID numbers were assigned to every individual annotation, and the information linked to each ID number was composed of the different ontology concepts and codes, describing the content of the annotation. All this information was saved in an Excel-file with following label headings: "organ," "SNOMED CT code organ," "sub-compartment," "SNOMED CT code sub-compartment," "SNOMED CT code combined organ and sub-compartment," "normal/abnormal including specific abnormalities," and "SNOMED CT codes for normal/abnormal including specific abnormalities." A link to the skin WSI in the software where the annotations were made also coupled with this information, to offer rapid accessibility.

## Results

From the ontology search at the BioPortal webpage, four biomedical "ontologies" were found out of a total number of 314 ontologies. These four ontologies are listed in Table 1: FMA, NCIt, MeSH, and SNOMED CT [Appendix A].[24-27,32-35]

SNOMED-CT was the only ontology fulfilling all the selection criteria and was chosen as the most suitable ontology for the project's purpose; it is a well-known, used, and evolving ontology with anatomical, histological, and pathological concepts. Disadvantages with SNOMED CT are as follows: it lacks formality and it is a clinical ontology and do not cover the laboratory process or imaging technology in the histopathology area.
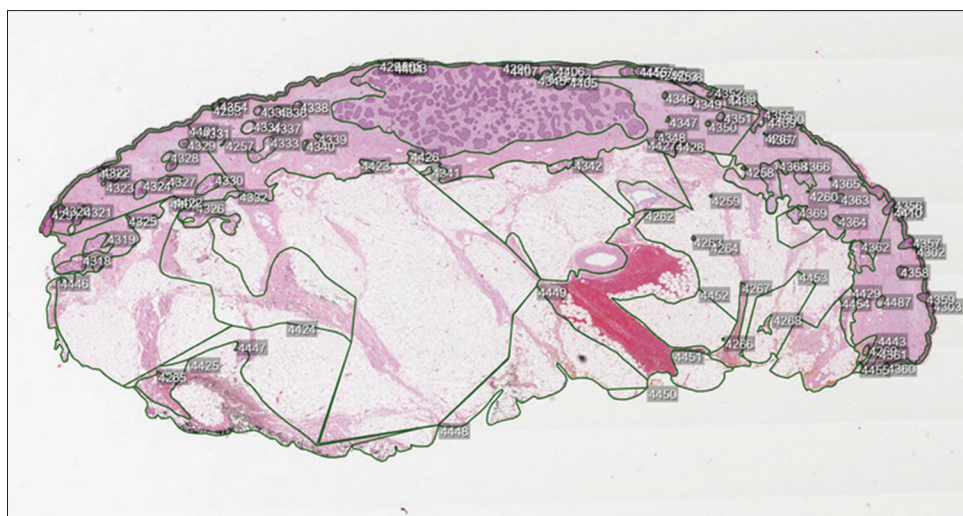
**Figure 3:** Annotated whole slide images from the skin. The annotations were done focusing on different subcompartments found in the skin, as well as pointing out abnormalities if present. When comparing the annotations from colon and skin whole slide images, the skin annotations were more complicated to perform

**Table 1: Summary of pros and cons of each ontology**

| Ontology | Pros | Cons |
|---|---|---|
| NCIt | Well-known and evolving reference terminology/biomedical ontology<br>Good coverage for cancer | Not a formal ontology<br>Do not cover all other abnormalities except cancer<br>Lack some morphological structures<br>No coverage for the pathology laboratory process |
| FMA | Well-known and evolving reference ontology<br>Good coverage for anatomical structures | No coverage for disease, terms as abnormal, normal, or artifact<br>No coverage for the pathology laboratory process |
| MeSH | Well-known and evolving vocabulary thesaurus<br>Good coverage for diseases, anatomy, and tissue subcompartments | No ontology<br>No coverage for all concepts needed<br>No logical hierarchy when it comes to anatomy and tissue sub compartments<br>Low coverage for the pathology laboratory process |
| SNOMED CT | Well-known and evolving hierarchy of concepts/ontology<br>Includes all concepts needed | Low formality<br>Low coverage for the pathology laboratory process |

SNOMED CT: Standardized nomenclature of medicine clinical terms, MeSH: Medical subject headings, FMA: Foundational model of anatomy, NCIt: National Cancer Institute thesaurus

The majority of the SNOMED CT concepts were taken from the hierarchy "body structure" [Figure 4] class except normal, abnormal ("qualifier value" class), and artifact ("clinical finding" class).

Fifty cases from different parts of the colon [Appendix B] were collected, giving 101 WSIs. Of these, 49 WSIs were assessed as normal and 52 as abnormal. A total of 756 annotations were made from colon, and 39 unique SNOMED CT concepts and codes were used [Appendix B].

Sixty-nine cases from skin from different parts of the body were collected, giving 99 WSIs, of which 50 were assessed as normal and 49 as abnormal. As for skin, a total of 16,741 annotations were made, and 56 unique SNOMED CT concepts and codes were used [Appendix C].

A total area of 302.19 cm², 127.7 giga pixel, was annotated.

In the colon, the highest magnifications used was 6 times magnification. In the skin, 40 times magnification, this because of the small adnexal structures in the skin. The smallest annotated finding in the colon measured 300 microns in largest diameter, and in the skin, 14 microns in largest diameter.

Each slide took from 45 s to 360 min to annotate. The total time taken to annotate all of the WSIs was approximately 360 h. The full summary of results can be seen in Table 2.

Digital object identifiers (DOIs) have been assigned to the datasets, for future reference.[36,37] The data (WSIs and annotation information in JSON format) are being shared within the AIDA consortium.[38] The dataset is not publicly available because of data regulation laws, but inquiries in access can be directed to AIDA management.
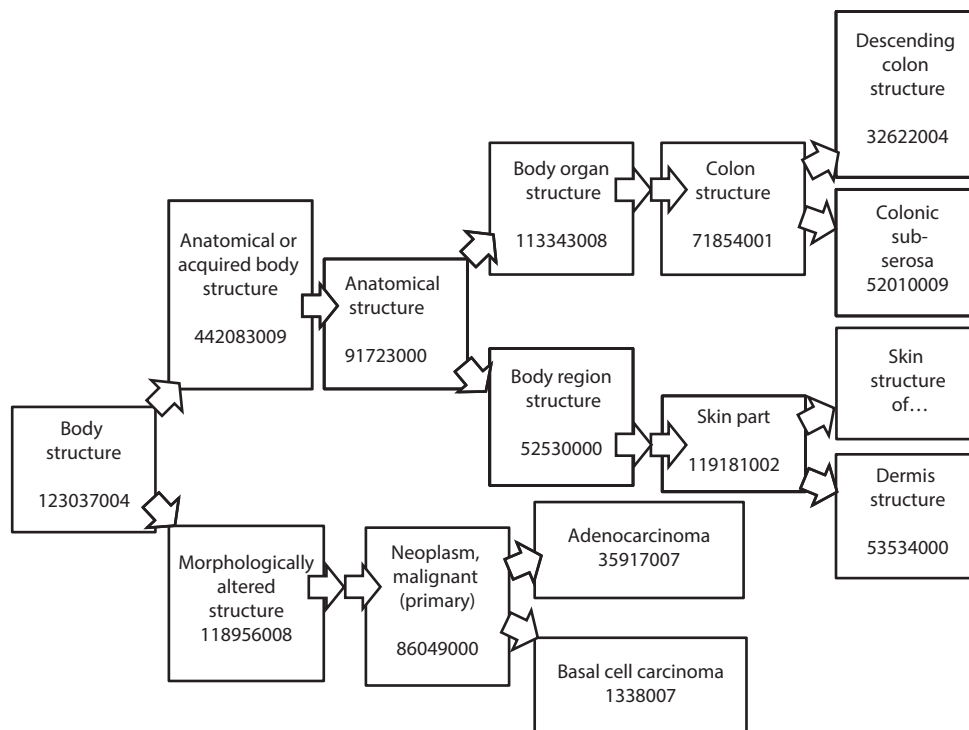
**Figure 4:** Image branch of the body structure hierarchy in systemized nomenclature of medicine clinical terms (SNOMED CT) including the systemized nomenclature of medicine-clinical terms' codes. The different structures and diagnoses are ordered in a structured way, originating from the "body structure" concept

**Table 2: Summary of results**

| Result type | Value |
| --- | --- |
| Total number of WSIs | 200 |
|    Normal WSIs | 99 |
|    Abnormal WSIs | 101 |
| Total number of SNOMED CT codes | 95 |
|    Anatomical locations | 37 |
|    Subcompartments | 10 |
|    Specific abnormality codes | 44 |
| Total number of annotations | 17497 |
| Total area annotated | 302.19 cm$^2$ 127.7 GP |
| Total time taken to annotate | 360 h (approximately) |

SNOMED CT: Standardized nomenclature of medicine clinical terms, WSI: Whole slide images

## DISCUSSION

The aim of this project was to develop a pilot dataset of exhaustively annotated WSIs of normal and abnormal human tissue and link the annotations produced by this process to appropriate ontological information. Our dataset is scalable, both in terms of adding new domains (diseases, tissues, etc.,) and scaling to generate a high-quality dataset. We strongly believe in the necessity of standardization of annotation and annotation labeling, in creating reusable processes to generate high-quality large scale annotations, which will enable future reusability and interoperability. We think this work has contributed to that.

Our results show the distribution of anatomical locations, diagnoses, and annotations from the colon and skin in our dataset.

The results also show a significant difference in the annotation number between colon and skin and in the time taken to annotate a WSI. These findings can be explained by the difference in architectural complexity between a WSI from the colon, with simple and layered subcompartment architecture and adnexal-rich skin, with more complex subcompartments architecture due to numerous and small adnexal structures.

To effectively annotate the important parts of the tissues, a structured methodology had to be implemented in the annotation process. The most effective way that the research group found was to start from the outermost layers of the tissues and then going into the innermost layers, and from the left side of the tissue to the right one. The ordered pattern in which the annotation process was performed allowed a fluent workflow, which at the same time avoided overlay and crossing of different annotations. Also by going from left to right and from the epithelium to the outer most layers, the process was faster since the location of every annotation would lead to the location of next annotation, improving the workflow even further. This highlights the importance of developing efficient workflow strategies and their implementation when creating exhaustive and detailed annotation databases.

To find, use, and share information in an annotated image database, it is essential to standardize the content.[2-4] In this study, we examined ontologies and structured vocabularies for this purpose, and found SNOMED CT to be the best ontology to use, since it covers both anatomical locations, histological subcompartments, histopathologic diagnoses,

histopathologic terms and generic terms such as normal, abnormal and artifact. SNOMED CT was originally created by the College of American Pathologists and is the world's largest clinical terminology with broad coverage of clinical medicine including disease and phenotypes. SNOMED CT is a class hierarchy of concepts which includes high-level categories such as body structure, clinical findings, and so on. Within each hierarchy, the concepts are organized from the general to the more detailed concept through "is-a" relationships. The concepts can also be linked to different hierarchies through attribute relationships. SNOMED CT is not a specific histopathology ontology, but we could find all the terms we needed for this project, and SNOMED CT is well-known and widely used, which means it is constantly improving, evolving, and with pathology input can develop and be more suitable for the histopathology area.[5,6] In addition, SNOMED CT terms do not cover the histopathological laboratory process or histopathological imaging technology which is of importance in sharing annotated datasets. To date, QHIO is the only ontology covering terms representing the different types and subtypes of histopathological images, imaging processes and techniques, and computational algorithms, although it is not yet ready for use.[30] We think a future combination of SNOMED CT and QHIO would be of great benefit.

One significant application of machine learning tools in the future could be to screen normal slides of tissue to minimize manual work, allowing pathologists to have more time to focus on the diagnosis and classification of disease. To train machine learning algorithms for image analysis software, a large amount of data is required to achieve a high accuracy rate; image databases like this one, which includes a high amount of exhaustive and detailed annotations, can be useful. Another important function is for educational purposes, that is, for both helping students and trainees in learning histology and histopathology.

This study had some challenges and limitations. One of the main challenges for the creation of the dataset was the time taken to annotate. The annotations were made manually, and complicated WSIs of adnexal-rich skin took a long time to be accurately annotated since they contain additional small structures that are not found in colon tissue. We believe that automatic to semiautomatic annotation tools have an important role in a more efficient annotation workflow, and we will investigate this in future research.

We had a small dataset of 200 WSI. One hundred and one WSIs were abnormal with a wide variety of diagnoses; many only covered by one WSI. Our study was a pilot project, exhaustive annotation takes time and effort, and our annotations were general and not made for specific machine learning tasks. 100 WSI from each organ were decided to be enough for the objective of this pilot project. Despite our rather small dataset, we believe that the images and annotations can be useful for others to contribute to or combine with other datasets.

Our lack of a specific machine learning task when creating our annotation rules made it difficult to define which level of detail the annotations should cover. We decided to aim at major subcompartments in the colon and skin and to annotate in as low a magnification as possible, due to the effort and time exhaustive annotation entails. The low magnification made it difficult to cover all pixels of the image and to do nonoverlapping annotations.

Each WSI was annotated by one annotator, and consensus annotations from multiple annotators would have enabled more objective and high-quality annotations.

Another challenge in the annotation process was the annotation of diffuse lesions, such as for example, dermatofibroma, where the exact extent of the lesion was hard to define. We delineated these lesions as precisely as we could, with lesion on one side and normal tissue on the other. Lesions involved by diffuse inflammation and where the borders were hard to distinguish, were also difficult to annotate. In these cases, we annotated the lesion as precisely as we could, including part of the inflammation. The inflammation outside the lesion was annotated separately.

An uncertainty in the annotation process was also the decision if some tissue changes were normal or abnormal, for example, reactive changes. These decisions were made by consensus referred to the annotator procedure in the method.

## Conclusion

This pilot project has resulted in a dataset of 200 exhaustive annotated WSI of normal and abnormal human tissues from the colon and skin. The project illustrates the process of building an exhaustive annotated dataset of WSIs. It also illustrates the usage of systemized nomenclatures in the labelling of the annotations, with the aim of facilitating future contribution to, and sharing of the annotated image data. The 200 gathered WSI from the colon and skin resulted in 17,497 ontology-linked annotations, covering anatomical location, histological subcompartments, normal/abnormal tissue, and more specific diagnoses as well as tissue abnormalities. SNOMED CT proved to be the best ontology for the objective of this project.

This work has informed plans to build a comprehensive library of annotated WSIs. This work also shows the need for future development of annotation tools to do the annotation process faster and more efficient.

### Financial support and sponsorship

### Conflicts of interest
There are no conflicts of interest.

### References
1. Evans AJ, Salama ME, Henricks WH, Pantanowitz L. Implementation of whole slide imaging for clinical purposes: Issues to consider from the perspective of early adopters. Arch Pathol Lab Med 2017;141:944-59.

2.  Pantanowitz L, Valenstein PN, Evans AJ, Kaplan KJ, Pfeifer JD, Wilbur DC, *et al.* Review of the current state of whole slide imaging in pathology. J Pathol Inform 2011;2:36.

3.  Webster JD, Dunstan RW. Whole-slide imaging and automated image analysis: Considerations and opportunities in the practice of pathology. Vet Pathol 2014;51:211-23.

4.  Bueno G, Fernández-Carrobles MM, Deniz O, García-Rojo M. New trends of emerging technologies in digital pathology. Pathobiology 2016;83:61-9.

5.  Veta M, Pluim JP, van Diest PJ, Viergever MA. Breast cancer histopathology image analysis: A review. IEEE Trans Biomed Eng 2014;61:1400-11.

6.  Mosquera-Lopez C, Agaian S, Velez-Hoyos A, Thompson I. Computer-aided prostate cancer diagnosis from digitized histopathology: A review on texture-based systems. IEEE Rev Biomed Eng 2015;8:98-113.

7.  Griffin J, Treanor D. Digital pathology in clinical use: Where are we now and what is holding us back? Histopathology 2017;70:134-45.

8.  Litjens G, Kooi T, Bejnordi BE, Setio AAA, Ciompi F, Ghafoorian M, *et al.* A survey on deep learning in medical image analysis. Med Image Anal 2017;42:60-88.

9.  Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, Blau HM, *et al.* Dermatologist-level classification of skin cancer with deep neural networks. Nature 2017;542:115-8.

10. Turner JA, Mejino JL, Brinkley JF, Detwiler LT, Lee HJ, Martone ME, *et al.* Application of neuroanatomical ontologies for neuroimaging data annotation. Front Neuroinform 2010;4. pii: 10.

11. Traore L, Kergosien Y, Racoceanu D. Bridging the semantic gap between diagnostic histopathology and image analysis. Stud Health Technol Inform 2017;235:436-40.

12. Armato SG 3rd, McLennan G, Bidaut L, McNitt-Gray MF, Meyer CR, Reeves AP, *et al.* The lung image database consortium (LIDC) and image database resource initiative (IDRI): A completed reference database of lung nodules on CT scans. Med Phys 2011;38:915-31.

13. Haug PJ, Ferraro JP, Holmen J, Wu X, Mynam K, Ebert M, *et al.* An ontology-driven, diagnostic modeling system. J Am Med Inform Assoc 2013;20:e102-10.

14. Kovacheva VN, Snead D, Rajpoot NM. A model of the spatial tumour heterogeneity in colorectal adenocarcinoma tissue. BMC Bioinformatics 2016;17:255.

15. Wang W, Luo J, Yang X, Lin H. Data analysis of the lung imaging database consortium and image database resource initiative. Acad Radiol 2015;22:488-95.

16. Pontén F, Schwenk JM, Asplund A, Edqvist PH. The human protein atlas as a proteomic resource for biomarker discovery. J Intern Med 2011;270:428-46.

17. Navani S. Manual evaluation of tissue microarrays in a high-throughput research project: The contribution of Indian surgical pathology to the human protein atlas (HPA) project. Proteomics 2016;16:1266-70.

18. Egevad L, Cheville J, Evans AJ, Hörnblad J, Kench JG, Kristiansen G, *et al.* Pathology imagebase-a reference image database for standardization of pathology. Histopathology 2017;71:677-85.

19. Kostopoulos S, Ravazoula P, Asvestas P, Kalatzis I, Xenogiannopoulos G, Cavouras D, *et al.* Development of a reference image collection library for histopathology image processing, analysis and decision support

20. Slide Library, Virtual Pathology at the University of Leeds. Available from: http://www.virtualpathology.leeds.ac.uk/slides/library/. [Last accessed on 2019 Jan 11].

21. Philips and LabPON Plan to Create World's Largest Pathology Database of Annotated Tissue Images for Deep Learning. Available from: https://www.philips.com/a-w/about/news/archive/standard/news/ᵗpress/2017/20170306-philips-and-labpon-plan-to-create-worlds-largest-pathology-database-of-annotated-tissue-images-for-deep-learning.html. [Last accessed on 2019 Jan 11].

22. Smith B, Arabandi S, Brochhausen M, Calhoun M, Ciccarese P, Doyle S, *et al.* Biomedical imaging ontologies: A survey and proposal for future work. J Pathol Inform 2015;6:37.

23. Mabotuwana T, Lee MC, Cohen-Solal EV. An ontology-based similarity measure for biomedical data-application to radiology reports. J Biomed Inform 2013;46:857-68.

24. Lipscomb CE. Medical subject headings (MeSH). Bull Med Libr Assoc 2000;88:265-6.

25. de Coronado S, Wright LW, Fragoso G, Haber MW, Hahn-Dantona EA, Hartel FW, *et al.* The NCI thesaurus quality assurance life cycle. J Biomed Inform 2009;42:530-9.

26. Rosse C, Mejino JL Jr. A reference ontology for biomedical informatics: The foundational model of anatomy. J Biomed Inform 2003;36:478-500.

27. García-Rojo M, Daniel C, Laurinavicius A. SNOMED CT in pathology. Stud Health Technol Inform 2012;179:123-40.

28. Spanier AB, Cohen D, Joskowicz L. A new method for the automatic retrieval of medical cases based on the RadLex ontology. Int J Comput Assist Radiol Surg 2017;12:471-84.

29. Wang KC, Kohli M, Carrino JA. Technology standards in imaging: A practical overview. J Am Coll Radiol 2014;11:1251-9.

30. Gurcan MN, Tomaszewski J, Overton JA, Doyle S, Ruttenberg A, Smith B, *et al.* Developing the quantitative histopathology image ontology (QHIO): A case study using the hot spot detection problem. J Biomed Inform 2017;66:129-35.

31. Kothari S, Phan JH, Stokes TH, Wang MD. Pathology imaging informatics for quantitative analysis of whole-slide images. J Am Med Inform Assoc 2013;20:1099-108.

32. Medical Subject Headings – Home Page. Available from: https://www.nlm.nih.gov/mesh/meshhome.html. [Last accessed on 2019 Jan 11].

33. SNOMED Home page. Available from: https://www.snomed.org. [Last accessed on 2019 Jan 11].

34. Foundational Model of Anatomy, Structural Informatics Group. Available from: http://www.si.washington.edu/projects/fma. [Last accessed on 2019 Jan 11].

35. NCI Thesaurus. Available from: https://www.ncit.nci.nih.gov/ncitbrowser/. [Last accessed on 2019 Jan 11].

36. Lindman K, Rose JF, Lindvall M, Stadler CB, Lundström C, Treanor D. Colon data from the Visual Sweden project DROID. Analytic Imaging Diagnostics Arena (AIDA). 2019. https://doi.org/10.23698/aida/drco

37. Lindman K, Rose JF, Lindvall M, Stadler CB, Lundström C, Treanor D. Skin data from the Visual Sweden project DROID. Analytic Imaging Diagnostics Arena (AIDA). 2019. https://doi.org/10.23698/aida/drsk

38. Svensson K. AIDA Overview (English) – Medtech4Health. Available from: https://medtech4health.se/aida-en/about-aida/. [Last accessed on 2019 Jan 25].

# APPENDICES

## APPENDIX A

### ONTOLOGIES

### National Cancer Institute Thesaurus

The National Cancer Institute Thesaurus (NCIt) is a reference terminology and biomedical ontology used by the NCI. It was first published in 2000. The purpose was to facilitate interoperability and data sharing by the various components of NCI. NCIt is updated on a regular basis.

NCIt covers vocabulary for clinical care, translational and basic research, public information, and administrative activities. The content is focused on cancer but also contains terminology not specific to cancer.

NCIt is a concept-based terminology, with 70,000 concepts hierarchically organized in 19 distinct domains. Each concept in the hierarchy has one or more is-a relation. It was not created as a pure ontology. NCIt does not cover all domains exhaustively, except in certain areas such as cancer diseases and cancer drugs.

The ontology has a good cover for cancer research domains, especially cancer-related drugs and chemotherapy regimens, and moderate coverage for more general health care research. The degree of formality varies by area of the terminology.[1,2]

+ Well-known and evolving reference terminology, good coverage for cancer.

− Not a formal ontology, do not cover other abnormalities except cancer, lack some morphological structures, no coverage for the pathology laboratory process.

### Foundational model of anatomy

The foundational model of anatomy (FMA) is developed and maintained by the Structural Informatics Group at the University of Washington. The ontology is a reference ontology covering human anatomy used in education and biomedical research. FMA is an evolving ontology for biomedical informatics and is concerned with "the representation of entities and relationships necessary for the symbolic modeling of the structure of the human body in a computable form that is also understandable by humans." All its classes are linked in a hierarchy way to a single root: Anatomical entity.[3,4]

+ Well-known and evolving reference ontology, good coverage for anatomical structures.

− No coverage for disease, terms as abnormal, normal, or artifact, no coverage for the pathology laboratory process.

### Medical subject heading

Medical subject headings (MeSH) terms were introduced at the National Library of Medicine (NLM) 1960. MeSH terms are the NLM's controlled terminology, primarily used to organize and index information and manuscripts found in common databases such as PubMed. It consists of over 87,000 entry terms naming descriptors in a hierarchical structure.

MeSH is not an ontology. Its concepts are not classes, and its hierarchical links are not subclass relations.[5,6]

+ Well-known and evolving vocabulary thesaurus, good coverage for disease, anatomy, and tissue subcompartments.

− No ontology, do not include all concepts needed, no logical hierarchy when it comes to anatomy and tissue subcompartments, low coverage for the pathology laboratory process.

### Standardized nomenclature of medicine clinical term

Standardized nomenclature of medicine clinical terms (SNOMED CT) was originally created by the College of American Pathologists. Since 2007, it is owned, maintained, and distributed by the International Health Terminology Standard Development Organization.

It is the world's largest clinical terminology and provides broad coverage of clinical medicine, including diseases and phenotypes. It is used in electronic health care records, clinical user interface, decision support system, knowledge access systems, and natural language processing. SNOMED CT is a class hierarchy of concepts which includes high-level categories such as body structure, clinical findings, and so on. Within each hierarchy, the concepts are organized from the general to the more detailed concept through "is-a" relationships. The concepts can also be linked to different hierarchies through attribute relationships. This structure allows each concept to have multiple relationships and zero to multiple synonyms.[7,8]

+ Well-known and evolving hierarchy of concepts/ontology, includes all concepts needed: anatomy, subcompartments, disease, and terms such as normal, abnormal, and artifact.

− Low formality, low coverage for the pathology laboratory process.

## REFERENCES

1. de Coronado S, Wright LW, Fragoso G, Haber MW, Hahn-Dantona EA, Hartel FW, *et al.* The NCI thesaurus quality assurance life cycle. J Biomed Inform 2009;42:530-9.
2. NCI Thesaurus. Available from: https://www.ncit.nci.nih.gov/ncitbrowser/. [Last accessed on 2019 Jan 11].
3. Rosse C, Mejino JL Jr. A reference ontology for biomedical informatics: The foundational model of anatomy. J Biomed Inform 2003;36:478-500.
4. Foundational Model of Anatomy. Structural Informatics Group. Available from: http://www.si.washington.edu/projects/fma. [Last accessed on 2019 Jan 11].
5. Lipscomb CE. Medical subject headings (MeSH). Bull Med Lib Assoc 2000;88:265-6.
6. Medical Subject Headings – Home Page. Available from: https://www.nlm.nih.gov/mesh/meshhome.html. [Last accessed on 2019 Jan 11].
7. García-Rojo M, Daniel C, Laurinavicius A. SNOMED CT in pathology. Stud Health Technol Inform 2012;179:123-40.
8. SNOMED Home page. Available from: https://www.snomed.org. [Last accessed on 2019 Jan 11].

## APPENDIX B

| SNOMED CT concept | Code |
| --- | --- |
| Abnormal | 263654008 |
| Acute and chronic inflammation | 75889009 |
| Acute inflammation | 4532008 |
| Adenocarcinoma | 35917007 |
| Artifact | 47973001 |
| Ascending colon | 9040008 |
| Atrophy | 13331008 |
| Cecum | 32713005 |
| Chronic inflammation | 84499006 |
| Colon | 71854001 |
| Colonic mucous membrane | 68502009 |
| Colonic muscularis propria | 41948009 |
| Colonic submucosa | 61647009 |
| Colonic subserosa | 52010009 |
| Descending colon | 32622004 |
| Diverticula | 31113003 |
| Diverticulitis | 18126004 |
| Dysplasia | 25723000 |
| Edema | 79654002 |
| Fibrosis | 112674009 |
| Granulation tissue | 61363009 |
| Hemorrhage | 50960005 |
| Hyalinization | 19010006 |
| Hyperplasia | 76197007 |
| Hyperplastic polyp | 62047007 |
| Ileum | 34516001 |
| Inflammation | 23583003 |
| Lymphoma | 21964009 |
| Mucinous adenocarcinoma | 72495009 |
| Necrosis | 6574001 |
| Normal | 17621005 |
| Rectum | 34402009 |
| Serrated adenoma | 128653004 |
| Sigmoid colon | 60184004 |
| Stasis | 19685008 |
| Transverse colon | 485005 |
| Tubular adenoma | 444408007 |
| Tubulovillous adenoma | 61722000 |
| Ulcer | 56208002 |

SNOMED CT: Standardized nomenclature of medicine clinical terms

**Table B1: Number of cases from different parts of colon or colon as a whole and corresponding standardized nomenclature of medicine clinical terms code**

| Part of colon | Number of cases | SNOMED CT code |
| --- | --- | --- |
| Right colon | 16 | 51342009 |
| Transverse colon | 3 | 485005 |
| Left colon | 7 | 55572008 |
| Sigmoid colon | 18 | 60184004 |
| Rectum | 2 | 34402009 |
| Colon unspecified | 4 | 71854001 |

SNOMED CT: Standardized nomenclature of medicine clinical terms

## Table B2: Number of different diagnoses and corresponding standardized nomenclature of medicine clinical terms code

| Diagnosis | Number of cases | SNOMED CT code |
|---|---|---|
| High-grade adenocarcinoma | 15 | 413447005 |
| Low-grade adenocarcinoma | 14 | 413449008 |
| Mucinous adenocarcinoma | 5 | 72495009 |
| Lymphoma | 1 | 21964009 |
| Tubulovillous adenoma | 1 | 61722000 |
| Tubular adenoma | 2 | 19665009 |
| Serrated adenoma | 2 | 128653004 |
| Ulcerative colitis | 1 | 64766004 |
| Morbus crohn | 1 | 34000006 |
| Hyperplastic polyp | 2 | 62047007 |
| Hyperplasia | 1 | 76197007 |
| Diverticulosis | 6 | 397881000 |
| Necrosis | 1 | 6574001 |
| Ulceration/hemorrhage | 1 | 56208002/50960005 |
| Inflammation | 1 | 23583003 |

SNOMED-CT: Standardized nomenclature of medicine clinical terms

# APPENDIX C

| SNOMED CT concept | Code |
|---|---|
| Acanthosis | 23620008 |
| Actinic keratosis | 856006 |
| Basal cell carcinoma | 1338007 |
| Dermatofibroma | 72079004 |
| Dermis | 53534000 |
| Dysplastic nevus | 61814002 |
| Epidermis | 55988001 |
| Fibrosis | 112674009 |
| Fibrin body | 45619005 |
| Granuloma | 45647009 |
| Inflammation | 23583003 |
| Inflammatory edema | 103619005 |
| Intradermal nevus | 112681002 |
| Keratoacanthoma | 416378000 |
| Lentigo maligna melanoma | 44474009 |
| Malignant melanoma | 2092003 |
| Melanocytic nevus | 400101001 |
| Malignant melanoma *in situ* | 77986002 |
| Normal skin | 225544001 |
| Perichondrium | 11881003 |
| Reactive cellular changes | 125513006 |
| Scar | 12402003 |
| Seborrheic keratosis | 25499005 |
| Skin and subcutaneous tissue structure | 127856007 |
| Skin and subcutaneous tissue structure of back | 417286006 |
| Skin and subcutaneous tissue structure of head | 389074000 |
| Skin and subcutaneous tissue structure of trunk | 389072001 |
| Skin appendage structure | 7748002 |
| Skin structure | 39937001 |
| Skin structure of back | 66643007 |
| Skin structure of breast | 82038008 |

*Contd...*

| SNOMED CT concept | Code |
|---|---|
| Skin structure of calf of leg | 51059006 |
| Skin structure of cheek | 36141000 |
| Skin structure of ear | 1902009 |
| Skin structure of eyebrow | 367577003 |
| Skin structure of eyelid and periocular area | 399996007 |
| Skin structure of face | 73897004 |
| Skin structure of female genitalia | 19938000 |
| Skin structure of forehead | 68698007 |
| Skin structure of foot | 60496002 |
| Skin structure of hand | 33712006 |
| Skin structure of head | 70762009 |
| Skin structure of lip | 88089004 |
| Skin structure of lower extremity | 371304004 |
| Skin structure of neck | 43081002 |
| Skin structure of nose | 113179006 |
| Skin structure of scalp | 43067004 |
| Skin structure of scapular region of back | 45980000 |
| Skin structure of shoulder | 76552005 |
| Skin structure of temple | 244081009 |
| Skin structure of thigh | 371305003 |
| Skin structure of upper extremity | 371311000 |
| Structure of cartilage of auditory canal | 83543000 |
| Squamous cell carcinoma | 28899001 |
| Squamous cell carcinoma *in situ* | 59529006 |
| Subcutaneous fatty tissue | 67769002 |
| Subcutaneous tissue | 71966008 |
| Surgical margins | 82868003 |

SNOMED CT: Standardized nomenclature of medicine clinical terms

## Table C1: Number of cases from different skin parts and corresponding standardized nomenclature of medicine clinical terms code

| Part of skin | Number of cases | SNOMED CT code |
| --- | --- | --- |
| Back | 5 | 66643007 |
| Head | 4 | 70762009 |
| Trunk | 3 | 389072001 |
| Breast | 2 | 82038008 |
| Calf of leg | 1 | 51059006 |
| Cheek | 6 | 36141000 |
| Ear | 9 | 1902009 |
| Eyebrow | 1 | 367577003 |
| Eyelid and periocular area | 3 | 399996007 |
| Face | 5 | 73897004 |
| Female genitalia | 1 | 19938000 |
| Forehead | 2 | 68698007 |
| Foot | 1 | 60496002 |
| Hand | 1 | 33712006 |
| Lip | 2 | 88089004 |
| Lower extremity | 6 | 371304004 |
| Neck | 2 | 43081002 |
| Nose | 6 | 113179006 |
| Scapular region of the back | 1 | 45980000 |
| Shoulder | 3 | 76552005 |
| Temple | 5 | 244081009 |
| Thigh | 2 | 371305003 |
| Upper extremity | 4 | 371311000 |

SNOMED CT: Standardized nomenclature of medicine clinical terms

## Table C2: Number of different diagnoses in the skin cases and corresponding standardized nomenclature of medicine clinical terms code

| Diagnosis | Number of cases | SNOMED CT code |
| --- | --- | --- |
| Actinic keratosis | 1 | 856006 |
| Basal cell carcinoma | 31 | 1338007 |
| Dermatofibroma | 2 | 72079004 |
| Dysplastic nevus | 1 | 61814002 |
| Intradermalt nevus | 2 | 112681002 |
| Keratoacanthoma | 1 | 416378000 |
| Lentigo maligna melanoma | 1 | 44474009 |
| Malignant melanoma | 3 | 2092003 |
| Malignant melanoma *in situ* | 1 | 77986002 |
| Scar | 3 | 12402003 |
| Seborrheic keratosis | 2 | 25499005 |
| Squamous cell carcinoma | 3 | 28899001 |
| Squamous cell carcinoma *in situ* | 3 | 59529006 |

SNOMED CT: Standardized nomenclature of medicine clinical terms

*Contd...*