ITU-T FG-AI4H Deliverable

TELECOMMUNICATION STANDARDIZATION SECTOR OF ITU

16 March 2023

PRE-PUBLISHED VERSION

DEL5.3

1-0-L

Data annotation specification



Summary

Data annotation would be one of the most dependable factors on model performance, it serves as one important aspect of data quality control on Artificial Intelligence for health. This Deliverable gives a general guideline of data annotation specification, including definition, background and goals, framework, standard operating procedure, scenario classifications and corresponding criteria, as well as recommended metadata, etc. A questionnaire is attached to seek input and collaboration with topic groups regarding data annotation.

Keywords

Artificial intelligence; health; data handling; data annotation; test data

Change Log

This document contains Version 1 of the Deliverable DEL5.3 on "*Data annotation specification*" approved on 16 March 2023 via the online approval process for the ITU-T Focus Group on AI for Health (FG-AI4H).

Editor:	Shan Xu CAICT, China	E-mail:	E-mail: <u>xushan@caict.ac.cn</u>		
	Sebastian Bosse Fraunhofer HHI, Germany	E-mail:	<u>sebastian.bc</u>	osse@hhi.fraunhofer.de	
	Jianrong Wu Tencent Healthcare, China	E-mail:	<u>edwinjrwu@</u>	encent.com	
Contributors:	(in alphabetical order)				
	Yanwu Xu Intelligent Healthcare Unit, Baidu,	China	E-mail:	xuyanwu@baidu.com	
	Nathan Guo ShuKun Techonology, China		E-mail:	guoning@shukun.net	
	Yajun Zhang Tencent Technology (Shenzhen), C	China	E-mail:	yajunzhang@tencent.com	
	Chunyang Niu Zhejiang Lab, China		E-mail:	niucy@zhejianglab.com	
	Guoqiang Li Zhejiang Lab, China		E-mail:	gli@zhejianglab.com	
	Huihui Fang Artificial Intelligence Innovation B Baidu, China	Business,	E-mail:	<u>fanghuihui@baidu.com</u>	
	Harpreet Singh ICMR, India		E-mail:	hsingh@bmi.icmr.org.in	

CONTENTS

Page

1	Scope				
2	References4				
3	Definition 3.1	ns4 Terms defined elsewhere4			
	3.2	Terms defined in this document			
4	Abbrevia	ations and acronyms5			
5	Convent	ions5			
6	Backgro	und and goals5			
7	Framewo	ork6			
8	Standard	operating procedure			
	8.1	Independent annotation7			
	8.2	Arbitration			
	8.3	Expert reviewing			
	8.4	Decision making box7			
	8.5	Annotators training and assessment			
	8.6	Variable description			
9	Consiste	ncy judgement			
	9.1	Input data type classification			
	9.2	Output requirement classification			
	9.3	Criteria option matrix			
	9.4	Post-processing of the annotations14			
10	10 Recommended metadata				
11	1 Output file				
12	12 File saving14				
Annex A Questionnaire on data annotation16					
Annex B Examples of endoscopic image metadata					
Bibliography21					

List of Tables

Page

Table 1 – Summary of common medical measurement modalities	9
Table 2 – Summary of input data modalities for AI4H tasks	9
Table 3 – Output requirements	10
Table 4 – Criteria options in different scenarios	10
Table 5 – Criteria calculation for classification	12
Table 6 – Criteria calculation for Image detection and segmentation	13
Table 7 – Recommended metadata	15

List of Figures

ITU-T FG-AI4H Deliverable DEL5.3

Data annotation specification

Summary

Data annotation would be one of the most dependable factors on model performance, it serves as one important aspect of data quality control on Artificial Intelligence for health. Deliverable 5.3 gives a general guideline of data annotation specification, including definition, background and goals, framework, standard operating procedure, scenario classifications and corresponding criteria, as well as recommended metadata, etc. A questionnaire is attached to seek input and collaboration with topic groups regarding data annotation.

1 Scope

Within the context of data quality for artificial intelligence applied in health, this deliverable gives a general guideline of data annotation specification, including inter alia definition, background and goals, framework, standard operating procedure, scenario classifications and corresponding criteria, as well as recommended metadata.

2 References

[ISO/IEC 2382]	ISO/IEC 2382:2015, Information technology — Vocabulary. https://www.iso.org/obp/ui/#iso:std:iso-iec:2382:ed-1:v1:en
[GHTF/SG1/N71]	IMDRF GHTF/SG1/N71:2012, Definition of the Terms 'Medical Device' and 'In Vitro Diagnostic (IVD) Medical Device', <u>https://www.imdrf.org/sites/default/files/docs/ghtf/final/sg1/technical-docs/ghtf-sg1-n071-</u> 2012-definition-of-terms-120516.pdf
[SAMD/N12]	IMDRF/SaMD WG/N12FINAL:2014, "Software as a Medical Device": Possible Framework for Risk Categorization and Corresponding Considerations", <u>https://www.imdrf.org/sites/default/files/docs/imdrf/final/technical/imdrf-tech-140918-samd-framework-risk-categorization-141013.pdf</u>

3 Definitions

3.1 Terms defined elsewhere

This document uses the following terms defined elsewhere:

3.1.1 Artificial intelligence [ISO/IEC 2382]: Branch of computer science devoted to developing data processing systems that perform functions normally associated with human intelligence, such as reasoning, learning, and self-improvement.

3.1.2 Machine learning [ISO/IEC 2382]: Automatic learning, process by which a functional unit improves its performance by acquiring new knowledge or skills, or by reorganizing existing knowledge or skills.

3.1.1 Medical device [GHTF/SG1/N71:2012]: Any instrument, apparatus, implement, machine, appliance, implant, reagent for in vitro use, software, material or other similar or related article, intended by the manufacturer to be used, alone or in combination, for human beings, for one or more of the specific medical purpose(s) of: a) diagnosis, prevention, monitoring, treatment or alleviation of disease; b) diagnosis, monitoring, treatment, alleviation of or compensation for an injury; c) investigation, replacement, modification, or support of the anatomy or of a physiological process; d) supporting or sustaining life; e) control of conception, f) disinfection of medical devices; g) providing information by means of in vitro examination of specimens derived from the human body; and does

not achieve its primary intended action by pharmacological, immunological or metabolic means, in or on the human body, but which may be assisted in its intended function by such means.

3.1.4 Software as a medical device [SaMD/N12]: Software intended to be used for one or more medical purposes that perform these purposes without being part of a hardware medical device.

3.2 Terms defined in this document

This document defines the following terms:

3.2.1 Controlled vocabulary: An organized arrangement of words and phrases used to index content and to retrieve content through browsing or searching.

3.2.2 Data annotation: Perform operations such as categorizing, sorting, editing, marking, and annotating on the data to be labeled such as images, and add tags to the data to generate machine-readable data codes that meet the requirements of machine learning training.

3.2.3 Metadata: Data that provides information about other data.

3.2.4 Supervised learning: The machine learning task of learning a function that maps an input to an output based on example input-output pairs.

3.2.5 Unsupervised learning: A type of machine learning that looks for previously undetected patterns in a data set with no pre-existing labels and with a minimum of human supervision.

4 Abbreviations and acronyms

This document uses the following abbreviations and acronyms:

AI	Artificial Intelligence
AI4H	Artificial Intelligence for health
FG-AI4H	Focus Group on Artificial Intelligence for health
JSON	JavaScript Object Notation
ML	Machine Learning
SOP	Standard operating procedure
XML	Extensible Markup Language

5 Conventions

None.

6 Background and goals

The great potential of digital technologies, especially Machine Learning (ML) and Artificial Intelligence (AI) are recognized to revolutionize the fields of medicine and public health in an unprecedented manner. While holding great promise, this rapidly developing field raises a number of uncertainties, for example if the model is poorly designed or the underlying training data are biased or incorrect, errors or problematic results can occur. AI technology can only be used with complete confidence if it has been quality controlled through a rigorous evaluation in a standardized way. Among all the quality controls, the data annotation would be one of the most dependable factors on model performance. In the case of mislabelled or inaccurate training instances, it is difficult for the supervised model to obtain the expected results. Many annotation tools exist, as exemplified in [8] to [27] but lack a consistent approach. The US Food and Drug Administration addressed some of these issues in [28].

Quality control on data annotation is a factor that is easily overlooked but crucial to the model performance. It is especially critical to models based on large-scale dataset. Therefore, this deliverable addresses the following points:

- To assist the quality control of data annotation from standard operating procedure.
- To reduce model performance problems caused by inconsistent data annotations.
- To enable large-scale dataset projects on high diversity of data formats and multi-annotators.
- To facilitate the training and education for non-professional annotators and improve common understandings.

7 Framework

Data annotation is one of the most dependable factors on the performance of supervised machine models. If the annotation for machine learning is incorrect, the decision rules built by the machine will be biased. As a part of the entire AI4H project, data annotation works as shown in Figure 1. During the testing and evaluation of the supervised machine models, unqualified annotation may be identified, which should be relabelled or deleted from the dataset.

With the help of annotators and annotation tools, a standard operating procedure of data annotation can convert input dataset into qualified annotations for supervised machine learning. This standard operating procedure is discussed in clause 8 in details.

The information from the data annotation process and the raw dataset can be used for training dataset for supervised machine learning and optimization, as well as testing dataset for the evaluation process. Therefore, data annotation has a very close relationship to the above core process of AI4H model, as a result, recognized as one of the most dependable factors on the model performance.



Figure 1 – Framework of data annotation and its external relations

8 Standard operating procedure

To establish a unified understanding and quality control mechanism, a standard operating procedure is recommended. Figure 2 illustrates a formulated process of data annotation, with much feasibility through variables and configurable threshold.



Figure 2 – Data annotation procedure

8.1 Independent annotation

The data annotation process starts with independent annotation, represented by the left grey box in Figure 2, each instance in the dataset needs to be labeled by all or part of annotators independently. To avoid bias in data distribution, it is suggested that the process is carried out by grouping and crossing, and ensure the effective resolution of inconsistencies. Several annotators (represented as variable n in the figure) are invited to label the raw dataset. Certain qualifications are required on the annotators, for example doctors and trained annotators in specific case domain.

However, for cost considerations, some projects will also set up one annotator (Set n to 1) in this parallel independent annotation part, and then goes to arbitration if encounter difficulties.

8.2 Arbitration

In the above independent annotation part, if there is an inconsistency that cannot be acceptable, or difficulties and uncertainties in single annotator setting, additional annotator with more experience should be introduced for the arbitration, represented by the upper right grey box in Figure 2. Stricter requirements on the arbitration annotator qualifications, for example, doctor with more than 3 years of experience in the case domain.

8.3 Expert reviewing

The expert reviewing is represented by the lower right grey box in Figure 2. This final review deals with some very tricky cases which cannot reach an agreement on the previous steps. Annotations confirmed by review experts will be marked as a final answer, and cases not approved could be considered to send back to arbitration process and arbitrated by another arbitration expert. Stricter requirements on the expert qualifications, for example, with 5 years of experience or more.

8.4 Decision making box

Represented by the blue boxes in Figure 2, the judgment and decision making on labelling consistency cannot be avoided anywhere in the independent annotation, arbitration and expert reviewing. The simple mechanism should be: If the consistency satisfies the specific requirements, like reaches a certain threshold or a combination of conditions, the annotation shall be saved with confidence; if the consistency dose not satisfy the specific requirements, like dose not reach a certain threshold or a combination of conditions, the annotation will be discarded. Therefore, the

criteria of consistency and corresponding requirements should be identified, and they are usually designed according to different scenarios, with more details in clause 9.

8.5 Annotators training and assessment

With the continuous popularity of the AI4H model afterwards, we may expect a future with more mature and extensive mechanisms for annotators' engagement. In addition to the experienced doctors mentioned above, candidates with no professional qualification but well-trained and quantitative assessed are also possible to be invited in the process of data annotation.

The training and assessment of annotators may include the following ways:

- Gold standard materials: Data annotation made by review experts or arbitration groups can be seen as gold standards, a unified document with examples can be developed as reference to teach candidates to how to achieve the tasks.
- Training courses: In addition to paper documents, training courses is also an effective way to
 educate candidates and reach a common understanding on data annotation tasks, especially in
 large-scale dataset and numerous annotators.
- Quantitative assessment: To evaluate the performance of different annotators, examinations and certifications with specific evaluation metrics can be conducted. Only after the corresponding evaluation metrics calculated with gold standard and annotator's results reaches a certain requirement, like beyond a certain threshold, candidate can be assigned to the annotation task being certificated.

8.6 Variable description

Variables and configurable threshold in this procedure are listed here for your convenience.

- Number of independent annotators
- Number of arbitration experts
- Number of review experts
- Different options on consistency criteria (usually the same in 8.1-8.3)
- Configurable requirement or threshold on consistency criteria in the independent annotation
- Configurable requirement or threshold on consistency criteria in the arbitration
- Configurable requirement or threshold on consistency criteria in the expert reviewing

9 Consistency judgement

For decision box in Figure 2, different criteria on consistency are selected according to different application scenarios. Main considerations are from two perspectives: one is input data type, elaborated in clause 9.1; the other is the output requirement for AI4H models, elaborated in clause 9.2. Under these two different classification dimensions, the options on consistency criteria will be different, elaborated in clause 9.3.

9.1 Input data type classification

Biomedical information evolved with the medicine practice and engineering technologies at an unprecedented speed through the medical images obtained by human body imaging, high-resolution viewing of cells, and pathological specimens. Modalities covered in common measurement include X-ray, ultrasound, magnetic resonance (MR), X-ray computed tomography (CT), nuclear medicine, and high-resolution microscopy, etc. Table 1 refers to their specific information.

	Dimensionality	Description	Anatomies
X-ray	2D, 2D+t	Produces images by measuring the attenuation of X-ray through the body, via a detector array [1]	Most organs
СТ	2D, 3D, 3D+t	Creates 2D cross-sectional images of the body by using a rotating X-ray source and detector [2]	Most organs
Ultrasound	2D, 2D+t, 3D, 3D+t	A transducer array emits acoustic pulses and measure he echoes from tissue scatters [1]	Most organs
MRI	3D, 3D+t	Use a magnetic field to align protons; RF and gradient pulses are used to selectively excite protons in tissues and blood in order to measure their spatially encoded unclear magnetic resonance signals [3]	Most organs
Nuclear	2D, 3D, 3D+t	Measures the emission of gamma rays through decay of radioisotopes introduced into the body via external detectors/Gamma cameras. [1]	All organs with radioactive tracer uptake
Microscop y	2D, 3D, 3D+t	Typically uses an illumination source and lenses to magnify specimens before capturing an image [1]	Primarily biopsies and surgical specimens

 Table 1 – Summary of common medical measurement modalities

Based on the above common medical measurement modalities, a classification of input data modalities for AI4H tasks are given in Table 2, with text and numbers added in specific cases of case history descriptions and blood pressure or respiratory rate, etc.

Data	Dimensionalit y	Description	Examples
Image	2D	Two-dimensional medical imaging	– Fundus photos
3D images	3D	Three-dimensional spatial imaging	 Sets of CT slices
4D	4D (3D+t)	3D space imaging changes over time	– Heart film imaging
Video	2D +t	Camera or monitor recording	- Falls among the elderly
Audio/ signal	1D +t	Sound or transmitted in signal form.	 Heart sound /ECG
Text	1D, 2D	Structured/ unstructured description in words	 Case history, diagnosis extraction
Single number	1D	Single measurement data	 Blood pressure or respiratory rate

Table 2 – Summary of input data modalities for AI4H tasks

9.2 Output requirement classification

When the final output requirements of models are different, even if it is the same input data format, data annotations will be different. Different output requirements include classification, detection, segmentation, localization, etc. Corresponding description and examples are given in Table 3.

Task	Description	Examples
Classification	The problem of classifying instances into two or more classes.	Identify abnormal tissueDiabetic retinopathy grade
Detection	Identify an object, usually marked with rectangle for further processing.	 Detect the position of a coronary plaque for further processing
Segmentation	separate certain lesions, and draw the specific outline of the lesion	– Tumour segmentation
Localization	Calculate the central coordinate of the anatomical structure	 Localize the optic disc or macular fovea for further analysis of ocular fundus diseases

Table 3 – Output requirements

9.3 Criteria option matrix

With the above two dimensions, a matrix can be developed according to different data input format and model output requirements. This matrix can act as a reference for the selecting criteria options. Details are shown in Table 4, and other scenarios are to be added to cover all possible use cases in the FG and the AI4H industry.

Task	Classification	Detection	Segmentation	Localization
Data type				
Image		Type 2: Detection and segmentation for images		Type 3: Localization
3D images	Type 1: Classification	 (a) slicing 3D data into different 2D views before fusing to obtain a final detection or segmentation regions (b) exploit the 3D data by using architectures that perform 3D convolutions and then train the network from scratch on 3D medical images [4][5][6][7] 		 (a) slicing 3D data into 2D views to obtain the regions of the target object before calculation of the final position coordinate (b) exploit the 3D data by using architectures that perform 3D convolutions and then train the network from scratch on 3D medical images
4D		condensing the 4D data into three dimensions		—
Video		condensing the 2D +t data into three dimensions		—
Audio/ signal				_
Text		_	_	_
Single number		_	_	_

 Table 4 – Criteria options in different scenarios

Type 1: Classification

For this type, criteria like Cohen's kappa, weighted kappa, Fleiss' kappa and Krippendorff's alpha are recommended to use for classification tasks. The detailed calculation methods are shown in Table 5.

- Cohen's kappa: Cohen's kappa coefficient (κ) is a statistic that is used to measure inter-rater reliability for qualitative items. It is generally thought to be a more robust measure than simple percent agreement calculation, as κ considers the possibility of the agreement occurring by chance.
- Weighted kappa: Weighted kappa allows disagreements to be weighted differently, and is especially useful when codes are ordered. Three matrices are involved, the matrix of observed scores, the matrix of expected scores based on chance agreement, and the weight matrix.
- Fleiss' kappa: Fleiss' kappa is a statistical measure for assessing the reliability of agreement between a fixed numbers of raters when assigning categorical ratings to a number of items or classifying items. This contrasts with other kappas such as Cohen's kappa, which only work when assessing the agreement between not more than two raters or the interrater reliability for one appraiser versus themselves.
- Krippendorf's alpha: Krippendorf's alpha is an assessment of inter-rate reliability dealing with missing data, various sample sizes, categories and numbers of raters, and any type of measurement level. It can be seen as a generalization of Fleiss' kappa (and others).

Type 2: Detection and segmentation for images

For this type, criteria like Jaccard index and Dice's coefficient are recommended to use for detection and segmentation for images. Detailed calculation methods are shown in Table 6.

- The Jaccard index: Jaccard index is also known as Intersection over Union (IoU) and the Jaccard similarity coefficient, which is a statistic used for gauging the similarity and diversity of sample sets.
- Dice's coefficient: Dice's coefficient is the quotient of similarity and ranges between 0 and 1. This coefficient is not very different in form from the Jaccard index, and they have a connection as J=D/(2-D), D=2J/(1+J)

Type 3: Localization

For this type, criteria like Euclidean Distance (ED) is recommended to use for localization. Euclidean distance is a commonly used definition of distance, which refers to the true distance between two points in the m-dimensional space.

- In 2D space: $ED((x_1, y_1), (x_2, y_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$.

- In 3D space:
$$ED((x_1, y_1, z_1), (x_2, y_2, z_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2}$$

Criteria	Situation	Calculation method	Parameter explanation
Cohen's kappa	Assessing the agreement between not more than two raters or the interrater reliability for one appraiser versus themselves.	$\kappa = \frac{p_o - p_e}{1 - p_e} = 1 - \frac{1 - p_o}{1 - p_e}$ If the raters are in complete agreement then kappa =1; If there is no agreement among the raters other than what would be expected by chance kappa =0. It is possible for the statistic to be negative which implies that there is no effective agreement between the two raters or the agreement is worse than random.	where p_o is the relative observed agreement among raters (identical to accuracy), and p_e is the hypothetical probability of chance agreement, using the observed data to calculate the probabilities of each observer randomly seeing each category $p_e = \frac{1}{N^2} \sum_k n_{k1} n_{k2}$
Weighted kappa	Allows disagreements to be weighted differently, and is especially useful when codes are ordered.	$\kappa = 1 - \frac{\sum_{i=1}^{k} \sum_{j=1}^{k} w_{ij} x_{ij}}{\sum_{i=1}^{k} \sum_{j=1}^{k} w_{ij} m_{ij}}$	where k is the number of codes and w_{ij} , x_{ij} , and m_{ij} are elements in the weight, observed, and expected matrices, respectively. The weights in the diagonal cells are all 1 (i.e., $w_{ij} = 1$, for all i), and the weights in the off-diagonal cells range from 0 to <1 (i.e., $0 \le w_{ij} < 1$, for all $i \ne j$).
Fleiss' kappa	Assessing the reliability of agreement between a fixed number of raters when assigning categorical ratings to a number of items or classifying items.	$\kappa = \frac{\overline{P} - \overline{P}_e}{1 - \overline{P}_e}$ If the raters are in complete agreement, then Fleiss' kappa =1. If there is no agreement among the raters (other than what would be expected by chance) then Fleiss' kappa <0.	The factor $1 - \overline{P}_e$ gives the degree of agreement that is attainable above chance, and $\overline{P} - \overline{P}_e$ gives the degree of agreement actually achieved above chance.
Krippendor f's alpha	Assessment of inter-rate reliability dealing with missing data, various sample sizes, categories and numbers of raters, and any type of measurement level. Generalization of Fleiss' kappa (and others)	$\alpha = 1 - \frac{D_o}{D_e} D_0 = \frac{1}{n} \sum_{c \in R} \sum_{k \in R} \delta(c, k) \sum_{u \in U} m_u \frac{n_{cku}}{P(m_u, 2)} D_e$ $= \frac{1}{P(n, 2)} \sum_{\substack{c \in R}} \sum_{\substack{k \in R}} \delta(c, k) P_{ck} P_{ck}$ $= \frac{n_c n_k i f c \neq k}{n_c (n_c - 1), i f c = k}$	D_o : Disagreement observed D_e : Disagreement expected by chance δ : Metric function n : Number of pairable elements m_u : Number of items per unit/sample $n_{c,k,u}$:Number of pairs in unit u P : Permutation function $P_{c,k}$: Number of permutations in pair (c, k)

Table 5 – Criteria	calculation	for	classification
--------------------	-------------	-----	----------------

Criteria	Calculation method	Graphical representation
Jaccard index	Numerator represents the area of overlap between two annotations; Denominator represents the area encompassed by two annotations. Dividing the area of overlap by the area of union yields our final score. $J(A,B) = \frac{ A \cap B }{ A \cup B } = \frac{ A \cap B }{ A + B - A \cap B }$	A AOB B
Dice's coefficient	Numerator represents the double area of overlap between two annotations; Denominator represents the sum of two annotation area. Dividing the area of overlap by the sum area yields our final score. $D = \frac{2 A \cap B }{ A + B }$	J=D/(2-D), D=2J/(1+J)

Table 6 – Criteria calculation for Image detection and segmentation

9.4 Post-processing of the annotations

After the criteria calculation and consistency judgment, different post-processing methods on annotations that are acceptable as consistent will also cause different result. For example, calculate the average value of the marked results (x, y, w, h) or a maximum area with a consistency above threshold in an image.

10 Recommended metadata

Metadata is considered to be the output of the data annotation process, all necessary information for the annotation process should be included in the metadata. A metadata format is to be given in Table 7, further details are for further study.

11 Output file

The output files include the annotation documents and the origin images. The formats of the annotation documents include but are not limited to XML, JSON, text, etc. The annotation documents should include at least three items: image path, object name and object coordinates. Supporting document may be given if it's necessary to interpret the annotation information. Annex B gives an example of the annotation document for endoscopic images.

12 File saving

Both the images and the documents should be named according to the same rules for easy querying. For example, they can be named with the number of the classification of the object, and the document's name is the same with the corresponding origin images.

Туре	Content						
General information	a. Institution and responsible or corresponding PI;						
	b. Construction dates of annotation dataset;						
	c. Regulatory aspects (e.g. Data privacy)						
	d. Description of use case						
Annotation procedure	e. Details on data annotation process (annotator number, experts group setting, tools, etc.)						
information	f. Achieved consistency and criteria employed						
	g. Post- processing method employed on annotations						
	h. Ontology employed						
	i. Label list or description						
Data acquisition information	j. collection device model						
	k. collection frame rate/ Sampling rate						
l. Instance Information	m. Instance identification code						
	n. patient information (age, gender)						
	o. diagnosis information (symptoms)						
Annotation information	Task	Classification	Detection	Segmentation	Localization		
	Data type						
	– Image	– signal instance		– label per instance	 coordinate label of signal instance 		
	- 3D images		 signal instance 				
	– 4D						
	– Video	class labels					
	 Audio/ signal Text 						
			_	Label per word,			
				intent, or sentence	_		
	– Single number	_	_	_	_		

Table 7 – Recommended metadata

Annex A Questionnaire on data annotation

By Google Form: https://forms.gle/3fYrm3SZSrNQu3eeA

The aim of this questionnaire is to gather insights into the current practices, the specific requirements of data annotation in the FG-AI4H topic groups and AI4H products. Your input and suggestion will be of great value for us in forming a data annotation specification together, as one of the deliverables with the FG-AI4H. We would appreciate it if you could take the time to complete the questionnaire, or if you have further ideas, please feel free to contact us. (xushan@caict.ac.cn; sebastian.bosse@hhi.fraunhofer.de, edwinjrwu@tencent.com) 1. To which topic group are you contributing? 10.1 TG-Cardio 10.2 TG-Derma 12.5% 12.5% 10.3 TG-Bacteria 10.4 TG-Falls 12.5% 12.5% 10.5 TG-Histo 10.6 TG-Malaria 12.5% 10.7 TG-MCH 25% 10.8 TG-Neuro 12.5% 🔺 1/3 🔻 2. Which annotation task category is relevant for your project within the topic group? Classification 7 (87.5%) 4 (50%) Detection Segmentation 4 (50%) 4 (50%) Prediction/ regression 0 2 4 6 8

16

3. Is there any gold standard (or state-of-the-art task intervention method) relevant for your project within the topic group?

Histology, Cross-image validation, human annotations

Pathological report, Cross annotation by doctors

Snake expert (herpetologist) identification

The gold standard photography method for the detection of Diabetic Retinopathy is stereoscopic color fundus photography in 7 standard fields (30°) as defined by the Early Treatment Diabetic Retinopathy Study (ETDRS) group.

Post-mortem pathology evaluation is the gold standard. unfortunately few data are available with postmortem evaluation. In the absence of such data, biological markers provide a more reliable alternative to the more uncertain clinical diagnostic based on symptoms only.

Yes

The average doctor opinion is the current gold-standard (even cases confirmed later with more evidence cannot be used to judge the "correct" answer for less evidence).

4. What is your data source for the training and testing dataset? (appreciated if you can give more info)



6. What is the nature of the annotation?



7. What kind of annotation procedure are you using? any annotation tool that you use?

all of the above, custom made tool

Cross annotation, Self-built annotation tool

Expert identification, crowdsourcing

For Diabetic Retinopathy and other conditions, annotation is usually provided by an ophthalmologist's diagnosis assigned to the image.

Manual

Structured data are used, thus simple R programming is used to recode structured data to required standardized labels.

Most companies have their own annotation tool. For the group we are developing a new case-creation tool where doctors enter the symptoms and expected conditions in a semantically structured way.

8. What annotation quality criterions are currently used in your topic group?



9. What kind of metadata do you consider relevant for data annotation?

Correspo Regulatory aspects (e.g. Data privacy) Achieved consistency criteria and thres Sample identification code hospital, patient information (age, gen Labels itself(types, annotation area)	-1 (12.5%) -1 (12.5%)	—2 (25%) —2 (25%)	—3 (37.5%)	-4 (50%)	5 (62.5%) 5 (62.5%) 5 (62.5%) 5 (62.5%) 5 (62.5%) 5 (62.5%)
hap	—1 (12.5%) 1	2	2	4	5
10. What type of ontology are yo	u using? (If any)	Z	5	4	5
	• ••••••••••••••••••••••••••••••••••••				
Medical ICD10 classification of dise	eases.				
N/A					
not finally decided yet, but most likely SnomedCT - maybe also HPO, ICD10 and we will in any case need to filter and also extend them. We will also need new ontologies e.g. for triage levels.					
11 Which additional information do you need to encode the actual meaning of the annotation?					
			tuur meaning	, of the unifor	ation.
ideally also dental history and tooth specific clinical findings					
-					
The following paper provides addit Observational Research Designs (D	ional indicators: "Gui AQCORD)" (https://t.	delines for Da .co/brhboFYI5	ta Acquisition, (4?amp=1)	Quality and Cur	ation for
N/A					

Annex B Examples of endoscopic image metadata

There is an example of metadata of Endoscopic image. The annotation results include:

Attributes	Examples			
json version	-			
folder number	Trial_Date_Hospital_Batch			
filename	PatientNum_Uniquecode.png			
file path	/Data/Endoscopic/Trial_Date_Hospital_Batch/ PatientNum_Uniquecode.png			
source	Hospital Equipment			
size: width/ height/depty	1280×720			
segmented object: name, <u>coordinates</u>	<pre>"label_classification": [</pre>			

Bibliography

NOTE - [8] onwards are references on some data annotation tools.

- [1] J. T. Bushberg et al., The Essential Physics of Medical Imaging. Philadelphia, PA, USA: Lippincott Williams Wilkins, 2011.
- [2] J. Hsieh, Computed Tomography: Principles, Design, Artifacts, and Recent Advances.Bellingham,WA, USA: Soc. Photo-Opt. Instrum. Eng.,2009.
- [3] Panayides A S, Amini A, Filipovic N D, et al. AI in Medical Imaging Informatics: Current Challenges and Future Directions[J] . IEEE Journal of Biomedical and Health Informatics, 2020, 24(7): 1837-1857.
- [4] Q. Dou et al., "Automatic detection of cerebral microbleeds from MR images via 3D convolutional neural networks," IEEE Trans. Med. Imag., vol. 35, no. 5, pp. 1182–1195, May 2016.
- [5] J. Ding, A. Li, Z. Hu, and L. Wang, "Accurate pulmonary nodule detection in computed tomography images using deep convolutional neural networks," in Proc. Int. Conf. Med. Image Comput. Comput.-Assisted Intervention, Sep. 2017, pp. 559–567.
- [6] S. Korolev, A. Safiullin, M. Belyaev, and Y. Dodonova, "Residual and plain convolutional neural networks for 3D brain MRI classification," Jan. 2017, online. Available: <u>https://arxiv.org/abs/1701.06643</u>
- [7] D. Nie, H. Zhang, E. Adeli, L. Liu, and D. Shen, "3D deep learning for multi-modal imagingguided survival time prediction of brain tumor patients," in Proc. MICCAI, 2016, pp. 212– 220.
- [8] Sangkuhl K, Whirl-Carrillo M, Whaley R M, et al. Pharmacogenomics Clinical Annotation Tool (Pharm CAT) [J]. Clinical Pharmacology & Therapeutics, 2020, 107(1): 203-210.
- [9] Hou R, Denisenko E, Forrest A R R. sc Match: a single-cell gene expression profile annotation tool using reference datasets [J]. Bioinformatics, 2019, 35(22): 4688-4695.
- [10] Philbrick K A, Weston A D, Akkus Z, et al. RIL-Contour: a Medical Imaging Dataset Annotation Tool for and with Deep Learning [J]. Journal of digital imaging, 2019, 32(4): 571-581.
- [11] Kraljevic Z, Bean D, Mascio A, et al. MedCAT Medical Concept Annotation Tool [J]. arXiv preprint arXiv:1912.10166, 2019.
- [12] Ramos A H, Lichtenstein L, Gupta M, et al. Oncotator: cancer variant annotation tool [J]. Human mutation, 2015, 36(4): E2423-E2429.
- [13] Iakovidis D K, Goudas T, Smailis C, et al. Ratsnake: a versatile image annotation tool with application to computer-aided diagnosis [J]. The Scientific World Journal, 2014, 2014.
- [14] Seifert S, Kelm M, Moeller M, et al. Semantic annotation of medical images [C]. Medical Imaging 2010: Advanced PACS-based Imaging Informatics and Therapeutic Applications. International Society for Optics and Photonics, 2010, 7628: 762808.
- [15] Russell B C, Torralba A, Murphy K P, et al. LabelMe: a database and web-based tool for image annotation [J]. International journal of computer vision, 2008, 77(1-3): 157-173.
- [16] Diete A, Sztyler T, Stuckenschmidt H. A smart data annotation tool for multi-sensor activity recognition [C]. 2017 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops). IEEE, 2017, 111-116.

- [17] Oronoz M, Casillas A, Gojenola K, et al. Automatic annotation of medical records in Spanish with disease, drug and substance names [C]. Iberoamerican Congress on Pattern Recognition. Springer, Berlin, Heidelberg, 2013: 536-543.
- [18] Ferreira P M, Mendonça T, Rozeira J, et al. An annotation tool for dermoscopic image segmentation [C]. Proceedings of the 1st International Workshop on Visual Interfaces for Ground Truth Collection in Computer Vision Applications. 2012: 1-6.
- [19] Kim E, Huang X, Tan G. Markup SVG An Online Content-Aware Image Abstraction and Annotation Tool [J]. IEEE Transactions on Multimedia, 2011, 13(5): 993-1006.
- [20] Rubin D L, Mongkolwat P, Kleper V, et al. Annotation and image markup: accessing and interoperating with the semantic content in medical imaging [J]. IEEE Intelligent Systems, 2009, 24(1): 57-65.
- [21] Rubin D L, Mongkolwat P, Kleper V, et al. Medical Imaging on the Semantic Web: Annotation and Image Markup [C]. AAAI Spring Symposium: semantic scientific knowledge integration. 2008: 93-98.
- [22] Karlgren K, Dahlström A, Ponzer S. Design of an annotation tool to support simulation training of medical teams [C]. European Conference on Technology Enhanced Learning. Springer, Berlin, Heidelberg, 2008: 179-184.
- [23] Karlgren K, Dahlström A, Lonka K, et al. A new educational annotation tool for supporting medical teams to improve their teamwork and communication [J]. ICEM/ILE, 2007: 20-22.
- [24] Yao B, Yang X, Zhu S C. Introduction to a large-scale general purpose ground truth database: methodology, annotation tool and benchmarks [C]. International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition. Springer, Berlin, Heidelberg, 2007: 169-183.
- [25] Kumar A, Jawahar C V. Content-level annotation of large collection of printed document images[C]. Ninth International Conference on Document Analysis and Recognition (ICDAR 2007). IEEE, 2007, 2: 799-803.
- [26] Hollink L, Schreiber G, Wielemaker J, et al. Semantic annotation of image collections [C]. Knowledge capture. 2003, 2.
- [27] Kipp M. Anvil-a generic annotation tool for multimodal dialogue [C]. Seventh European Conference on Speech Communication and Technology. 2001.
- [28] US Food and Drug Administration Proposed Regulatory Framework for Modifications to Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD) - Discussion Paper and Request for Feedback. <u>https://www.fda.gov/files/medical%20devices/ published/US-FDA-Artificial-Intelligence-and-Machine-Learning-Discussion-Paper.pdf</u>

22