

1 Multimedia Appendix 1: Gold Standard Development

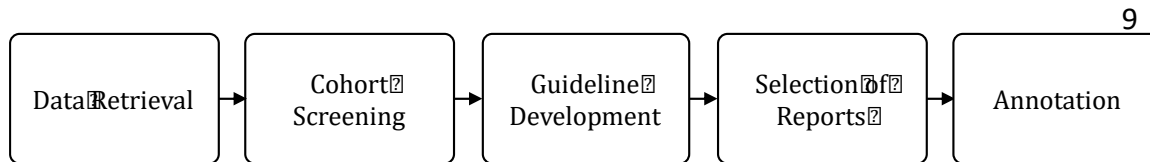
2

3 Gold Standard Development

4 The process of generating the gold standard reports followed a standardized data abstraction
5 process with four major steps: data retrieval, cohort screening, guideline development,
6 selection of reports and annotation (Figure 1).

7

8 Figure 1. Overview of gold standard creation process



10

11 Data Retrieval

12 At TMC, the data comprised ICD-9 codes, CPT-4 codes, clinical notes and neuroimaging reports,
13 and were obtained from three EHRs: General Electric Logician (outpatient general medicine),
14 eClinicalWorks (outpatient specialties including neurology and neurosurgery), and Cerner
15 Soarian (all inpatient encounters). At Mayo Clinic, the same types of data were retrieved from
16 the Mayo Unified Data Platform (UDP).

17 Cohort Screening

18 A screening protocol, comprised of diagnosis codes and problem lists, was used to identify
19 index neuroimaging reports (the patient's first MRI or CT in the EHR) for individual patients
20 without clinically-evident stroke, transient ischemic attack (TIA), and dementia any time before
21 or up to 60 days after the index imaging exam. The protocol was prototyped by a Vascular
22 Neurologist (LYL) leveraging literature review(1, 2) and refined as a multi-disciplinary team
23 effort including an internist (DMK), a neuroradiologist (PHL) and a radiologist (DFK). The final
24 version of the protocol was vetted by both clinical and informatics teams from two sites. The
25 protocol included both terminology codes and the corresponding text descriptions. For
26 example, the text description of the code 434.10 is "cerebral embolism without cerebral
27 infarction", which was used for text search. We excluded TIA because TIA is sometimes
28 incorrectly assigned by clinicians as the diagnosis in the setting of transient neurologic
29 symptoms and positive evidence of brain infarction on neuroimaging(3). Dementia was an
30 exclusion criterion because of a projected future application of the NLP algorithm in identifying
31 patients for comparative effectiveness studies or clinical trials for which both stroke and
32 dementia could be outcomes of interest. To conduct cohort screening, an exclusion was applied
33 when a patient has exclusionary ICD codes or his/her clinical note contained more than one
34 type of keywords from the key terms list. An NLP system MedTagger(4) was utilized at Mayo
35 Clinic to capture the mention from clinical notes. At TMC, this process was conducted through
36 manual review. The screening result was validated manually by comparing the screened cohort
37 with EMRs (clinicians: Mayo: PL, DI; TMC: LYL).

38 *Guideline Development*

39 A baseline annotation guideline and online educational videos were created to facilitate
40 annotation guideline development. To develop the guideline, a preliminary annotation was
41 performed individually by two senior radiologists, one senior neurologist and four medical
42 residents. All annotators were asked to annotate the same 40 reports generated by both Mayo
43 and Tufts. Feedback collected from the annotation was used to revise and update the guideline.
44 The guideline includes task definition, annotation instruction, annotation concept and
45 examples. The annotation task is to tag the findings of SBI and WMD lesion in both body
46 (Findings) and summary (Impression and Assessment) sections of radiology reports. The
47 definitions for SBI and WMD on both CT and MRI are mentioned in the guideline. For example,
48 positive SBI findings on CT are defined as discrete, focal hypodense lesion greater than 3 mm in
49 size conforming to a vascular distribution in the white matter, gray matter, or both(1, 2). For
50 WMD, findings on MRI are defined as Confluent or poorly marginated T2 hyperintensities
51 involving the white matter that do not meet criteria for brain infarction(1, 2). We defined likely
52 but uncertain findings as SBI_Indeterminate. The annotation concept for SBI and WMD contains
53 two categories: mentions of SBI and WMD related findings and attributes of the findings (e.g.
54 Acuity for SBI and Grade for WMD).

55 *Selection of Reports*

56 At Mayo, 262,061 reports were obtained from Mayo EHR based on the CPT inclusion criteria.
57 4000 reports were randomly sampled for cohort screening. 910 were eligible for annotation
58 after applying the ICD exclusion criteria (structured and unstructured). At TMC, 63,419 reports
59 were obtained from TMC EHR based on CPT inclusion criteria. 12,092 reports remained after
60 applying the ICD exclusion criteria (structured). 1000 reports were randomly selected for text
61 screening. 773 reports were eligible for annotation.

62
63 Each site randomly selected 500 radiology reports (Total 1000) from the eligible samples
64 (Mayo: 910, TMC: 773). 400 reports were randomly sample from the total 1000 reports (500
65 Mayo, 500 TMC) and duplicated for double reading in order to calculate inter-annotator
66 agreement (IAA). Each resident was assigned 350 reports stratified by report type (CT, MR) and
67 site (Mayo, TMC).

68 *Annotation*

69 Two third-year residents (KAK, MSC) from Mayo and two first-year residents (AOR, KN) from
70 TMC performed the annotation. The annotation was organized into two phases. The first phase
71 extended from the finalization of the preliminary guideline until the midpoint when one-half of
72 the reports were annotated. The primary goal for the first phase was to identify new problems
73 that were not captured in the sample data. This allowed the guideline to be more robust in
74 accommodating new data. After the first phase, all problematic cases were reviewed by the two
75 senior clinicians, and the guidelines were updated. The second phase annotation then
76 commenced using the updated guidelines. A consensus meeting was organized to resolve all
77 disagreements after the annotation. All conflicting cases were then adjudicated by the two
78 senior clinicians. The annotation process was supervised by two senior clinicians (LKL, PL). All

79 the issues during the process were documented. The final gold standard corpus consisted of
80 1,000 annotated neuroimaging reports (500 TMC, 500 Mayo).

81

82 [References](#)

83

84 1. Fanning JP, Wesley AJ, Wong AA, Fraser JF. Emerging spectra of silent brain infarction.
85 *Stroke*. 2014;45(11):3461-71.

86 2. Wardlaw JM, Smith EE, Biessels GJ, Cordonnier C, Fazekas F, Frayne R, et al.
87 Neuroimaging standards for research into small vessel disease and its contribution to ageing
88 and neurodegeneration. *The Lancet Neurology*. 2013;12(8):822-38.

89 3. Sorensen AG, Ay H. Transient ischemic attack: definition, diagnosis, and risk
90 stratification. *Neuroimaging Clinics*. 2011;21(2):303-13.

91 4. Liu H, Bielinski SJ, Sohn S, Murphy S, Waghlikar KB, Jonnalagadda SR, et al. An
92 information extraction framework for cohort identification using electronic health records.
93 *AMIA Summits on Translational Science Proceedings*. 2013;2013:149.

94